Are Police Racially Biased in the Decision to Shoot?*

Tom S. Clark[†] Elisha Cohen[‡] Adam Glynn[§] Michael Leo Owens[§] Anna Gunderson[¶] Kaylyn Jackson Schiff[‡]

September 13, 2021

Abstract

We present a theoretical model predicting racially biased policing produces 1) more use of potentially lethal force by firearms against Black civilians than against White civilians and 2) lower fatality rates for Black civilians than White civilians. We empirically evaluate this second prediction with original officer-involved shooting data from eight local police jurisdictions from 2010 to 2017, finding that Black fatality rates are significantly lower than White fatality rates and that this significance would survive an omitted covariate three times as strong as any of our observed covariates. Furthermore, using outcome test methodology and a comparability assumption, we estimate that at least 30% of Black civilians shot by the police would not have been shot had they been White. An omitted covariate would need to be at least three times as strong as any of our observed covariates to eliminate this finding. Finally, any omitted covariate would have to affect Black fatality rates and not Hispanic fatality rates in order to be consistent with the data.

^{*}We thank Chuck Cameron, Anna Harvey, Hakeem Jefferson, John Kastellec, Beatriz Magaloni-Kerpel, Pablo Montagnes, John Patty, Zac Peskowitz, Dan Thompson, Georg Vanberg, and Kirsten Widner for helpful comments. Earlier versions of this paper were presented at the 2019 meeting of the American Political Science Association, the Politics of Policing Conference at Princeton University, and seminars at Emory University.

[†]Charles Howard Candler Professor of Political Science, Emory University

[‡]Ph.D. Student, Emory University

[§]Associate Professor of Political Science, Emory University

[¶]Assistant Professor, Louisiana State University

1 Introduction

Police are agents of the state, exercising a high degree of autonomy and discretion when implementing policy (M. K. Brown, 1981; Wilson, 1978). But, unlike other domestic agents of the state, "the police are... a mechanism for the distribution of situationally justified force in society" (Bittner, 1970, p. 39). Consequently, the character of their interactions with the public differ greatly from those of other "street-level bureaucrats" (Lipsky, 1980): Policecivilian encounters are more unpredictable, with greater potential for violence and death, for civilians and police. Accordingly, policing is "profoundly involved with the most significant questions facing any political order, those pertaining to justice, order, and equity" (M. K. Brown, 1981, pp. 6-7). It is especially true when police use their discretion to shoot civilians.

While police use force against civilians more in some nations than others, police shootings of civilians are more common in the United States relative to other advanced, liberal democracies (Zimring, 2017). Furthermore, racial disparities in police use of force in the U.S. seem common and particularly wide between Blacks and Whites, and a marker of racial disparities in policing, generally, including deployment, surveillance, involuntary contact by stop-and-frisk, arrest, and jailing (Bittner, 1970; Soss and V. Weaver, 2017; R. A. Brown, 2019). Given the fraught history and contemporary realities of race in the U.S., racial disparity in police shootings raise concerns about racial bias influencing officers' discretion to shoot during police-civilian encounters. Whether racial bias causes racial disparities in policing, and how much, however, remains an academic and civic puzzle.

It is empirically difficult to discern how many police shootings of Black Americans result from their disproportionate contact with police versus disproportionate use of force by police against them versus racial bias by patrol officers and their departments (e.g., Knox, Lowe, and Mummolo, 2020; Knowles, Persico, and Todd, 2001; Fryer Jr., 2016). Further, neither police departments nor agencies overseeing them track or report all lethal and non-lethal police shootings of civilians, especially by race (Zimring, 2017). Consequently, depending on the data, measures, and methods, studies draw contradictory conclusions, ranging from significant differences in the likelihood and speed of shooting Black civilians compared to other civilians (Mekawi and Bresin, 2015) to no racial differences in fatal shootings of civilians by police (Johnson et al., 2019). Therefore, even when relatively good data are available for social scientists to observe and describe racial patterns in policing, scholarly consensus on whether and how much police discriminate by race of civilian when using lethal force, let alone nonlethal force, remains elusive.

To better assess whether there is evidence of racial bias in the use of force by police against civilians, measured by shootings, lethal and non-lethal, we develop a model of policecivilian encounters that yields empirical implications for evaluating racial bias in officerinvolved shootings (OIS). In our model, informed by studies of the transactional nature and iterative process of police-civilian encounters (Binder and Scharf, 1980; Terrill, 2005; Kahn, Steele, et al., 2017), civilians and police engage in behaviors, include actions that may and can escalate their encounters towards harm, including police violence against civilians (and civilian violence against police). Ultimately, our model predicts that racially biased police officers will be more likely to use force against Black civilians than against White civilians. Moreover, police shootings of Black civilians should result in more non-fatalities than fatalities.

We test the implication of our model with OIS data from eight local police jurisdictions in the U.S.. Our data, covering 2010 through 2017, and obtained through public records requests, include all instances of police reporting they shot civilians—fatally and non-fatally and the race of civilians, along with other attributes of the police-civilian encounters. Consistent with our theoretical expectation, we find that Black civilians are significantly more likely to survive an OIS, reflecting, we posit, a higher degree of racial bias in the decisions by officers to shoot Black civilians compared to non-Black civilians. Furthermore, we show that this disparity would survive an omitted covariate three times as strong as any of our observed covariates.¹

¹Strength is defined in terms of the bias factor of VanderWeele and Ding (2017).

Additionally, we estimate a lower bound on the magnitude of racial bias in the decision to shoot a civilian, guided by Knox, Lowe, and Mummolo (2020) and Cohen (2021). Borrowing their techniques, we conceptually divide Black civilians who were shot by police into two groups—1) Black civilians that would have been shot had they been White and 2) Black civilians that would not have been shot had they been White. The proportional size of the second group is our parameter of racial bias. To estimate a lower bound for this quantity, we evaluate the difference in fatality rates of White and Black civilians shot by the police in the eight local jurisdictions relative to their White fatality rates, where we posit fatal shootings are more likely to be justified as "reasonable" shootings from the perspective of police departments, and that non-fatal shootings are more prevalent among Black civilians compared to other groups. Using the techniques from VanderWeele and Ding (2017) and Cohen (2021), we estimate that at least 30% of Black civilians shot would not have been shot had they been White and that to eliminate this estimate, an omitted covariate would again need to be three times a strong as any of our observed covariates. Finally, such an omitted covariate would have to affect Black fatality rates and not Hispanic fatality rates in order to be consistent with the data.²

Our theory and findings demonstrate that identifying racial bias in police decision-making is possible, buttressing other research (Knowles, Persico, and Todd, 2001; Persico and Todd, 2006; Knox and Mummolo, 2020; Knox, Lowe, and Mummolo, 2020). That alone is important in light of the continuing need to understand discretion by the police as "street-level bureaucrats" and how much race affects policing, including use and severity of force. Plus, our theory and findings about the most extreme form of police use of force bear on classic concerns in political science, including but not limited to the exercise of power by the state, democratic accountability, and equality under the law (M. K. Brown, 1981).

 $^{^{2}}$ We lack data on all instances of police drawing their weapons, but including moments where police drew guns without firing would likely increase the estimate of the lower bound (Worrall et al., 2018).

2 Police Discretion in Use of Force

Encounters with the police are among the most common encounters civilians have with government agents (Jacob, 1972; M. K. Brown, 1981; Soss and V. Weaver, 2017). A key contrast with other civilian encounters with government agents is that police-civilian contact, whether initiated by police or initiated by civilians, has the potential for violence. How officers exercise their discretion to use force and violence during police-civilian encounters and why it may cause racial disparities are important considerations (e.g., Terrill, 2011). "In the police shooting context," in particular, "there is a concern that officers, despite their best intentions and/or conscious beliefs, will subconsciously let preconceived ideas about certain individuals influence their decision processes" (Worrall et al., 2018, p. 1176). This includes their racial beliefs, which may bias their behaviors during police-civilian encounters. Inferring racial bias, however, is challenging.

2.1 Racial Disparities in Use of Force

Generally, social scientists expect police are more likely to use force and more of it against Black civilians than against White civilians (L. James, Vila, and Daratha, 2013; Philip Atiba Goff et al., 2016; Jetelina et al., 2017). Whether police do is well-studied experimentally and observationally, often finding that officers *are* more willing to use force against Black civilians than against White civilians (Correll, Park, Judd, Wittenbrink, et al., 2007; Mekawi and Bresin, 2015; Eberhardt et al., 2004; Buehler, 2017; Sikora and Mulvihill, 2002; Johnson et al., 2019; Robert E. Worden, 2015; Engel and Calnon, 2004; Schuck, 2004; Terrill, 2005; Baumgartner, D. A. Epp, and Shoub, 2018). Furthermore, the recent availability of "big data" on police-civilian encounters at incident-level (e.g., New York City's Stop, Question, and Frisk program) has enabled rigorous social science to deepen evidence of racial disparities in police use of force (e.g., Fryer Jr., 2016; Voigt et al., 2017; Pierson, Simoiu, Overgoor, Corbett-Davies, Ramachandran, et al., 2017; Gelman, Fagan, and Kiss, 2007; Goel, Rao, and Shroff, 2016; Mummolo, 2018).

However, some studies temper or contradict claims and the expectation of racial bias in police use of force, particularly shootings (e.g., Worrall et al., 2018). In other words, racial bias in policing may not necessarily increase the likelihood of use of force against Black civilians. Some evidence, drawn typically from observational studies, and limited by concerns about unmeasured confounding and/or misapplied methods (J. H. Garner, Schade, et al., 1995; J. Garner and C. Maxwell, 1999; J. H. Garner, C. D. Maxwell, and Heraux, 2002; Alpert and Dunham, 2004; Fryer Jr., 2016; Johnson et al., 2019), suggests we should expect and observe either smaller-scale or no racial disparities in police use force (e.g., shootings). Plus, a "counter bias" may exist, inducing officers to be extra sensitive to the potential negative consequences of using force against racial minorities, especially Black civilians (L. James, Vila, and Daratha, 2013). The negative consequences of using force and more of it against Black civilians might be *higher*, not lower, than they are for using force against White civilians, even as the strength of evidence of that effect is debatable (Johnson et al., 2019; Knox and Mummolo, 2020).)

2.2 Challenges to Inferring Racial Bias

Different conceptions of racial bias can exist. On the one hand, we could focus on the bias of the patrol officer that shoots a civilian. On the other hand, we could focus on the police department (and supervisors) of the officer. As Bittner, p. 10 posited, "The ecological deployment of police work at the level of departmentally determined concentrations of deployment, as well as in terms of the orientations of individual police officers, reflects a whole range of public prejudices." For this study, we focus on bias by the patrol officer, acknowledging the potential of administrative control and bureaucratic bias to affect the context of police-civilian encounters (M. K. Brown, 1981). However, we must acknowledge that the "race" of an individual is not randomly realized during police encounters with

civilians.³ As a consequence, any inference about the causal effect of the race of a civilian on police use of force, or other police behaviors (e.g., driver or pedestrian stops) depends on the comparability of incidents.

Confounds in the use of force can be difficult to measure. Even if one can account for the lack of observed outcomes for officer-civilian encounters that never take place, empirical tests for racial bias still require accounting for confounds affecting contact and use of force (Knox, Lowe, and Mummolo, 2020). Race, for example, may be correlated with other characteristics (e.g., income, education, geography, employment, social networks) that might cause disparate rates of contact with police, thereby influencing civilian exposure to police use of force. Therefore, racially disparate patterns in the use of force and its severity may spuriously relate to characteristics of police-civilian encounters that explain use of force (e.g., Jetelina et al., 2017; Worrall et al., 2018; Knowles, Persico, and Todd, 2001; Cesario, Johnson, and Terrill, 2019). To best study the effect of race on the propensities of civilians to experience police use of force requires conditioning on a range of civilian characteristics that may confound the relationship. Furthermore, there is the matter of selection into contact with police and how it challenges inference-drawing about racial bias during citizen-police interactions (Johnson et al., 2019; Knox and Mummolo, 2020; Knox, Lowe, and Mummolo, 2020).

Assuming racial bias in police shootings exists, there are at least two theoretical mechanisms, one circumstantial and the other psychological (for a brief discussion, see Ross, 2015, p. 3). The first mechanism is that racial minorities, especially Black Americans, are circumstantially associated with conditions that give rise to police using greater force against them: They are more likely to come into contact with police because police officers racially profile them⁴ or they are more proximate to high-crime and/or highly-policed environments. The second mechanism is that police officers differentially perceive the stakes for using force against civilians depending on the race of the civilians. Officers might, for example, antici-

³By "race" of civilian, we mean the officer's perception of their race.

⁴Racial profiling as a mechanism of racial disparities in use of force, however, is potentially circular.

pate differential downstream consequences from using force against Black civilians than from using force against White civilians, or interpret behaviors differently for Black and White civilians. In its most nefarious expression, regardless of the race of the officer, police may devalue the lives of Black civilians relative to the lives of White civilians.

3 A Racial Bias Model of Police Shootings

Our racial bias model of police shootings stems from the model Knowles, Persico, and Todd (2001) employ to examine police stops of drivers. It seeks to capture "the transactional, or step-by-step unfolding, of police–public encounters" and the "micro process of the police-suspect encounter," in which civilian noncompliance, be it actual or perceived, can be pivotal to the decisions and discretion of police officers to use force (e.g., Terrill, 2005).

The first stage of our model is a selection stage. It allows for disparate rates of civilian encounters with police officers across civilian racial groups. Such an allowance is important. Encounters with police where civilians are "suspect" are unequal. Differences in the deployment of and exposure to police in the United States are historic, with some races (and racialized places) more than others receiving greater surveillance, intervention, and statesanctioned violence by the police, even when unmerited. In particular, studies from across the social and public health sciences of police contact with civilians, drawing on varied data from police records, public opinion surveys, face-to-face interviews and focus groups, demonstrate that, generally, police devote and Black civilians receive greater—often needless—attention relative to white civilians for the same activities (e.g., traffic and pedestrian stops and outcomes of searches for contraband) (Pierson, Simoiu, Overgoor, Corbett-Davies, Jenson, et al., 2020; Baumgartner, D. A. Epp, and Shoub, 2018; C. R. Epp, Maynard-Moody, and Haider-Markel, 2014; Prowse, V. M. Weaver, and Meares, 2020).

Modeling the first stage allows us to make empirical predictions about behavior *implied* by racial bias that should manifest even in the presence of selection into encounters with

the police. The selection stage captures, conceptually, every element of the police-civilian interaction that takes place up until the civilian and the officer reach the point of violence. It includes quotidian inequalities such as "attentional biases" to Black civilians in public and differential perceptions of "suspicious" and "threatening" civilian behavior by race of civilian (Eberhardt et al., 2004), along with the social construction of the "Black symbolic assailant" (Bell, 2017) and differences in civilian experiences with police discretion by skin color and phenotype (e.g., Monk, 2019; Kahn, Phillip Atiba Goff, et al., 2016).

In the second stage, we model a conflict subgame. It seeks to capture the kinds of splitsecond choices that police make at the point of using force. "During high-pressure situations, including some police-citizen encounters," however, "officers may not have the luxury of making slow, considered analytical decisions and, instead, rely on intuition and experience" (Hine et al., 2018). The same may be true for civilians. Nevertheless, the heightened pace of decision making, the urgency with which individuals, both civilian and police, respond to real or perceived threats to their dignity and physical safety, and the uncertainty about each other (e.g., does the civilian have a gun or a wallet), suggest this process is accurately captured by simultaneous structure.

In our model, conflict takes the form of escalating or accumulating aggression in the demeanor and deed of the civilian (actual or perceived by the officer) and the use of force by the officer, following initial interaction(s) between the civilian and officer (e.g., stopping the civilian, civilian disregard of verbal commands, etc.). We use "escalation" in a specific way civilian demeanor or deed perceived by a police officer to be threatening, where the real or misperceived aggression could "harm another person who is motivated to avoid that harm" (Allen and Anderson, 2017). It includes nonphysical non-compliance with police directives, inclusive of verbal hostility and antagonism (e.g., cursing or berating an officer) and physical non-compliance (e.g., turning from or striking an officer). Escalation by the civilian risks the dignity, respect, authority, and/or safety of an officer (or another civilian).

Although the choices of police during police-civilian encounters may partially result from

the demeanors and deeds of civilians, not all use of force, especially shootings by police, or civilian deaths by police are entirely or at all affected by civilian behavior. A civilian may comply with a directive from an officer, displaying neither defiance nor belligerence, but an officer may mistake or misperceive the behavior of the civilian and use deadly force. Examples include the 1967 and 2014 non-fatal shootings of Huey Newton and Levar Smith, and the 1999 and 2016 fatal shootings of Amadou Diallo and Philando Castille. Additionally, a civilian may be impaired by intoxicants or untreated mental illness, preventing them from making decisions or acting to reduce their appearance of threat to an officer (or other civilian), inclusive of non-response to police directives, resulting in civilian harm, inclusive of death (e.g., the fatal police shootings of Eleanor Bumpurs in 1984 and Daniel Prude in 2020). Plus, situational factors beyond the influence and control of civilians may influence shootings by police and civilian deaths by police. Informational priming by 911 dispatchers or other civilians, for instance, may exaggerate the degree of threat a "suspect" civilian poses for police, quickening lethal use of force by police when none was necessary (e.g., Tamir Rice, Breonna Taylor, and John Crawford). Also, training and socialization of police officers to expect immediate compliance and looming violence against them may influence the use of force (Oberfield, 2012; Sierra-Arévalo, 2021). Lastly, differences in the demeanor of police (e.g., tone, tenor, courtesy, and respect) when dealing with different civilians (Voigt et al., 2017; C. R. Epp, Maynard-Moody, and Haider-Markel, 2014) may play a role in civilian aggression, or at least test civilian patience and increase their aggravation with police during encounters.

From the perspective of the "objectively reasonable" officer, civilian escalation may heighten the stakes of conflict between civilians and officers during encounters. At a minimum, escalation can create "a type of strain that may also have situational effects, increasing officers' anger and frustration toward specific civilians within individual encounters" (Nix et al., 2017, p. 615), as well as strengthening their assumptions that escalation signifies danger and "a greater likelihood of violence" to occur during their encounters with civilians they deem aggressive.⁵

Together, perception, emotion(s), and assumptions likely account, in part, for the scholarly consensus that "noncompliant citizens face a greater likelihood of being treated disrespectfully by the police...[and] are more likely to experience other negative outcomes, such as arrest and the use of force" (Nix et al., 2017, p. 1155). We assume, therefore, that if civilian aggression may increase the severity of police use of force, it, *in part*, should increase the likelihood of death following police shooting a civilian. Studies that statistically associate the degree of civilian non-compliance (e.g., resistance) with police directives and the degree of police use of force against civilians buttress our assumption (e.g., Engel, Sobol, and Robert E Worden, 2000; J. H. Garner, C. D. Maxwell, and Heraux, 2002; Sun, Payne, and Wu, 2008; L. James, S. James, and Vila, 2018; McCluskey and Terrill, 2005; McElvain and Kposowa, 2008; Wheeler et al., 2017).

We model the possibility of racial bias by allowing officer perceptions of the cost of fatally shooting a civilian to vary by race of civilian. Our formal representation captures emotional reactions, anxiety and threat perception associated with racial bias and the use of force (e.g. Kleider, Parrott, and King, 2010; Nieuwenhuys, Savelsbergh, and Oudejans, 2012; Welch, 2007; Correll, Park, Judd, and Wittenbrink, 2002), along with a more dispassionate cost-benefit analysis by the officer about the *anticipated* consequences of killing a civilian.

3.1 Primitives

Players, sequence of play, and strategies. The model is played between a civilian, C, and an officer, O. The civilian is characterized by a type, which is a pair, $\tau = \langle \kappa, \rho \rangle$. This pair includes a racial identity, $\rho \in \{B, W\}$, and observable civilian characteristics, denoted $\kappa \in \mathbb{R}$. The latter include dress, demeanor, location, time, or any other characteristic.

⁵The likelihood of police use force may be greater, too, when officers have evidence a crime occurred (McCluskey Terrill, 2005; McCluskey, Terrill, Paoline, 2005; Sun Payne, 2004) and civilians possess weapons (Johnson, 2011; McCluskey et al., 2005; Sun Payne, 2004). However, there is conflicting evidence that the seriousness of an offense or crime greatly influences the likelihood of police use of force (Friedrich, 1980; Lawton, 2007).

We denote the probability density function of κ , conditional on ρ , as $g(\kappa|\rho)$. That is, the distribution of observable characteristics in the population can be different for any racial group. When we turn to the empirical implications of our model, we consider a population of civilians, \mathcal{P} , characterized by the density function, $g(\cdot)$, from whom the civilian in the interaction is drawn.

Figure 1 summarizes the play sequence. The game begins with the civilian, who engages in behavior the "objectively reasonable" officer could perceive questionable or suspicious. Crucially, the behavior the civilian engages in need not actually be suspicious; it may be any kind of activity that an officer has the ability to further investigate (e.g., "loitering" or "furtive movement"). Let $s \in \{0,1\}$ denote that choice, where s = 1 indicates the choice to engage in an activity, which could *potentially* be perceived as questionable or suspicious by an officer (or another civilian). If the "suspect" civilian chooses s = 0, the game ends. However, if the "suspect" civilian chooses s = 1, then the officer must use their discretion to decide whether to engage the civilian for purposes of order maintenance or law enforcement (e.g., stop-question-frisk). Let $l \in \{0,1\}$ denote this choice, with l = 1denoting engaging the civilian. If the officer chooses l = 0, the game ends; if he chooses l = 1, the game proceeds to the next stage, with simultaneous interactions by civilian and officer. Specifically, both players must decide how to engage the other, whereby each can choose behaviors that could escalate to violence. The civilian must choose to escalate or not, $t \in \{0, 1\}$, where t = 1 denotes escalating. (Reiterating an earlier point, escalation can be in the eye of the beholder, especially that of the police officer, influenced by different factors). The officer must choose whether to use lethal force or not, $f \in \{0, 1\}$, where f = 1denotes lethal force. If the officer chooses lethal force, the civilian dies with probability $\delta(t)$, where we assume $1 \ge \delta(1) > \delta(0) \ge 0$. That is, the probability the civilian dies when an officer uses lethal force is strictly greater when the civilian is escalating than when he is not, recognizing there can be exceptions, which we identified earlier. If neither player escalates conflict (i.e., t = 0 and f = 0), then less adverse, non-fatal outcomes follow. In either event,



Figure 1: Sequence of play in the model.

the game ends after these choices are made and payoffs are realized.

Let $\pi(\tau)$ denote a probability distribution over r, conditional on the civilian's type, $\tau = \langle \kappa, \rho \rangle$, and let $\sigma(\tau)$ denote a probability distribution over f conditional on the civilian's observable characteristics and race. A strategy profile for the civilian is, therefore, a tuple, $\mathcal{C} = \langle s, \pi(\tau) \rangle$, and a strategy profile for the officer is a tuple, $\mathcal{O} = \langle l, \sigma(\tau) \rangle$.

Preferences and utilities. The civilian has preferences over their behavior and the outcome of their interaction with the officer. Specifically, we assume that a civilian of type τ who chooses to engage in suspicious behavior, s = 1, receives a payoff $c(\tau) > 0$ if the officer chooses not to engage in law enforcement activity (i.e., l = 0). This source of utility represents the value of engaging in whatever kind of behavior a citizen of type τ would like to engage in, without having to deal with the police. This payoff can depend on the individual's type (i.e., her race and observable characteristics). If the officer chooses to engage, though, l = 1, then we assume the civilian's payoff depends on whether the officer chooses to apply lethal force or not, as well as whether the civilian chooses a behavior that escalates conflict. If the officer chooses l = 1, then the civilian pays a cost, $-w(\tau)$, where we assume $w(\tau) > 0$, $\forall \tau$. This source of utility represents the cost of being subjected to policing and, as with the value of potentially suspicious behavior, can depend on the civilian's type. In addition to the cost of being subjected to policing, we assume the civilian pays a cost $-d(\tau)$ if he dies. That is, if the officer chooses to use lethal force (i.e., f = 1), then the civilian pays, in expectation, $-\delta(r) \cdot d(\tau)$, where we assume $d(\tau) > 0$. This source of utility represents the cost associated with the loss of life, which can depend on civilian type—i.e., some civilians may value living more than others such as the suicidal. To avoid considering unreasonable situations, we assume that the cost of dying is worse than the cost of being subjected to policing for all types of civilians.

Assumption 1 (Civilians prefer not to die). $d(\tau) > w(\tau), \forall \tau$.

If the civilian escalates, and the officer chooses less-than-lethal force, we assume the civilian receives positive utility $b(\tau) > 0$. The source of utility represents the value of engaging in escalation against an officer and can vary by type. The civilian's expected utility function is given by:

$$EU_C(s,t|\tau) = \begin{cases} 0 & \text{if } s = 0\\ c(\tau) & \text{if } s = 1 \& l = 0\\ -w(\tau) & \text{if } s = 1 \& l = 1 \& t = 0 \& f = 0\\ b(\tau) - w(\tau) & \text{if } s = 1 \& l = 1 \& t = 1 \& f = 0\\ -w(\tau) - \delta(r) \cdot d(\tau) & \text{if } s = 1 \& l = 1 \& f = 1 \end{cases}$$

The officer has preferences over conducting policing work, stopping suspects and criminals, fatally wounding civilians, and his own physical well-being. Specifically, we assume the officer pays a cost $-c_O(\tau)$, where $c_O(\tau) \in (0, 1)$, whenever the civilian chooses to engage in potentially suspicious activity (i.e., s = 1) and the officer does not engage in law enforcement (i.e., l = 0). This cost represents the cost of allowing potentially criminal activity to go overlooked or a forsaking of duty. Importantly, we allow this cost to vary by civilian type, allowing an officer's disutility from permitting potentially criminal activity to occur is a function of everything the officer can observe about the civilian. In addition, the officer pays a cost $-k_{\rho}$, where we assume $k_{\rho} \in (0,1) \forall \rho$, whenever he fatally wounds a civilian of race ρ . By contrast, the officer pays a cost, $-d_O$, where $d_O > 0$ whenever a civilian is escalating and he does not use lethal force, (i.e., f = 0). Substantively, this cost can represent injury to the officer, disutility from not stopping a criminal who is acting aggressively, or another adverse consequence. Finally, we assume the officer receives positive utility 1 from using force to stop a civilian who is engaged in criminal activity and escalating the conflict. This represents the utility of exercising authority, maintaining order, and stopping a potentially dangerous person. Therefore, the officer's expected utility function is given by:

$$EU_O(\gamma, \lambda | \tau) = \begin{cases} -c_O(\tau) & \text{if } s = 1 \& l = 0\\ -d_O & \text{if } s = 1 \& l = 1 \& t = 1 \& f = 0\\ -w_O - \delta(0) \cdot k_\rho & \text{if } s = 1 \& l = 1 \& t = 0 \& f = 1\\ 1 - \delta(1) \cdot k_\rho & \text{if } s = 1 \& l = 1 \& t = 1 \& f = 1\\ 0 & \text{otherwise} \end{cases}$$

3.2 Analysis

We characterize a mixed-strategy subgame perfect Nash equilibrium. There can exist a pure strategy equilibrium if officers are never willing to use lethal force, which we rule implausible by assumption. For the officer to be willing to play a mixed strategy, the civilian must choose a probability distribution over her decision to escalate that makes the officer indifferent between using lethal force and not. There is a probability that satisfies this requirement:

$$\pi^*(\tau) = \frac{w_O + \delta(0)k_\rho}{1 + w_O + d_O - (\delta(1) - \delta(0))k_\rho} \tag{1}$$

Notice that $\pi^*(\tau)$ is increasing in k_{ρ} . As an officer perceives it to be costlier to kill a civilian of race ρ , the civilian will be more likely to escalate. In addition, $\pi^*(\tau)$ is decreasing in $(\delta(1) - \delta(0))$. Hence, as the civilian's behavior has a larger impact on the probability of dying when the officer uses force, the equilibrium probability of a civilian threatening will decrease. Intuitively, this makes sense, for if the civilian's behavior does not matter, fatality becomes irrelevant for his calculation, and fatality is the major factor deterring him from aggressive behavior. At the same time, the officer's equilibrium probability distribution over using lethal force, $\sigma^*(\tau)$, must make the civilian indifferent between choosing to escalate. That probability is given by:

$$\sigma^*(\tau) = \frac{b(\tau)}{b(\tau) + d(\delta(1) + \delta(0))} \tag{2}$$

Thus, in any equilibrium that reaches the conflict subgame, there exists a mixed-strategy subgame perfect Nash equilibrium where civilians probabilistically escalate and officers probabilistically use lethal force.⁶

Proposition 1. In any subgame perfect Nash equilibrium where players reach the aggressive behavior subgame, the civilian and officer play mixed strategies whereby a civilian of type $\tau = \langle \kappa, \rho \rangle$ chooses to escalate with probability $\pi^*(\tau)$, and the officer chooses use lethal force with probability $\sigma^*(\tau)$.

3.3 Empirical Implications

How does racial bias by police officers affect equilibrium behavior? We offer a simple definition of bias, guided by Knowles, Persico, and Todd, 2001. *Specifically, we say that an officer is racially biased if he perceives the cost of shooting an individual to vary by racial groups*: If an officer thinks it is less costly to shoot a Black civilian than a White civilian, then we say the officer is biased against Black civilians.

 $^{^{6}\}mathrm{In}$ the appendix, we show that the civilian and officer reach the conflict subgame under intuitive conditions.

Definition 1. An officer is racially biased if $k_B \neq k_W$. An officer is racially unbiased if $k_B = k = k_W$.

With this definition in hand, Proposition 1 is instructive about evidence of racial bias by police in OIS. Given Definition 1, we can identify the probability that a civilian should die, conditional upon being involved in an officer-involved shooting, when the police are not racially biased, and when they are racially biased.

Importantly, the model yields implications for how we can infer bias without having to make judgments about how to measure group traits, benefits to crime, or the distribution of traits in a group. That is, we are able to draw inferences from OIS outcomes *among those* who are actually involved in a shooting, without having data on the selection process that leads individuals into OIS events. Specifically, let $K(\rho)$ represent the set of characteristics for which an individual of race ρ would choose s = 1. Then, the fatality rate among people who are shot is given by

$$\mathcal{F}(\rho) = \int_{K(\rho)} \left(\delta(1) \cdot \pi^*(\tau) + \delta(0) \cdot (1 - \pi^*(\tau)) \right) \frac{\sigma^*(\tau) g(\kappa|\rho)}{\int_{K(\rho)} \sigma^*(z|\rho) g(z|\rho) dz} d\kappa$$
(3)

Notice that this fatality rate is not the fatality rate for all civilians of a given race but only for those who are shot by a police officer. Notice that by Definition 1, if an officer is not racially biased, then $k_B = k = k_W$. Given the civilian's equilibrium strategy, $\pi(\tau)^* = \frac{w_O + \delta(0)k_{\rho}}{1 + w_O + d_O - (\delta(1) - \delta(0))k_{\rho}}$, from above, then we can substitute $\frac{w_O + \delta(0)k_{\rho}}{1 + w_O + d_O - (\delta(1) - \delta(0))k_{\rho}}$ for $\pi^*(\tau)$. Because this quantity is independent of κ , Equation (3) reduces to

$$\mathcal{F}(\rho) = \delta(0) + (\delta(1) - \delta(0)) \left(\frac{w_O + \delta(0)k_\rho}{1 + w_O + d_O - (\delta(1) - \delta(0))k_\rho} \right)$$
(4)

Notice the only way this quantity varies with civilian race is if the officer's perceived cost of taking a civilian life varies by race. Therefore, differential fatality rates can only arise as a result of racially biased policing.

Proposition 2. In equilibrium, different fatality rates by racial groups arise only when the officer is racially biased.

The consequence is that if police are not racially biased then the probability a civilian is killed in an OIS, conditional on being involved in a shooting, should be independent of her race, even accounting for all other observable characteristics that might influence her incentive to engage in noncompliance or resistance, as well as the officer's incentive to use force in the first instance. That is, Equation 3 provides the theoretical foundations for a sufficient test of racial bias in the use of lethal force in OIS. It is important to underscore that this implication of our model allows us to evaluate evidence of racial bias, even taking into account unobservable behavioral differences across racial groups that might take place during a police-civilian encounter. This result is parallel in logic to the way Knowles, Persico, and Todd (2001) study racial disparities in traffic stops and Alesina and La Ferrara (2014) study bias in capital sentencing. It allows us to assess evidence of racial bias without having to measure observable or behavioral characteristics of either civilians or officers. It is sufficient to evaluate variation in ultimate consequences—namely, patterns of fatality.

Implication 1. If police officers are racially biased in favor of shooting Black civilians, then, conditional upon being involved in an officer-involved shooting, Black civilians will be less likely to die than will non-Black civilians.

The core logic underlying this implication is that officers will be more likely to use force in less dangerous situations involving Black civilians than in similar situations involving White civilians. As a consequence, a greater proportion of OIS involving Black civilians will not lead to a fatal outcome.

A corollary implication of our model is that White civilians should be more likely than Black civilians to engage in escalation and aggression towards officers. That implication is important, because it helps clarify the underlying theoretical mechanism we posit. Black civilians are induced to be more cautious during an interaction with police than are White civilians. **Implication 2.** If police officers are racially biased against shooting White civilians, then, conditional upon being subjected to law-enforcement activity, White civilians will be more likely to engage in threatening behavior, such as resisting arrest, disobeying officer commands, or behaving belligerently than will non-White civilians.

It is beyond the limits of this paper to fully investigate that implication due to insurmountable data limitations, particularly data on the perceptions and/or degree of civilian aggression, but its verisimilitude is important for establishing the mechanism that drives the analysis we present. To that end, we note that beyond anecdotal support for the mechanism, there is some evidence from extant literature to support the implication. Kavanagh (1997) studies more than 1,000 encounters between civilians and officers in New York City's Port Authority Bus Terminal between 1990 and 1991 and finds suggestive evidence that White civilians are more likely to resist arrest than are non-White civilians. Matrofski, Snipes, and Supina (1996) compare civilian-officer race combinations as predictors of civilian compliance with officer requests for orderly behavior. They find that, compared to White civilians interacting with White officers, White civilians interacting with minority officers are *less* likely to comply with officer instructions. At the same time, they find that minority civilians interacting with White officers are *more* likely to comply with officer instructions. They also find that minority civilians interacting with minority officers are more likely to comply, though this difference is not statistically significant. Finally, according to the FBI's Law Enforcement Officers Killed & Assaulted data, as of July 2017, 55% of officers killed by civilians were killed by White civilians and 58% of officers assaulted by civilians were assaulted by White civilians. While far from constituting a systematic evaluation, those descriptive findings provide initial evidence to corroborate the underlying mechanism we posit. However, for the remainder of the paper, we evaluate the primary implication of the mechanism articulated above.

4 Empirical Assessment

Our empirical assessment of the implications for racial bias in police shootings proceeds in four steps. First, we describe our method—the outcome test. Second, we describe an original dataset we built that includes all OIS (fatal and non-fatal) in eight local police jurisdictions but, due to data limitations imposed by police reporting, excludes sufficient data on civilian behavior (or police perceptions of it) during the police-civilian encounter. Third, we focus on an evaluation of Implication 1 that predicts that racial bias among police officers will produce disparities in fatalities across racial groups. We underscore that this prediction is not intended to estimate the effect of civilian race on the decision to use force; it is designed to demonstrate evidence *implied* by any such bias. In the fourth step, therefore, we directly engage the issue of causal effect size. Taking our evidence as consistent with the presence of racial bias as a starting point, we calculate a lower bound for the magnitude of the effect of racial bias in the decision of an officer to shoot a civilian in our sample of localities.

4.1 Discerning Racial Bias: The Outcome Test Method

To evaluate Implication 1, we employ an outcome or "hit rate" test, which is capable of observing disparate impact and identifying bias in decision-making (e.g., Knowles, Persico, and Todd, 2001; Persico and Todd, 2006; Alesina and La Ferrara, 2014). Mortgage lending illustrates the general logic of the approach. Mortgage lenders may care about timely repayment of loans. If we observe that non-White lendees repay mortgages on time at higher rates than Whites lendees, then that would suggest that qualified non-White applicants are being denied loans (Ayres, 2002). If the same standard were applied for mortgage lending, independent of borrowers' race, we should expect similar default rates across racial categories. However, because lenders were willing to lend to less qualified White borrowers than to Black borrowers, the default rate would be higher for White borrowers. For policing, we may see similar systematic differences by race, in the other direction. Stops may be considered successful, for instance, if they lead to arrest, perhaps because of the discovery of contraband or the harmful behavior of drivers. Gelman, Fagan, and Kiss (2007), for example, found that 1 in 7.9 Whites police stopped were arrested, compared to 1 in 9.5 Blacks. That suggests the discretion threshold police use to decide whom to stop is lower or more indiscriminate for Black drivers than for White drivers. Our logic similarly implies that if officers have a lower threshold for deciding to shoot Black civilians than White civilians, then there will be a greater proportion of Black civilians who will choose to not threaten and, therefore, survive an officer-involved shooting. Importantly, in many traditional settings, hit-rate tests are used to evaluate the presence of a latent trait to uncover evidence of bias. In our setting, as in Knowles, Persico, and Todd (2001), the latent trait is a choice by another player. In Knowles et al., drivers strategically choose whether to carry contraband; in our model, civilians strategically decide how to behave during police-civilian encounters. Anticipating bias by officers, Black civilians will be less likely in equilibrium to behave in ways that escalate a confrontation towards police-civilian violence than will White civilians. That feature is a result, not an assumption. The motivating assumption, as we noted above, is that the risk of death should be higher during a police-civilian encounter involving civilian aggression than one without civilian aggression.

A note on causality. Before presenting our analysis, we underscore the causal pathway at the heart of our argument. Our claim is not that racial bias directly causes differential fatality rates. Our argument is instead that racial bias causes officers to use force differently in different situations across racial groups. Anticipating that, civilians interact differently with officers in a way correlated with the civilian racial identities. The effects of those behaviors in conjunction is a distribution of force-civilian action combinations that vary by civilian race. Our analysis reveals that differential fatality rates is evidence consistent with that effect, not the effect itself. Just as we would not argue that differential default rates by race are a direct effect of racial bias in mortgage lending, we do not argue that differential fatality rates are a direct effect of racial bias in the decision to use force. Thus, as we proceed to our empirical analysis, we do not set out to demonstrate a causal effect of bias on fatality rates because the path from bias to fatality rates runs through myriad immeasurable intermediate mechanisms.

4.2 Data on Officer-Involved Shootings

To evaluate racial disparities in fatality rates among different racial groups, we require data on every single officer-involved shooting, not just fatal shootings. Data on OIS—even just fatal ones—are notoriously difficult to acquire (Zimring, 2017). Recent efforts have begun to compile extensive data on fatal encounters between officers and civilians. They typically rely on media reports and crowd-sourced data, making it difficult to assess how comprehensive and systematic the data are. Moreover, existing data typically do not include instances of OIS that do not include a fatality. Thus, we collected original OIS data by filing public records requests with individual police departments.

We sent public records requests to police departments and sheriffs' offices in the 50 largest local jurisdictions in the U.S., measured by population. We requested records of every single instance of an officer discharging their weapon between 2010 and 2017. Although most policing agencies were positively responsive to our requests, most policing agencies that responded with data did not provide racial information about civilians involved in OIS. Our data, therefore, comprise eight jurisdictions—Charlotte, Houston, King County, WA, Los Angeles, Orlando, San Jose, Seattle, and Tucson—that provided comprehensive racial information in response to our public records requests.⁷ The unit of analysis for each incident is the civilian/officer pair.⁸

We constructed all civilian/officer pairs, yielding 1,274 total pairs, representing 748

⁷Unfortunately, the departments could not provide objective data on observed officer interactions with civilians or civilian behavior during all interactions with police officers, and often not even subjective data for interactions involving use of force by officers and civilian behavior leading up to it.

⁸San Antonio also provided such information but the sample size was too small to make black/white comparisons. Results are available upon request.



Figure 2: Number of officer-involved shootings per month in eight cities, 2010-2017. The figure plots the (logged) number of officer-involved shootings each month in each city.

unique incidents. Overall, 48% of our OIS incidents represent fatal shootings, varying considerably by department. Charlotte had the lowest rate of fatalities from OIS, where 9 out of 45 observations were fatal (20%). Los Angeles had the highest number of reported OIS (663), where 58% of them were fatal. Our data demonstrate we have considerable variation in officer-involved shooting incidents, not just by department and by time (see Figure 2) but by fatality, too.

Figure 2 shows the frequency of OIS in each of jurisdictions. Because there is considerable variation in the size of the jurisdictions, there is considerable variation in the total number of OIS. The most come from Los Angeles, the second-largest police jurisdiction in the country. Therefore, we log the number of observations per month to prevent scale differences from skewing the temporal patterns and cross-jurisdiction variation. Notably, with the exception of an increase in OIS in Houston at the end of the series, there is little within-city variation in the frequency of OIS.

Furthermore, the spatial distribution and concentration of OIS within jurisdictions show



Figure 3: Locations of fatal shootings (black dots) and non-fatal shootings (white dots) in our sample of eight locations. The black triangles mark Level I Trauma Centers.

intuitive but instructive patterns. Figure 3 shows the distribution of fatal and non-fatal shootings in our eight jurisdictions. Los Angeles and Houston, by far the largest localities in our dataset, experience the most OIS, whereas cities like Charlotte and Tucson experience relatively few. Additionally, it appears there is a higher fatality rate among OIS in localities like Los Angeles and Houston, which is less of an issue in jurisdictions like San Antonio and Charlotte. Overall, Figure 3 highlights the geographical diversity in these fatal OIS, that they do not appear to systematically occur in only certain parts of certain localities, and that fatality rates vary across geographies.

4.3 Analysis and Results

We begin our empirical analysis of Implication 1 by simply comparing the distribution of fatality across racial groups, conditional on being involved in an officer-involved shooting. Table 1 summarizes the frequencies among the observations in our data. The columns break down OIS by the race of the civilian involved, and the rows distinguish between fatal and non-fatal OIS.

	White	Black	Hispanic	Asian
Not Fatal	118 (48%)	329 (67%)	208 (42%)	11 (27%)
Fatal	126~(52%)	162 (33%)	290~(58%)	30~(73%)

Table 1: Summary of officer-involved shootings by race and fatality. $\chi^2 = 76.888, p \leq 0.001.$

The evidence is startling, revealing considerable dependence between fatalities and the race of the civilian ($\chi^2 = 76.888$, $p \leq 0.001$). In particular, a majority of Black civilians survive OIS, whereas a majority of civilians of all other races do not. Of course, demographics and police behavior both vary across jurisdictions, and we might worry that the correlation detected in Table 1 is spurious. To speak to this we estimate a series of logistic regression specifications on all observations of OIS for which the departments we sampled provided race information. The unit of analysis is the civilian involved in an officer-involved shooting, and the the outcome variable is an indicator for whether the civilian was fatally wounded.

For 17 observations, the outcome was recorded as "Undetermined" or "Unknown." We treat these observations as missing data. Our primary explanatory variable of interest is the race of the civilian involved.

We also consider specifications where we include as explanatory variables the distance from each officer-involved shooting to the nearest trauma center as well as year fixed effects (see Table 6 for the specifications with year fixed effects).⁹ We also include fixed effects for the cities from which we have data, which are likely correlated with the distance to trauma center and the racial indicator. This is because trauma centers have fixed locations in cities, and demographic characteristics of populations vary across cities. Unfortunately, for 24 of our 1,274 observations, the address of the officer-involved shooting was too imprecise to calculate a reliable distance measure. We consider specifications both with and without this control variable.

The main results of our analysis are reported in Table ??. The primary result appears in the top row. In each of our specifications, among those civilians shot by an officer, Black civilians are less likely to die than are White civilians. This difference is statistically significant in each specification. In our main specification, reported in the first column of results, White civilians have a predicted probability of 0.52 of dying, whereas Black civilians have a predicted probability of 0.33—a 19 percentage point decrease. This relationship supports the primary empirical implication of our theoretical model of racial bias. It is consistent with the claim that police officers have a lower threshold for deciding to use lethal force against Black civilians than against White civilians. Notably, the magnitude of the relationship between being a Black civilian and the probability of dying *increases* once we include jurisdiction fixed-effects, and maintains when we include year fixed-effects. What is more, the relationship between being a Hispanic civilian and a reduced probability of dying does not emerge even after we include jurisdiction and year fixed effects. This functions as

⁹Some observations lacked adequate location information to calculate the distance to the nearest trauma center, which has been shown to be a particularly important factor for the chances of survival of a gunshot wound (Crandall et al., 2013). Therefore, in the models including distance to the nearest trauma center as a control variable, we only have 1250 observations, covering 729 unique incidents.

a placebo test and implies that any problematic unmeasured covariates would have to have different relationships for Black and Hispanic civilians (e.g., concerns about characteristics that affect the probability of death—such as police behavior, training, and medical attention would be largely ruled out by this analysis).

	Model 1	Model 2	Model 3	Model 4
Black	-0.77^{*}	-0.70^{*}	-0.74^{*}	-0.67^{*}
	(0.32)	(0.33)	(0.32)	(0.33)
Hispanic	0.27	0.10	0.29	0.13
	(0.31)	(0.32)	(0.31)	(0.31)
Asian/AI/AN/PI	0.94	0.82	0.97	0.91
	(0.62)	(0.58)	(0.63)	(0.58)
Distance			0.14	0.41
			(0.19)	(0.21)
Houston		0.00		-0.21
		(0.57)		(0.63)
King County		0.27		-0.12
		(0.76)		(0.81)
Los Angeles		1.26^{*}		1.15
		(0.57)		(0.62)
Orlando		0.59		0.52
		(0.65)		(0.70)
San Jose		0.18		0.16
		(0.65)		(0.69)
Seattle		1.25		1.27
		(0.73)		(0.77)
Tucson		1.61^{*}		1.64^{*}
		(0.69)		(0.73)
Intercept	0.07	-0.80	-0.04	-1.00
	(0.25)	(0.59)	(0.28)	(0.64)
Num. obs.	1274	1274	1250	1250
**** 0.001 ** 0.01				

***p < 0.001; **p < 0.01; *p < 0.05

Table 2: Estimated relation	tionship between ci	vilian race and pr	obability of fatal	ity conditional
upon being involved in a	n officer-involved sl	nooting. Cells sho	w logit coefficient	s with cluster-
robust standard errors.	Omitted category	is White civilians	and Charlotte.	Distance is in
tens of miles.				

As we do not observe a depression of the relationship between being a Black civilian and the probability of survival after we include jurisdiction and time fixed effects, a spurious correlation between race and jurisdiction does not drive the observed relationship. This pattern—while not necessarily causal—is precisely what we expect if police are racially biased in favor of shooting Black civilians, given the logic of our model. In order to explore the possibility that the relationship would be eliminated by an omitted covariate, we conduct a number of sensitivity analyses based on the methodology presented in VanderWeele and Ding (2017) and Cinelli and Hazlett (2020).¹⁰ These analyses consider how strong an unmeasured confounding variable would have to be in order to wipe out the effects we are finding for Black civilians. One method to measure such strength is to benchmark any potential unmeasured confounder against measured covariates in the model. In analysis presented below and in the Appendix, we show that in order to eliminate the effect, there would need to be an unmeasured confound that is more than three times as strong as any of the variables currently included in the model (jurisdiction fixed effects, time fixed effects, and distance to trauma center). We have not been able to identify any such missing variable that would affect fatality rates for Black civilians and not Hispanic civilians.

4.4 How Big of an Effect Could Racial Bias Have on Officer-Involved Shootings?

Our analysis revealed evidence consistent with racial bias, per our definition, in the decision of police officers to use lethal force. However, we have not directly estimated a causal effect of a civilian's race on the decision to use force. That means we still have to quantify the size of the bias, substantively. Accordingly, we estimate a lower bound on the magnitude of racial bias in OIS, relying on logic paralleling Knox, Lowe, and Mummolo (2020), Cohen (2021), and Cohen and Glynn (2021) for identifying racial bias in police contact with civilians. The approach we adopt has two steps. First, we define the fatality rate for Black civilians that police shot, comprising two components—those that would not have been shot had they

¹⁰Although the Cinelli and Hazlett (2020) analysis is based on a linear probability model, we generally find small differences for this data between analyses based on the logit model and the linear model.

been White and those would would have also been shot were they White. Second, we define the fatality rates of groups relative to each other.

The magnitude of racial bias in the decision to shoot a civilian is the proportion of Black civilians shot who would not have been shot had they been White. The intuition behind this is that the observed fatality rate of Black civilians is made up of two components — Black civilians who were shot but *would not* have been shot had they been White and Black civilians who *would* have been shot had they instead been White. Our quantity of interest p, is the proportion that are in the former i.e. the proportion of Black civilians shot, who would not have been shot had they been White. By using the principal strata defining these groups we can derive a lower-bound for p as the ratio of a difference of fatality rates: rate for Blacks shot who would have been shot if White minus rate for Blacks who would not have been shot if White minus rate for Blacks who would not have been shot if White.

$$p = \frac{\mathcal{F}_{s(b)=s(w),b} - \mathcal{F}_b}{\mathcal{F}_{s(b)=s(w),b} - \mathcal{F}_{s(b)>s(w),b}}.$$
(5)

Equation (5) contains counterfactual quantities, so to derive an empirically estimable lower-bound for p, we assume that the fatality rate for Whites is no greater than the fatality rate for Blacks who would have been shot if White (Cohen and Glynn, 2021). Furthermore, we do not observe $\mathcal{F}_{s(b)>s(w)}$, however, we know if there is no racial bias then $\mathcal{F}_{s(b)>s(w)}$ would be 0, as there would be no civilians shot because they were Black. Therefore, substituting 0 for $\mathcal{F}_{s(b)>s(w)}$ yields a lower bound on the true value of p.

$$p \ge \frac{\mathcal{F}_w - \mathcal{F}_b}{\mathcal{F}_w - \mathcal{F}_{s(b) > s(w), b}} \ge \frac{\mathcal{F}_w - \mathcal{F}_b}{\mathcal{F}_w}.$$
(6)

Equation (6) expresses p as a function of \mathcal{F}_w and \mathcal{F}_b , the observed fatality rates among White and Black civilians shot. See the Appendix for formal assumptions, definitions and derivation of our quantity of interest p. To estimate the lower bound on p, we first estimate a logistic model. We estimate it with a subset of officer-involved shooting data containing only Black and White civilians, including our main covariate of interest, namely race (White equal to 1, Black equal to 0), along with binary indicator variables for locality and a continuous variable of distance to closest trauma center in miles. Using this model we estimate the regression coefficient on White to be 0.70 (see the Appendix full regression specification results in Table 11), the associated fatality difference between White civilians and Black civilians controlling for city fixed effects. The lower bound estimate p follows from the estimated risk ratios (Cohen, 2021) as in equation 7 (see the appendix for the derivation).

$$p \ge 1 - \frac{1}{\hat{RR}} \tag{7}$$

. Thus, we estimate 30% is the lower bound on the proportion of Black civilians that police would not have shot had they been White. Potentially there are unmeasured confounders that would affect both race and the likelihood of being fatally shot. As a sensitivity analysis we use the techniques of VanderWeele and Ding (2017) and Cohen (2021). These analyses indicate that an omitted covariate would have to produce a bias factor three times stronger than any covariate in our data set.

Substantively, our estimate of 30% is considerable and given it is a lower bound, may be higher. Our estimate implies that police would not have shot 156 Black civilians had they been White, from the 497 Black civilians in our eight localities over the years we study. Extrapolating this estimate to the larger population of the United States, however, is beyond the limits of our data. Moreover, significant intra-locality variation suggests police behavior, measured by OIS, is not uniform across the country. Additionally, comparing Hispanic civilians and Asian civilians to White civilians yielded no statistically significant differences. That is consistent with what we would expect—police officers differentially exercise discretion against Black civilians as compared to all other groups. Given the extant debate about whether the use of force by police is tainted with racial bias, these findings suggest there is a substantively significant problem. Quantifying the magnitude of its effect, though, requires richer administrative data beyond what police departments, generally, in the U.S. currently provide. Specifically, the important matter of how much police violence is attributable to racial bias requires knowing how often police fire their weapons, as well as how often they draw their weapons (e.g., Worrall et al., 2018), which is not universally known across local police departments.

5 Discussion

A significant challenge to credible inferences about the influence of racial bias in policing is that empirical observations typically need to condition on a wide range of difficult-to-measure confounds. For example, if civilian race is correlated with factors that directly affect contact with police—such as income, locality, employment rates or sectors, education level, or any possible factor—then it will be challenging to disentangle the causal effect of one's race from the effects of those other confounding forces. However, our approach helps overcome that challenge by identifying an empirical implication of racial bias in the use of force that is conditional on contact with the police, allowing social scientists to sidestep the challenges of selection bias due to racial rates of police contact with civilians (e.g., Knox, Lowe, and Mummolo, 2020).

What is more, our theory, analysis, and results help make better sense of seemingly contradictory findings in the contemporary use of force literature. For example, some studies show that the probability of being Black, conditional on being shot, is not statistically different from the probability of being White, conditional on being shot (Johnson et al., 2019). In our theoretical model, however, this pattern is completely consistent with racial bias by officers in favor of shooting Black civilians. Such a pattern could emerge because Black civilians are aware of such bias and systematically avoid aggression during encounters with the police that could lead to fatal OIS. Therefore, the probability of being shot, conditional on being Black, might still be higher than it is conditional on being White, even while the observed rates of being fatally wounded are the same. Similarly, our analysis can reconcile the distinction Fryer Jr. (2016) documents between lethal and non-lethal force against civilians.¹¹ If Black civilians are aware (or believe) that police officers are biased in favor of using force against them, then they should be less likely to engage in threatening behavior that would escalate a situation from a non-lethal outcome to a lethal outcome. We would expect, then, that Black civilians should be disproportionately subject to non-lethal force but not necessarily disproportionately represented in lethal encounters with police.

At the same time, while our analysis helps explain racial differences across the observed patterns in police use of force, all we can demonstrate is evidence consistent with racial bias and calculate a lower bound on the magnitude of the effect. The primary implication of our model, and the one we subject to empirical scrutiny, is a statement of an empirical regularity that is implied if civilians and officers behave as though the latter are racially biased. Lower fatality rates among Black civilians shot by the police than among White civilians shot by the police are a secondary form of evidence—a pattern implied by racial bias in the decision to shoot in the first instance. Those rates, however, do not in-and-of-themselves tell us anything about the magnitude of the effect of bias.

However, given what we know about the existence of racial bias, we are able to calculate a lower bound on the effect size. Still, the bounds we estimate cannot tell us about the upper limit on the effect. Doing that would require we overcome the aforementioned confounding and selection challenges to inference. While not necessarily an impossible task, undertaking it remains one of the most salient limitations research on the subject faces. As we document in the appendix, our model also helps identify the kinds of assumptions or data that would be necessary to make further progress on narrowing the estimated size of racial bias in the decision to shoot a civilian.

As we noted above, we have not investigated Implication 2. Doing so would require objective data on observed officer interactions with civilians. In particular, we would need

 $^{^{11}}$ Of course, Knox, Lowe, and Mummolo (2020) also suggest that the analysis in Fryer Jr. (2016) is flawed due to selection bias.

data on civilian behavior during all interactions with police officers, not just those involving use of force by officers. Such data are difficult to come by. However, it bears noting that there is some evidence in the extant literature that is potentially consistent with the expectation. It predicts that, if officers are racially biased against Black civilians, White civilians will be more likely to engage in escalating behavior than will Black civilians. While doing so would require the collection of rich new data that are not currently available, we believe it is a worthy endeavor as scholars continue to work out the mechanisms underlying disparate outcomes in civilian-officer interactions.

6 Conclusion

Police-civilian encounters have special implications for the study of democratic governance and equality of citizenship. Police are perhaps the most common government official with whom civilians have contact (e.g., Jacob, 1972) and, distinct from other bureaucrats, interactions with police officers always have the potential for violence. Consequently, the modal contact a civilian has with police relative to other government agents in the United States is one that might involve the use of physical force, including fatal and non-fatal shootings. Yet, whether justified or not, whether garnering mass and elite attention or not, whether we know enough or not about correlates and causes, police shootings (and other forms of police use of force such as use of compliance holds, pepper-spray, and canines) are moments that "raise fundamental questions of governmental responsiveness and state power, and they are frequently at the heart of grievances that generate political demands and protests" (Soss and V. Weaver, 2016, p. 83). Police shootings, along with predatory and extractive policing (Sances and You, 2017), police "militarization" (Lawson Jr., 2019), and broader practices of policing, inclusive of surveillance, order maintenance, and arrests, coupled with choices by local prosecutors and judges (e.g., requiring bail and jailing arrestees for low-level offenses), invite political scientists to ask "questions about police authority, state projects of social control, and daily encounters with local governance" (Soss and V. Weaver, 2017, p. 568). They also invite questions about the influence of bias, especially racial bias.

Racial bias on the part of government officials has the distinct potential to undermine the legitimacy of the state and civilian cooperation and engagement with government. To the extent, then, that police officers engage in racially biased use of force, that behavior has potentially profound consequences for the maintenance of a well-functioning democratic order. In light of these observations, recent analyses of racial disparities in the use of force by police officers have set out to address whether and how much racial bias influences policing in the United States. The implications of the findings are far-reaching.

Our results raise concern about racial bias in the use of force by police. They also highlight the need for more research and more comprehensive data about OIS, including, among other things, officer attributes and situational and contextual factors. For example, to understand the mechanisms by which racial bias affects civilian and police behavior, scholars need to study all civilian interactions with police, not just those encounters ending in fatalities, or even just the encounters where the use of force occurred. Of course, as others have pointed out (e.g., Knox, Lowe, and Mummolo, 2020) and as our model considers, there is potentially racial bias in the initial selection of civilians into contact with police. To the extent racial bias systematically affects not just how police interact with civilians but which civilians they interact with, our analysis underscores the extent to which training, recruiting, and monitoring of police officers have implications beyond public and officer safety.

Although our empirical study provides evidence *consistent* with racial bias in the use of force and a lower bound on the magnitude of racial bias in the decision to shoot, more research is necessary to assess the magnitude of the effect. We also need more research on racial bias in policing to assess the efficacy of policies designed to minimize racial disparities in policing, as well as to determine the underlying mechanisms that produce such racial bias. While normatively we might believe that, independent of its cause, racial disparities are problematic, what to do about them depends on identifying the root cause. In particular, whether racial disparities are a result of circumstantial factors or systematic bias by police officers affects what kinds of remedies are desirable and the implications of the disparities for the legitimacy and integrity of the police as a key law enforcement institution.

But better research will require richer administrative data on police practices, ranging across both the use of force continuum (e.g., no guns, guns drawn, guns fired) and outcomes (i.e., lethal and non-lethal consequences), as well as civilian behavior (e.g., resistance). The current nature and contents of use of force and consequences record-keeping by many police departments, however, presents serious challenges to improving research and establishing consensus in weighting across the varied factors associated with officer-involved shootings. Decentralization of law enforcement and varied discretion across localities in the United States further complicates research. Nonetheless, police departments, elected officials, and institutions of civilian oversight of police departments may become more interested in research about policing practices and outcomes, more anticipatory of scholarly needs, more transparent about and willing to share data with scholars and others through digitization and open-access, and interested in replication and extension of academic studies. If so, causal research on police behavior, from the spectacular to the mundane, may flourish, perhaps improving policymaking for public safety and improving policing (and police legitimacy) in the United States.

References

- Alesina, Alberto and Eliana La Ferrara (2014). "A Test of Racial Bias in Capital Sentencing".
 In: American Economic Review 104.11, pp. 3397–3433.
- Allen, Johnie J and Craig A Anderson (2017). "Aggression and violence: Definitions and distinctions". In: The Wiley handbook of violence and aggression, pp. 1–14.
- Alpert, Geoffrey P. and Roger G. Dunham (2004). Understanding Police Use of Force: Officers, Suspects, and Reciprocity. Cambridge University Press.

- Ayres, Ian (2002). "Outcome Tests of Racial Disparities in Police Practices". In: Justice Research and Policy 4.1-2, pp. 131–142.
- Baumgartner, Frank R., Derek A. Epp, and Kelsey Shoub (2018). Suspect Citizens: What 20 Million Traffic Stops Tell Us About Policing and Race. Cambridge University Press.
- Bell, Jeannine (2017). "Dead canaries in the coal mines: The symbolic assailant revisited".In: Ga. St. UL Rev. 34, p. 513.
- Binder, Arnold and Peter Scharf (1980). "The Violent Police-Citizen Encounter". In: *The* ANNALS of the American Academy of Political and Social Science 452.1, pp. 111–121.
- Bittner, Egon (1970). The functions of the police in modern society: A review of background factors, current practices, and possible role models. 2059. National Institute of Mental Health, Center for Studies of Crime and Delinquency.
- Brown, Michael K. (1981). Working the street: Police discretion and the dilemmas of reform.Russell Sage Foundation.
- Brown, Robert A. (2019). "Policing in American History". In: Du Bois Review: Social Science Research on Race 16.1, pp. 189–195.
- Buehler, James W. (2017). "Racial/Ethnic Disparities in the Use of Lethal Force by US Police, 2010–2014". In: American Journal of Public Health 107.2, pp. 295–297.
- Cesario, Joseph, David J. Johnson, and William Terrill (2019). "Is There Evidence of Racial Disparity in Police Use of Deadly Force? Analyses of Officer-Involved Fatal Shootings in 2015–2016". In: Social Psychological and Personality Science 10.5, pp. 586–595.
- Cinelli, Carlos and Chad Hazlett (2020). "Making sense of sensitivity: Extending omitted variable bias". In: Journal of the Royal Statistical Society: Series B (Statistical Methodology) 82.1, pp. 39–67.
- Cohen, Elisha (2021). "Sensitivity Analysis for Outcome Tests with Binary Data".
- Cohen, Elisha and Adam Glynn (2021). "Estimating Bounds on Selection Bias with Outcome Tests".

- Correll, Joshua, Bernadette Park, Charles M. Judd, and Bernd Wittenbrink (2002). "The Police Officer's Dilemma: Using Ethnicity to Disambiguate Potentially Threatening Individuals". In: Journal of Personality and Social Psychology 83.6, p. 1314.
- Correll, Joshua, Bernadette Park, Charles M. Judd, Bernd Wittenbrink, et al. (2007). "Across the Thin Blue Line: Police Officers and Racial Bias in the Decision to Shoot". In: *Journal* of Personality and Social Psychology 92.6, p. 1006.
- Crandall, Marie et al. (2013). "Trauma Deserts: Distance from a Trauma Center, Transport Times, and Mortality from Gunshot Wounds in Chicago". In: American Journal of Public Health 103.6, pp. 1103–1109.
- Eberhardt, Jennifer L. et al. (2004). "Seeing Black: Race, Crime, and Visual Processing".In: Journal of Personality and Social Psychology 87.6, p. 876.
- Engel, Robin Shepard and Jennifer M. Calnon (2004). "Examining the Influence of Drivers' Characteristics During Traffic Stops with Police: Results from a National Survey". In: *Justice Quarterly* 21.1, pp. 49–90.
- Engel, Robin Shepard, James J Sobol, and Robert E Worden (2000). "Further exploration of the demeanor hypothesis: The interaction effects of suspects' characteristics and demeanor on police behavior". In: Justice quarterly 17.2, pp. 235–258.
- Epp, Charles R, Steven Maynard-Moody, and Donald P Haider-Markel (2014). Pulled over: How police stops define race and citizenship. University of Chicago Press.
- Fryer Jr., Roland G. (2016). An Empirical Analysis of Racial Differences in Police Use of Force. Tech. rep. National Bureau of Economic Research.
- Garner, Joel H., Christopher D. Maxwell, and Cedrick G. Heraux (2002). "Characteristics Associated with the Prevalence and Severity of Force Used by the Police". In: Justice Quarterly 19.4, pp. 705–746.
- Garner, Joel H., Thomas Schade, et al. (1995). "Measuring the Continuum of Force Used by and Against the Police". In: *Criminal Justice Review* 20.2, pp. 146–168.

- Garner, Joel and Christopher Maxwell (1999). Measuring the Amount of Force Used by and Against the Police in Six Jurisdictions. Tech. rep. Bureau of Justice Statistics and National Institute of Justice.
- Gelman, Andrew, Jeffrey Fagan, and Alex Kiss (2007). "An Analysis of the New York City Police Department's "Stop-and-Frisk" Policy in the Context of Claims of Racial Bias".
 In: Journal of the American Statistical Association 102.479, pp. 813–823.
- Goel, Sharad, Justin M. Rao, and Ravi Shroff (2016). "Precinct or Prejudice? Understanding Racial Disparities in New York City's Stop-and-Frisk Policy". In: *The Annals of Applied Statistics* 10.1, pp. 365–394.
- Goff, Philip Atiba et al. (2016). The Science of Justice: Race, Arrests, and Police Use of Force. Center for Policing Equity.
- Hine, Kelly A et al. (2018). "Exploring police use of force decision-making processes and impairments using a naturalistic decision-making approach". In: *Criminal justice and Behavior* 45.11, pp. 1782–1801.
- Jacob, Herbert (1972). "Contact with Government Agencies: A Preliminary Analysis of the Distribution of Government Services". In: Midwest Journal of Political Science, pp. 123– 146.
- James, Lois, Stephen James, and Bryan Vila (2018). "Testing the impact of citizen characteristics and demeanor on police officer behavior in potentially violent encounters". In: *Policing: An International Journal.*
- James, Lois, Bryan Vila, and Kenn Daratha (2013). "Results from Experimental Trials Testing Participant Responses to White, Hispanic and Black Suspects in High-Fidelity Deadly Force Judgment and Decision-Making Simulations". In: Journal of Experimental Criminology 9.2, pp. 189–212.
- Jetelina, Katelyn K et al. (2017). "Dissecting the complexities of the relationship between Police Officer-civilian race/Ethnicity dyads and less-than-Lethal Use of Force". In: American Journal of Public Health 107.7, pp. 1164–1170.

- Johnson, David J. et al. (2019). "Officer Characteristics and Racial Disparities in Fatal Officer-Involved Shootings". In: Proceedings of the National Academy of Sciences, pp. 15877– 15882.
- Kahn, Kimberly Barsamian, Phillip Atiba Goff, et al. (2016). "Protecting whiteness: White phenotypic racial stereotypicality reduces police use of force". In: Social Psychological and Personality Science 7.5, pp. 403–411.
- Kahn, Kimberly Barsamian, Joel S. Steele, et al. (2017). "How suspect race affects police use of force in an interaction over time." In: *Law and human behavior* 41.2, p. 117.
- Kavanagh, John (1997). "The Occurrence of Resisting Arrest in Arrest Encounters: A Study of Police-Citizen Violence". In: *Criminal Justice Review* 22.1, pp. 16–33.
- Kleider, Heather M., Dominic J. Parrott, and Tricia Z. King (2010). "Shooting Behaviour: How Working Memory and Negative Emotionality Influence Police Officer Shoot Decisions". In: Applied Cognitive Psychology 24.5, pp. 707–717.
- Knowles, John, Nicola Persico, and Petra Todd (2001). "Racial Bias in Motor Vehicle Searches: Theory and Evidence". In: *Journal of Political Economy* 109.1, pp. 203–229.
- Knox, Dean, Will Lowe, and Jonathan Mummolo (2020). "Administrative records mask racially biased policing". In: *American Political Science Review* 114.3, pp. 619–637.
- Knox, Dean and Jonathan Mummolo (2020). "Making Inferences about Racial Disparities in Police Violence". In: Proceedings of the National Academy of Sciences 117 (3), pp. 1261– 1262.
- Lawson Jr., Edward (2019). "TRENDS: Police Militarization and the Use of Lethal Force". In: Political Research Quarterly 72.1, pp. 177–189.
- Lipsky, Michael (1980). Street-Level Bureaucracy: Dilemmas of the Individual in Public Service. Russell Sage Foundation.
- Matrofski, Stephen D., Jeffrey B. Snipes, and Anne E. Supina (1996). "Compliance on Demand: The Public's Response to Specific Police Requests". In: Journal of Research in Crime and Delinquency 33.3, pp. 269–305.

- McCluskey, John D. and William Terrill (2005). "Departmental and Citizen Complaints as Predictors of Police Coercion". In: *Policing: An InterNational Journal of Police Strategies* & Management 28.3, pp. 513–529.
- McElvain, James P. and Augustine J. Kposowa (2008). "Police Officer Characteristics and the Likelihood of Using Deadly Force". In: *Criminal Justice and Behavior* 35.4, pp. 505– 521.
- Mekawi, Yara and Konrad Bresin (2015). "Is the evidence from racial bias shooting task studies a smoking gun? Results from a meta-analysis". In: Journal of Experimental Social Psychology 61, pp. 120–130.
- Monk, Ellis P (2019). "The color of punishment: African Americans, skin tone, and the criminal justice system". In: *Ethnic and Racial Studies* 42.10, pp. 1593–1612.
- Mummolo, Jonathan (2018). "Modern Police Tactics, Police-Citizen Interactions, and the Prospects for Reform". In: *The Journal of Politics* 80.1, pp. 1–15.
- Nieuwenhuys, Arne, Geert J.P. Savelsbergh, and Raôul R.D. Oudejans (2012). "Shoot or Don't Shoot? Why Police Officers Are More Inclined to Shoot When They Are Anxious". In: *Emotion* 12.4, p. 827.
- Nix, Justin et al. (2017). "Demeanor, race, and police perceptions of procedural justice: Evidence from two randomized experiments". In: *Justice quarterly* 34.7, pp. 1154–1183.
- Oberfield, Zachary W (2012). "Socialization and self-selection: How police officers develop their views about using force". In: *Administration & Society* 44.6, pp. 702–730.
- Persico, Nicola and Petra Todd (2006). "Generalising the Hit Rates Test for Racial Bias in Law Enforcement, with an Application to Vehicle Searches in Wichita". In: *The Economic Journal* 116.515, F351–F367.
- Pierson, Emma, Camelia Simoiu, Jan Overgoor, Sam Corbett-Davies, Daniel Jenson, et al. (2020). "A large-scale analysis of racial disparities in police stops across the United States". In: Nature human behaviour 4.7, pp. 736–745.

- Pierson, Emma, Camelia Simoiu, Jan Overgoor, Sam Corbett-Davies, Vignesh Ramachandran, et al. (2017). "A Large-Scale Analysis of Racial Disparities in Police Stops across the United States". Stanford University Working Paper.
- Prowse, Gwen, Vesla M Weaver, and Tracey L Meares (2020). "The state from below: Distorted responsiveness in policed communities". In: Urban Affairs Review 56.5, pp. 1423– 1471.
- Ross, Cody T. (2015). "A Multi-Level Bayesian Analysis of Racial Bias in Police Shootings at the County-Level in the United States, 2011–2014". In: *PLoS One* 10.11, e0141854.
- Sances, Michael W. and Hye Young You (2017). "Who Pays for Government? Descriptive Representation and Exploitative Revenue Sources". In: *The Journal of Politics* 79.3, pp. 1090–1094.
- Schuck, Amie M. (2004). "The Masking of Racial and Ethnic Disparity in Police Use of Physical Force: The Effects of Gender and Custody Status". In: Journal of Criminal Justice 32.6, pp. 557–564.
- Sierra-Arévalo, Michael (2021). "American policing and the danger imperative". In: Law & Society Review 55.1, pp. 70–103.
- Sikora, Andrew G. and Michael Mulvihill (2002). "Trends in Mortality Due to Legal Intervention in the United States, 1979 Through 1997". In: American Journal of Public Health 92.5, pp. 841–843.
- Soss, Joe and Vesla Weaver (2016). "Learning from Ferguson". In: The Double Bind: The Politics of Racial & Class Inequalities in the Americas. Report of the APSA Task Force on Racial & Class Inequalities in the Americas.
- (2017). "Police Are Our Government: Politics, Political Science, and the Policing of Race-Class Subjugated Communities". In: Annual Review of Political Science 20, pp. 565–591.
- Sun, Ivan Y, Brian K Payne, and Yuning Wu (2008). "The impact of situational factors, officer characteristics, and neighborhood context on police behavior: A multilevel analysis".
 In: Journal of criminal justice 36.1, pp. 22–32.

- Terrill, William (2005). "Police Use of Force: A Transactional Approach". In: Justice Quarterly 22.1, pp. 107–138.
- (2011). "Police Coercion: Application of the Force Continuum". Ph.D. Dissertation.
- VanderWeele, Tyler J and Peng Ding (2017). "Sensitivity analysis in observational research: introducing the E-value". In: Annals of internal medicine.
- Voigt, Rob et al. (2017). "Language from Police Body Camera Footage Shows Racial Disparities in Officer Respect". In: Proceedings of the National Academy of Sciences 114.25, pp. 6521–6526.
- Welch, Kelly (2007). "Black Criminal Stereotypes and Racial Profiling". In: Journal of Contemporary Criminal Justice 23.3, pp. 276–288.
- Wheeler, Andrew P. et al. (2017). "What Factors Influence an Officer's Decision to Shoot? The Promise and Limitations of Using Public Data". In: Justice Research and Policy 18.1, pp. 48–76.
- Wilson, James Q. (1978). Varieties of Police Behavior. Harvard University Press.
- Worden, Robert E. (2015). "The 'Causes' of Police Brutality: Theory and Evidence on Police Use of Force". In: Criminal Justice Theory: Explaining The Nature and Behavior of Criminal Justice 2, pp. 149–204.
- Worrall, John L. et al. (2018). "Exploring Bias in Police Shooting Decisions With Real Shoot/Don't Shoot Cases". In: Crime & Delinquency, pp. 1171–1192.

Zimring, Franklin E. (2017). When Police Kill. Harvard University Press.

A Supplemental Theoretical Results

In any equilibrium in which officers choose to engage in law-enforcement activity—i.e., any equilibrium that reaches the aggressive behavior subgame—there can exist one pure strategy equilibrium to the game, but it requires a particularly strong condition. Specifically, there can exist a pure strategy equilibrium, where the aggressive behavior subgame involves the civilian always choosing to threaten (t = 1) and the officer choosing never to use lethal force (f = 0) if the officer views the cost of killing a civilian is much larger than the cost of losing his own life—formally, if $\delta(1)k_{\rho} - d_O > 1$. In other words, the officer must regard the differential between the value of his own life and the life of the civilian as being greater than the value of stopping crime. That is, the officer would never be willing to use force to stop a violent criminal.

Lemma 1. Any pure strategy equilibrium to the aggressive behavior subgame involves the civilian always threatening (t = 1) and the officer never using lethal force (f = 0). This equilibrium can only hold for $\delta(1)k_{\rho} - d_O > 1$; that is, when the officer regards the difference between value of the civilian's life and his own to be greater the the value of stopping a violent criminal.

Lemma 1 shows that any pure strategy equilibrium is substantively uninteresting. It can only occur under conditions where an officer is completely unwilling to use lethal force. In addition, pure strategy equilibria are not substantively interesting insofar as we observe variation in civilian and police officer behavior, conditional on both observable characteristics and racial categories. In the observed world, officers and civilians appear to be playing mixed strategies. We, therefore, focus the remainder of our analysis on characterizing a mixed strategy equilibrium. We assume, then, that the officer is always willing to use lethal force, if necessary (Alpert and Dunham, 2004).

Assumption 2 (Officers willing to use force). The officer is always weakly willing to use lethal force, $\delta(1)k_{\rho} - d_O \leq 1$, $\forall \rho$

In addition to our theoretical model's implications for observable implications of racial

bias, it also yields insights about the extent to which racial bias might affect what we can learn from studying police-civilian contact altogether. Starting with the officer's decision to engage in law-enforcement activity, the model reveals a number of factors that are important. First, as we described, in equilibrium we will only observe interactions between individuals whose behavior is sufficiently costly to ignore and the police. Specifically, it must be the case that $c_O(\tau)$ —the cost of overlooking potentially criminal activity behavior—is sufficiently large in order to observe law-enforcement activity. While we allow this parameter to vary by the civilian type, we do not assume that observable characteristics and race are orthogonal. That means that evidence of racial bias from the rate of police engagement with a population must take the form of racial disparities conditional on all observable character*istics.* Moreover, because once an officer decides to engage in law-enforcement activity, the probability the officer escalates (i.e., uses force) will also be driven by civilian characteristics. (Recall, the officer's strategy must keep the civilian indifferent between threatening and not threatening.) Therefore, absolute rates of contact between officers and the use of force across racial categories cannot in and of themselves demonstrate racial bias by the police (see also, Knox, Lowe, and Mummolo, 2020).

To assess the effect of racial bias on observed police-civilian interactions, we therefore need to evaluate how changing the value of k_{ρ} affects each stage of the game. we can rearrange Condition (9) as follows:

$$\pi(\tau)^* \le \frac{w(\tau) + \sigma^* d\delta(0)}{b(\tau)(1 - \sigma^*) - \sigma^* d(\delta(1) - \delta(0))}$$

Because $\frac{\partial \pi(\tau)^*}{\partial k_r ho} > 0$, as the cost of killing a civilian of race ρ increases, the civilian of that race, holding constant his observable characteristics, κ , is more willing to play s = 1. And, as we saw, Condition (8) does not depend on k_{ρ} .

Therefore, as police become increasingly racially biased against civilians of race ρ , we should see the pool of individuals interacting with the police shift. In particular, if $k_W > k_B$ —

that is, if the police are racially biased against civilians of race B, then, conditional on observable characteristics, we should see more civilians of race W being subjected to lawenforcement and ultimately killed by police. The intuition here is that individuals of race B will censor their behavior in anticipation of a lower threshold by police for using force against them.

Proposition 3. Racial bias by police induces selection bias in the observed population of civilian-police interactions. All else equal, civilians of the race against whom police are biased censor their behavior and are less likely to engage in behavior that could trigger law-enforcement activity than are civilians of another race.

Of course, Proposition 3 is an all-else-equal statement. It must be the case not only that κ is held constant, but so, too, are the other parameters, especially $w(\tau)$, $b(\tau)$, and $c(\tau)$. This means that the costs of being stopped by the police, the benefit of resisting police, and the benefits of engaging in potentially suspicious activity must be held constant. There is good reason to believe that these quantities are correlated with race, even holding constant observable characteristics, because of the unique social, political, and economic experiences that people of different races have in the US. The consequence is that even attempting to control for every observable characteristics of civilians who could potentially be subjected to law-enforcement activity will not alleviate selection bias (cf, Knox, Lowe, and Mummolo, 2020).

Now we turn to the first stage of the game to assess the conditions under which the players reach the subgame where they decide whether to threaten and use lethal force. For the officer to choose l = 1, it must be the case that the equilibrium expected utility from reaching the aggressive behavior stage is better than the cost of letting a suspicious civilian or possible criminal go undeterred. Abusing notation, this condition is given by:

$$c_O(\tau) \ge \frac{d_O b(\tau)}{b(\tau) + d(\delta(1) + \delta(0))} \tag{8}$$

This condition can be met under a variety of conditions. Most simply, the officer will be willing to engage in law enforcement activity (i.e., play l = 1) for any arbitrarily large value of $c_O(\tau)$ —that is, if the officer finds it too costly to ignore potentially suspicious activity by a civilian of type τ . Later, we consider the empirical implications of this result, but we note that if officers are more suspicious of individuals of a particular race or with a particular set of observable characteristics, we might expect disproportionate engagement by the police with those types of individuals, even for the same type of behavior. This pattern could give rise to observed disparities in the racial makeup of individuals stopped by police and selection bias in the observed set of civilian-officer interactions (see, for example, Knox, Lowe, and Mummolo, 2020).

Finally, then, we show that the civilian can have a positive expected utility from the subgame, which is sufficient to induce her to be willing to enter into the confrontation in the first instance. This condition is met whenever the following is satisfied:

$$\sigma^* \left[\pi^* \left(\tau \right) \left(\delta(1) - \delta(0) \right) + \delta(0) \right] \le \frac{\pi^* \left(\tau \right) b \left(\tau \right) \left(1 - \sigma^* \right) - w \left(\tau \right)}{d} \tag{9}$$

or, alternatively, whenever the officer is unwilling to engage in law-enforcement activity—i.e., when Condition (8) is not satisfied.

While Condition (9) is algebraically messy, it has an intuitive interpretation. It can be satisfied when either (a) $\delta(1) - \delta(0)$, (b) d, or (c) $w(\tau)$ is small enough, relative to other parameters. Substantively, that means a civilian is willing to engage in potentially suspicious behavior whenever Recall, the behavior need not actually be suspicious, and we may think of this as a minimal condition. If a civilian compares essentially remaining cloistered to living a free life and finds the value of engaging in his life's activities to be sufficiently valuable relative to the inconvenience of potentially being stopped by the police, along with the subsequent potential outcomes, then he will be willing to engage in said activity.

Proposition 4. There exists a unique subgame perfect Nash equilibrium in which a civilian of type $\tau = \langle \kappa, \rho \rangle$ initiates conflict for sufficiently low $w(\tau)$ or large d_C , the officer chooses

to engage the civilian for sufficiently low values to either the officer or the civilian of life, and the civilian and officer probabilistically threaten and use force, with probabilities $\pi(\tau)^*$ and $\sigma(\tau)^*$, respectively.

Proposition 4 provides a number of insights into the nature of officer-involved shootings. First, it suggests there can be disparities in who is involved in officer-involved shootings that may or may not be driven by racial bias. A sufficient condition for a racial disparity is that the distribution of characteristics that make the value of life lower for civilians is higher for individuals of one racial group than another. Similarly, it might be that the distribution of characteristics that make the value of crime high is larger for one group than the other. In other words, Proposition 4 reveals that the source of racial disparities in the prevalence of officer-involved shootings (fatal and non-fatal) cannot be ascertained without identifying the distribution of characteristics in the population associated with how individuals value life and crime. However, as we show in the next section, there are implications that follow from this model that allow us to assess racial bias without having to measure such concepts.

At the same time, Proposition 2 reveals a related, but distinct, empirical implication. Note that civilian behavior is designed to maintain indifference by the officer. Therefore, civilians of a race for whom officers are biased in favor of using lethal force should be less likely to threaten officers, *ceteris paribus*. That is, for the same reason we expect Black civilians to be more likely to survive a police officer's use of force, we expect, conditional on being subjected to law-enforcement activity, White civilians will be more likely to engage in threatening behavior, such as resisting arrest, disobeying officer commands, or behaving belligerently. There is no definitive evidence, however, from empirical studies of suspect resistance to support the expectation. While some studies find that Whites are more likely to escalate their behavior during encounters with the police, some studies suggest civilian of color are more likely to escalate towards harm.¹²

¹²Empirical studies directly assessing whether race of civilian is associated with civilian non-compliance and resistance during encounters with police are few. Their conclusions, derived from a variety of sources, including post-encounter narratives of police, use of force case reports, and surveys of victimized and nonvictimized police officers, are mixed. Some studies observe that Whites may be quicker than non-Whites to display resistance during encounters with police (Kahn et al., 2017). Others observe no differences by race

To see this, note that

$$\frac{\partial \pi^*\left(\tau\right)}{\partial k_{\rho}} = \frac{\delta(0)}{1 + d_O + w_O - \left(\delta(1) - \delta(0)\right)k_{\rho}} + \frac{\left(\delta(1) - \delta(0)\right)\left(w_O + \delta(0)k_{\rho}\right)}{\left(1 + d_O + w_O - \left(\delta(1) - \delta(0)\right)\right)^2} > 0$$

That is, as the cost of taking a civilian's life increases, so too does the probability that a civilian of that race threatens the officer. Therefore, given Definition 1, if an officer is racially biased towards killing a civilian of race ρ (i.e., $k_{\rho} < k_{\neg\rho}$), then that civilian should be less likely to engage in behaviors during the encounter that threaten the officer or elevate the use of force by the officer. Conversely, if an officer is racially biased against killing civilian of some race, then that civilian should be more likely to engage in behaviors during the encounter that threaten the officer.

B Proofs of Formal Results

Proof of Lemma 1. The proof proceeds in two steps. First, we show that no pure strategy pair other than $\langle t = 1, f = 0 \rangle$ can be an equilibrium. Second, we show that $\langle t = 1, f = 0 \rangle$ can be an equilibrium on for $k_{\rho} - d_O > 1$.

Consider the strategy pair $\langle t = 1, f = 1 \rangle$. The civilian receives the payoff $-w(\tau) - d(\tau)$. By deviating to t = 0, the civilian receives $-w(\tau)$, which is strictly greater and so would not play $\rho = r$ in equilibrium. Next, consider the strategy pair $\langle r = 0, f = 1 \rangle$. The officer receives a payoff $-w_O - c_O$. By deviating to f = 0, the officer receives 0, which is strictly greater, and so this strategy pair cannot be an equilibrium. Now, consider the pair $\langle t = 0, f = 0 \rangle$. The civilian receives $-w(\tau)$. By deviating to t = 1, she receives $b(\tau) - w(\tau)$, which is strictly greater, and therefore the strategy pair cannot be an equilibrium.

Next, consider the strategy pair $\langle t = 1, f = 0 \rangle$. The officer receives utility $-d_O$, and the

in civilian resistance to the police (e.g., Bierie, Detar, and Craun, 2016). The remainder claim race-based differences, with Blacks being more likely to be (or be perceived) as resistant during police-civilian encounters (Belvedere, John L Worrall, and Tibbetts, 2005; Bierie, 2017). In sum, there is no scholarly consensus about race as an explanation for civilian aggression, especially resistance, during police encounters.

civilian receives utility $b(\tau) - w(\tau)$. By deviating to t = 0, the civilian receives $-w(\tau)$ and so has no incentive to deviate. By deviating to f = 0, the officer receives $1 - k_{\rho}$. Thus, he has an incentive to deviate if only if $k_{\rho} - d_O > 1$.

Proof of Proposition 1. In order for the civilian and the officer to play mixed strategies in a subgame perfect Nash equilibrium, each player's strategy must make the other indifferent between the elements of her choice set. This means the civilian's probability distribution over t must satisfy:

$$EU_O(f = 1|\pi^*, \tau) = EU_O(f = 0|\pi^*, \tau)$$
$$\pi^*(\tau)(1 - \delta(1)k_\rho) - (1 - \pi^*(\tau))(w_O + \delta(0)k_\rho) = -\pi^*(\tau) d_O$$
$$\pi^*(\tau) = \frac{w_O + \delta(0)k_\rho}{1 + w_O + d_O - (\delta(1) - \delta(0))k_\rho}$$

Similarly, the officer's equilibrium probability distribution over f must make the civilian indifferent over the elements of her choice set, t. That is, $\sigma^*(\tau)$ must solve

$$EU_{i}(t = 1|\sigma, \tau) = EU_{i}(t = 0|\sigma, \tau)$$
$$-\sigma^{*}(\tau)(w(\tau) + \delta(1)d(\tau)) + (1 - \sigma^{*}(\tau))(b(\tau) - w(\tau)) = -\sigma^{*}(\tau)(w(\tau) + \delta(0)d(\tau)) - (1 - \sigma^{*}(\tau))w(\tau)$$
$$\sigma^{*}(\tau) = \frac{b(\tau)}{b(\tau) + d(\delta(1) + \delta(0))}$$

Assumption $2 \implies \pi^*(\tau) \in (0, 1)$, and Assumption $1 \implies \sigma^*(\tau) \in (0, 1)$. Therefore, for all parameter values, the players can make each other indifferent and can therefore play mixed strategies in equilibrium.

Proof of Proposition 2. Fatality rates are the proportion of civilians who die among those for whom O chooses to play f = 1; therefore, it is directly proportional to $\pi^*(\tau)$. In order for there to be different fatality proportions among racial groups, $\pi^*(\tau)$ must vary by ρ . The only parameter in $\pi^*(\tau)$ that is a function of ρ is k_{ρ} . Notice that by Definition 1, if an officer is not racially biased, then $k_B = k = k_W$. Given the equilibrium probability of choosing to to threaten is $\pi(\tau)^* = \frac{w_O + \delta(0)k_{\rho}}{1 + w_O + d_O - (\delta(1) - \delta(0))k_{\rho}}$, from above, then we can substitute $\frac{w_O + \delta(0)k_{\rho}}{1 + w_O + d_O - (\delta(1) - \delta(0))k_{\rho}}$ for $\pi^*(\tau)$ in Equation (3) and re-arrange as follows:

$$\mathcal{F}(\rho) = \int_{K(\rho)} \left(\delta(1)\pi^{*}(\tau) + \delta(0)(1 - \pi^{*}(\tau)) \right) \frac{\sigma^{*}(\tau) g(\kappa|\rho)}{\int_{K(\rho)} \sigma^{*}(z|\rho)g(z|\rho)dz} d\kappa$$
$$\mathcal{F}(\rho) = \delta(0) + \left(\delta(1) - \delta(0)\right) \left(\frac{w_{O} + \delta(0)k_{\rho}}{1 + w_{O} + d_{O} - \left(\delta(1) - \delta(0)\right)k_{\rho}} \right)$$

This implies that

$$\mathcal{F}(B) = \delta(0) + (\delta(1) - \delta(0)) \left(\frac{w_O + \delta(0)k_\rho}{1 + w_O + d_O - (\delta(1) - \delta(0))k_\rho} \right) = \mathcal{F}(W)$$
(10)

Therefore, differential fatality rates can only arise if $k_{\rho} \neq k_{\neg \rho}$, which, by Definition 1 means O is racially biased.

Proof of Proposition 3. Notice that Condition (9) characterizes which types of civilians are willing to engage in potentially suspicious behavior and therefore be candidates for law-enforcement activity. That condition can be re-written as

$$d\left[\pi^{*}(\tau) \left(\delta(1) - \delta(0)\right) + \delta(0)\right] \leq \frac{\pi^{*}(\tau) b(\tau) \left(1 - \sigma^{*}\right) - w(\tau)}{\sigma^{*}}$$

Note that $\frac{\partial \sigma^*}{\partial k_{\rho}} > 0$. Therefore, all else equal, as k_{ρ} decreases—as the officer becomes increasingly biased towards using lethal force against a civilian of race ρ , then this inequality is less likely to hold.

Proof of Proposition 4. Proposition 1 demonstrates that in the subgame where the officer has decided to engage in law-enforcement activity (i.e., l = 1), there exists a mixed strategy subgame perfect Nash equilibrium. To complete the proof, we must show that the subgame can be reached in equilibrium and that the mixed strategies characterized by Proposition 1 constitute the unique equilibrium.

In order to reach the subgame, the officer must be willing to engage in law-enforcement

activity, and the civilian must be willing to engage in potentially suspicious behavior. A sufficient condition for the civilian to play s = 0 is $EU_i [s = 1, \pi | \tau] \ge EU_i [s = 0, \pi | \tau]$ and a sufficient condition for the officer to play l = 1 is $EU_O [l = 1, \sigma | \tau, s] \ge EU_O [l = 0, \sigma | \tau, s]$ Notice that if s = 0, the game ends. Therefore, we it is sufficient to show

$$EU_{O} [l = 1 | \tau, s = 1] \ge EU_{O} [l = 0 | \tau, s = 1]$$
$$-\sigma^{*} (\pi^{*} (1 + \delta(1)k_{\rho}) + (1 - \pi^{*}) (w_{O} + \delta(0)k_{\rho})) - (1 - \sigma^{*}) d_{O} \ge -c_{O} (\tau)$$
$$\frac{d_{O}b(\tau)}{b(\tau) + d(\delta(1) + \delta(0))} \le c_{O} (\tau)$$

which can be true for an arbitrarily large value of $c_O(\tau)$. Finally, we must show that, given these constraints, the civilian is willing to engage in potentially suspicious behavior. Formally, it must be the case that

$$EU_{i} [s = 1|\tau] \ge EU_{i} [s = 0|\tau]$$
$$-\sigma^{*} [\pi^{*} (w(\tau) + \delta(1)d(\tau)) - (1 - \pi^{*}) (w(\tau) + \delta(0)d)] + (1 - \sigma^{*}) [\pi^{*} (b(\tau) - w(\tau)) - (1 - \pi^{*}) w(\tau)] \ge 0$$
$$w(\tau) \le d \left[\frac{\pi^{*} (\tau) b(\tau) (1 - \sigma^{*})}{d} - \sigma^{*} [\pi^{*} (\tau) (\delta(1) - \delta(0)) + \delta(0)] \right]$$

which can be true for an arbitrarily small value of $w(\tau)$. Now to see that the equilibrium is unique, notice that Assumption 2 and Lemma 1 imply the mixed strategies $\pi^*(\tau)$ and $\sigma^*(\tau)$ are the unique equilibrium strategies in the aggressive behavior subgame. Further, notice the earlier stage of the game involves perfect and complete information, and the players cannot be indifferent between their strategies choices. Therefore, the pure strategies given by Conditions (8) and (9) characterize the unique equilibrium.

C Data Cleaning and Organization

The data the law enforcement agencies provided contained neither unique incident ID numbers nor the total number of civilian or officers involved in any incident. Therefore, to construct unique civilian/officer pairs we first made the assumption that any observations from the same city, location and date comprised the same incident. Table 3a gives an example of data we received from Tucson. Given that these two observations are both from Tucson and are on the same date of "02/07/2013" with the same location of "925 E Mill St." we assume they are the same incident. Second, the data contained neither civilian nor officer identifying information other than race. Therefore, we must assume each row of data represents unique people. The incident in Table 3a would then be considered to have three total people involved: one White civilian, one Hispanic civilian and one White officer. Third, we made as many civilian/officer pairs for each incident as there were possible combinations of people. Table 3b shows how we completed the civilian/officer pairs. We now have, for example, the White civilian matched to the White officer and the Hispanic civilian matched to the White officer,¹³

Obs.	Date	Location	City	Civilian Race	Officer Race
1	02/07/2013	925 E Mill St.	Tucson	White	White
2	02/07/2013	$925 \to Mill St.$	Tucson	Hispanic	NA
(a) Example of raw data received from FOIA requests					
	(a) Lx	ample of raw data	i iteeiveu	ITOIL FOIA TEQUES	G
Obs.		Location	City	Civilian Race	Officer Race
Obs. 1	Date 02/07/2013	Location 925 E Mill St.	City Tucson	Civilian Race White	Officer Race White

Table 3: Data Example

(b) Example of complete civilian/officer pairs

¹³The modal incident only has one officer-civilian pair. However, we have also estimated our model clustering observations by incident to account for within-incident correlation between observations, and our results do not change appreciably.

D Lower Bound Derivation for Magnitude of Racial Bias

Following the derivations in Cohen and Glynn (2021), given two groups, Black civilians and White civilians, where $\rho_i \in \{b, w\}$ indicates civilian *i*'s race, and the decision to shoot civilian *i* is defined as $S_i \in \{0, 1\}$. If $Y \in \{0, 1\}$ is the outcome, we can define the probability of police fatally shooting a civilian of race ρ as follows:

Definition 2. The rate of of being fatally shot among a racial group, ρ , is given by $\mathcal{F}_{\rho} = E[Y|S(\rho) = 1, \rho].$

Per our model, let $\rho = b$ indicate a Black civilian. The fatality rate among Black civilians who were shot, \mathcal{F}_b , can be expressed as a combination of Black civilians who would have been shot had they been White $(\mathcal{F}_{s(b)=s(w),b})$, and Black civilians who would *not* have been shot had they been White $(\mathcal{F}_{s(b)>s(w),b})$. The magnitude of racial bias in the decision to shoot a civilian is is the proportion of Black civilians shot who would not have been shot had they been White, $p = Pr[S(w) = 0, S(b) = 1|\rho = b]$. This can be formally written as:

$$\mathcal{F}_b = p \cdot \mathcal{F}_{s(b) > s(w), b} + (1 - p) \cdot \mathcal{F}_{s(b) = s(w), b}$$
(11)

and can be rearranged to solve for p:

$$p = \frac{\mathcal{F}_b - \mathcal{F}_{s(b)=s(w),b}}{\mathcal{F}_{s(b)>s(w),b} - \mathcal{F}_{s(b)=s(w),b}}.$$

If we further assume that the fatality rate among Whites is at least the fatality rate among Blacks that would have been shot had they been White (i.e., $\mathcal{F}_w \geq \mathcal{F}_{s(b)=s(w),b}$), then we can substitute \mathcal{F}_w for $\mathcal{F}_{s(b)=s(w),b}$ —the fatality rate among all White civilians shot can stand in for the fatality rate for Black civilians who were shot and who would have been shot were they White. Note that the logic of our formal analysis and our empirical findings both indicate $\mathcal{F}_w > \mathcal{F}_b$. Therefore, we can rewrite this expression of p:

$$p = \frac{\mathcal{F}_w - \mathcal{F}_b}{\mathcal{F}_w - \mathcal{F}_{s(b) > s(w), b}}.$$
(12)

Our p is expressed as the difference between the observed rate of being fatally shot among White and Black civilians, divided by the difference between the observed rate of being fatally shot among White civilians and the rate of being fatally shot among Black civilians who would not have been shot were they White, $\mathcal{F}_{s(b)>s(w),b}$. That quantity is not observed because we do not know precisely who would not have been fatally shot had they been White. However, if there is no racial bias, we know that $\mathcal{F}_{s(b)>s(w)}$ would be 0, as there would be no civilians shot because they were Black. Therefore, substituting 0 for $\mathcal{F}_{s(b)>s(w)}$ yields a lower bound on the true value of p:

$$p \ge \frac{\mathcal{F}_w - \mathcal{F}_b}{\mathcal{F}_w - 0} = \frac{\mathcal{F}_w - \mathcal{F}_b}{\mathcal{F}_w}.$$
(13)

D.1 Estimating the lower bound

Following the derivations in Cohen (2021), to estimate the lower bound on the proportion of Black civilians that would not have been shot had they been White, we begin by estimating a Poisson model:

$$P[\hat{Y} = 1|X] = exp^{\hat{\alpha} + \hat{\beta}_w D + \mathbf{X}\hat{\gamma} + \hat{\epsilon}}$$
(14)

In equation 14, $\hat{\beta}_w$ is the coefficient of interest on the binary treatment variable D. We estimate it with a subset of officer-involved shooting data containing only Black and White civilians, where our main covariate of interest is an indicator for race (White equal to 1, Black equal to 0). **X** is a vector of covariates which include city fixed effects and distance in miles to closest trauma center. From this we can estimate the risk ratio of White compared to Black as

$$\widehat{RR}_{DY|\mathbf{X}} = \frac{E[P(\hat{Y} = 1|D = 1, \mathbf{X} = x)]}{E[P(\hat{Y} = 1|D = 0, \mathbf{X} = x)]}$$
(15)

Using this model we estimate the regression coefficient on White to be 0.37 (see Appendix Table 11 for the full regression specification results) making the relative risk 1.45. White civilians have 45% increase in the likelihood of fatality compared to Black civilians controlling for city fixed effects and distance to closest trauma center. Using the relative risk we can estimate p (equation 16, see Cohen (2021) for the full derivation) the lower bound on the proportion of Black civilians that would not have been shot had they been white.

$$p \ge 1 - \frac{1}{\hat{RR}} \tag{16}$$

With a relative risk of 1.45 we estimate a lower bound of 0.31. We are using a Poisson model because we can directly and easily estimate the relative risk but as additional robustness checks we also estimate a logistic model and a linear probability model (LPM). The coefficient on White from the logistic model is 0.70 (see Table 13). This gives a relative risk of 1.46 and a lower bound estimate of 0.32.

The results from the LPM can be used to estimate relative risks as with the Poisson and logistic or the coefficient on White can be used to directly estimate the lower bound (equation 17). The coefficient estimate on White is the numerator and the coefficient plus the observed mean fatality rate among Black civilians is the denominator. In the linear probability model we estimate $\hat{\beta}_w = 0.16$ and a lower bound of at least 0.32.

$$p \ge \frac{\hat{\beta}_w}{F_b + \hat{\beta}_w} \tag{17}$$

Substantively, our estimate of 31% is considerable and given it is a lower bound, may be much higher. Our estimate implies that police would not have shot 148 Black civilians had they been White, from the 477 Black civilians in our nine localities over the years we study. Extrapolating this estimate to the larger population of the United States, however, is beyond the limits of our data. Moreover, significant intra-locality variation suggests police behavior, measured by officer-involved shootings, is not uniform across the country. Additionally, comparing Hispanic civilians and Asian civilians to White civilians yielded no statistically significant differences. That is consistent with what we would expect—police officers differentially exercise discretion against Black civilians as compared to all other groups. Given the extant debate about whether the use of force by police is tainted with racial bias, these findings suggest there is a substantively significant problem. Quantifying the magnitude of its effect, though, requires richer administrative data beyond what police departments, generally, in the U.S. currently provide. Specifically, the important matter of how much police violence is attributable to racial bias requires knowing how often police fire their weapons, as well as how often they draw their weapons (John L. Worrall et al., 2018; Wheeler et al., 2017), which is not universally known across local police departments.

E Assessing the Mechanism

We argue the mechanism at work in the empirical evidence we have shown is that raciallybiased officers have a lower threshold for using force against racial minorities than against White civilians. The evidence we have shown is consistent with the consequences of such bias. We now step back to assess broader evidence, outside the context of officer-involved shootings, to corroborate our claim about the underlying mechanism. In particular, we consider whether we do in fact observe lower thresholds for using force by when officers encounter Black civilians, as compared to White civilians.

To do so, we marshal several datasets on officer-civilian interactions and show a consistent pattern across a varied set of local police jurisdictions. We note at the outset, these data comprise documented incidents of officer-civilian interactions from the perspective of the police and they may suffer from problems of selection bias and unobservable counterfactuals. However, our goal here is to demonstrate there are patterns beyond those we have documented that are consistent with the claim that officers have a lower threshold for using force against Black civilians than against White civilians. To the extent we find evidence consistent with that claim, we can be more confident that the patterns in civilian fatalities in officer-involved shootings are caused by the theoretical model we have proposed, as opposed to some other process.

Jurisdiction	Years covered	Brief description
New York City	2006-2015	Stops with indicators for different kinds of force
Washington, DC	4 weeks in 2019	Data on all stops, outcome is whether a pat-down was conducted

Table 4: Summary of civilian contact data used to assess racial disparities in the use of force.

Table 4 summarizes the police data we have assembled. From New York City, we have the widely studied Stop, Question, and Frisk data, which comprise 11 years of data on incidents in which officers stop civilians and contain detailed information about actions taken by the officer during the encounter. One very widely studied source of variation in these data are indicators for different kinds of force that an officer may have used during a stop. From Washington, DC, we have a recently-released dataset that comprises just four weeks of stops during 2019 but include an indicator for whether an officer conducted a pat-down of the civilian involved in the stop.

Table 5 reports the results of a series of fixed effects linear regression models in which the dependent variables are indicators of force or the conducting of a pat-down. The explanatory variables are fixed effects for civilian race as well as other fixed effects, depending on what is available and feasible from each city. For example, we have more than 5,000,000 observations from New York, and so we include fixed effects for every month-year pair during our window, as well as the precinct in which the stop took place. For Washington, DC, we only have a few weeks' data, and so we include date-specific fixed effects, along with the district in which the stop took place. In each of the models we specify, we use White civilians as the baseline category, so the table entries can be interpreted as differences between the groups in the table and White civilians.

	Washington, DC	NYC
Dependent variable:	Pat-down	Force Used
Dlask Civilian	0.09***	0.04***
DIACK CIVIIIAII	(0.01)	(< 0.01)
Hispanic Civilian	0.01	
mspanic Orvinan	(0.01)	
Black-Hispanic Civilian		0.04^{***}
black mspanic Oreman		(< 0.01)
White-Hispanic Civilian		0.02***
white hispanie ervinan		(< 0.01)
Asian Civilian	0.00	-0.01***
	(0.02)	(< 0.01)
Native American		0.00
		(< 0.01)
Multiple Race Civilian	-0.01	
	(0.05)	
Other Race Civilian	-0.02	0.01***
• • • • • • • • • • • • • • • • • • • •	(0.14)	(< 0.01)
Unknown Civilian Race	-0.01	-0.06^{***}
	(0.02)	(0.01)
N		5029789
Fixed effects	Date, district	Month-year, precinct

Table 5: Racial disparities in the use of force in selected datasets. Entries are linear regression coefficients, standard errors in parentheses. *** $p \leq .001$, ** $p \leq .01$, * $p \leq .05$

Across all of our specifications, Black civilians are more likely to be subjected to force than are White civilians. These data are consistent with the theoretical mechanism underlying our model—that officers might have a lower threshold for using force against a Black civilian than a White civilian. We caution, though, these data are less closely-connected to our model and so should be interpreted with caution. However, we believe they do provide at least preliminary additional evidence of the theoretical mechanism we contemplate.

F Additional Model Specifications

Table 8 repeats the original specification from the main paper with city fixed effects as Model 1. The additional specifications include year fixed effects. Notably, the magnitude of the relationship between being a Black civilian and the probability of dying *increases* once we include jurisdiction fixed-effects, and maintains when we include year fixed-effects.

	Model 1	Model 2	Model 3	Model 4
Black	-0.67^{*}	-0.82^{***}	-0.74^{**}	-0.77^{***}
	(0.26)	(0.25)	(0.25)	(0.23)
Hispanic	0.13	0.28	0.11	0.12
	(0.16)	(0.18)	(0.16)	(0.13)
Asian/AI/AN/PI	0.91^{*}	1.04^{**}	0.96^{*}	0.98^{*}
	(0.39)	(0.36)	(0.39)	(0.50)
Closest Trauma (10s miles)	0.41^{***}	0.19	0.45^{***}	0.50^{***}
	(0.06)	(0.19)	(0.07)	(0.09)
Houston	-0.21^{***}		-0.19^{*}	15.17^{***}
	(0.05)		(0.09)	(0.13)
King County	-0.12		-0.17	14.02^{***}
	(0.14)		(0.18)	(0.28)
Los Angeles	1.15^{***}		1.18^{***}	16.96^{***}
	(0.06)		(0.09)	(0.14)
Orlando	0.52^{***}		0.50^{***}	15.56^{***}
	(0.04)		(0.04)	(0.18)
San Jose	0.16		0.18	15.38^{***}
	(0.09)		(0.16)	(0.08)
Seattle	1.27^{***}		1.20^{***}	16.66^{***}
	(0.05)		(0.06)	(0.24)
Tucson	1.64^{***}		1.54^{***}	32.32^{***}
	(0.07)		(0.08)	(0.49)
2011		0.51^{*}	0.33^{*}	0.88^{***}
		(0.20)	(0.16)	(0.23)
2012		0.25	0.33	33.51^{***}
		(0.29)	(0.25)	(0.16)
2013		0.34	0.30	16.43^{***}
		(0.45)	(0.45)	(0.20)
2014		0.63^{**}	0.63^{*}	15.24^{***}
		(0.19)	(0.25)	(0.22)
2015		0.20	0.27	17.38^{***}
		(0.45)	(0.48)	(0.22)
2016		0.13	0.28	1.29^{***}
		(0.23)	(0.21)	(0.37)
2017		-0.13	-0.09	1.29^{***}
		(0.21)	(0.22)	(0.03)
Intercept	-1.00^{***}	-0.31	-1.27^{***}	-16.83^{***}
	(0.17)	(0.41)	(0.23)	(0.05)
City*Year				\checkmark
Num. obs.	1250	1250	1250	1250

 $^{***}p < 0.001; \ ^{**}p < 0.01; \ ^*p < 0.05$

Table 6: Estimated relationship between civilian race and probability of fatality conditional upon being involved in an officer-involved shooting. Cells show logistic coefficients with clustered-robust standard errors. Omitted category is white civilians, Charlotte and the year 2010.

N	fodel 1	Model 2 $$	Model 3	Model 4
Black -	-0.15^{*}	-0.19^{*}	-0.17^{*}	-0.16^{*}
	(0.06)	(0.06)	(0.06)	(0.05)
Hispanic	0.03	0.07	0.03	0.03
_	(0.04)	(0.04)	(0.04)	(0.03)
Asian/AI/AN/PI	0.19^{*}	0.24^{**}	0.20^{*}	0.18
	(0.07)	(0.06)	(0.07)	(0.08)
Closest Trauma (10s miles) ().09***	0.05	0.10^{***}	0.10^{***}
	(0.01)	(0.04)	(0.01)	(0.02)
Houston –	-0.05**		-0.04^{*}	0.22^{***}
	(0.01)		(0.01)	(0.03)
King County	-0.03		-0.04	-0.05
	(0.03)		(0.03)	(0.05)
Los Angeles ().25***		0.26^{***}	0.59^{***}
	(0.02)		(0.02)	(0.03)
Orlando (0.10^{***}		0.10^{***}	0.30^{***}
	(0.01)		(0.01)	(0.04)
San Jose	0.02		0.03	0.23^{***}
	(0.02)		(0.03)	(0.02)
Seattle).28***		0.27^{***}	0.55^{***}
	(0.02)		(0.02)	(0.04)
Tucson).37***		0.35^{***}	0.85^{***}
	(0.02)		(0.02)	(0.08)
2011		0.12^{*}	0.07	0.18^{**}
		(0.05)	(0.04)	(0.05)
2012		0.06	0.07	1.07^{***}
		(0.07)	(0.05)	(0.02)
2013		0.08	0.07	0.45^{***}
		(0.10)	(0.09)	(0.04)
2014		0.15^{*}	0.14	0.27^{***}
		(0.05)	(0.06)	(0.05)
2015		0.05	0.06	0.67^{***}
		(0.10)	(0.10)	(0.05)
2016		0.03	0.06	0.26^{**}
		(0.05)	(0.04)	(0.07)
2017		-0.03	-0.02	0.29^{***}
		(0.05)	(0.05)	(0.01)
Intercept ().29***	0.43^{**}	0.23^{**}	-0.06^{**}
	(0.04)	(0.10)	(0.05)	(0.01)
City*Year				\checkmark
Num. obs.	1250	1250	1250	1250

***p < 0.001; **p < 0.01; *p < 0.05

Table 7: Estimated relationship between civilian race and probability of fatality conditional upon being involved in an officer-involved shooting. Cells show linear probability model coefficients with clustered-robust standard errors. Omitted category is white civilians, Charlotte and the year 2010.

	Model 1	Model 2	Model 3	Model 4
Black	-0.35^{**}	-0.46^{**}	-0.39^{**}	-0.37^{***}
	(0.14)	(0.15)	(0.13)	(0.11)
Hispanic	0.05	0.12	0.04	0.04
-	(0.07)	(0.08)	(0.06)	(0.04)
Asian/AI/AN/PI	0.30***	0.40***	0.31***	0.27**
	(0.08)	(0.07)	(0.07)	(0.10)
Closest Trauma (10s miles)	0.18***	0.09	0.19***	0.21***
	(0.02)	(0.08)	(0.03)	(0.04)
Houston	-0.06^{**}	. ,	-0.04	14.75***
	(0.02)		(0.04)	(0.06)
King County	0.10		0.08	14.31***
	(0.06)		(0.08)	(0.06)
Los Angeles	0.72^{***}		0.73***	15.78***
	(0.04)		(0.05)	(0.06)
Orlando	0.33^{***}		0.32^{***}	14.98^{***}
	(0.02)		(0.02)	(0.17)
San Jose	0.21^{***}		0.22^{**}	14.89^{***}
	(0.05)		(0.08)	(0.03)
Seattle	0.78^{***}		0.75^{***}	15.72^{***}
	(0.03)		(0.04)	(0.06)
Tucson	0.91^{***}		0.86^{***}	16.10^{***}
	(0.04)		(0.04)	(0.10)
2011		0.24	0.15	0.42^{***}
		(0.13)	(0.08)	(0.10)
2012		0.12	0.14	16.43^{***}
		(0.15)	(0.13)	(0.04)
2013		0.17	0.13	15.44^{***}
		(0.23)	(0.22)	(0.09)
2014		0.31^{**}	0.28^{**}	14.47^{***}
		(0.10)	(0.10)	(0.10)
2015		0.10	0.12	15.99^{***}
		(0.23)	(0.23)	(0.10)
2016		0.06	0.13	0.49^{***}
		(0.11)	(0.11)	(0.09)
2017		-0.06	-0.05	0.77^{***}
		(0.10)	(0.10)	(0.01)
Intercept	-1.34^{***}	-0.85^{***}	-1.46^{***}	-16.42^{***}
	(0.08)	(0.22)	(0.12)	(0.02)
City*Year				\checkmark
Num. obs.	1250	1250	1250	1250

 $^{***}p < 0.001; \ ^{**}p < 0.01; \ ^*p < 0.05$

Table 8: Estimated relationship between civilian race and probability of fatality conditional upon being involved in an officer-involved shooting. Cells show Poisson coefficients with clustered-robust standard errors. Omitted category is white civilians, Charlotte and the year 2010.

Table 7 includes linear probability models. In each of our specifications, among those civilians shot by an officer, Black civilians are less likely to die than are White civilians. This difference is statistically significant in each specification. Notably, as with the logistic model discussed in the paper, the magnitude of the relationship between being a Black civilian and the probability of dying *increases* once we include jurisdiction fixed-effects, and maintains when we include year fixed-effects. What is more, the relationship between being a Hispanic civilian and a reduced probability of dying does not emerge even after we include jurisdiction-level and year fixed effects. This functions as a placebo test and implies that any problematic unmeasured covariates would have to have different relationships for Black and Hispanic civilians (e.g., concerns about characteristics that affect the probability of death—such as police behavior, training, and medical attention would be largely ruled out by this analysis).

G Sensitivity Analysis

We conduct the sensitivity analysis using this specification from table 9 to estimate the relative risks (Cohen, 2021). We use an observed covariate as a comparison (Distance to closest trauma center, Los Angeles and Seattle) and estimate risk ratios of the associated effect of the covariate with the outcome (RR_{UY}) and with treatment (RR_{DU}) . The top two rows of Table 10 show these relative risks for all three comparison covariates. From these risk ratios we calculate a bias factor (VanderWeele and Ding, 2017), adjust the main risk ratio on White and estimate a lower bound. This lower bound is now an an estimate of the proportion of Black civilians who would not have been shot had they been white, adjusted for a possible unobserved confounder of the same strength of association and as the benchmark covariate.

	Model 1
White	0.37^{*}
	(0.15)
Houston	-0.04
	(0.02)
King County	0.10
	(0.09)
Los Angeles	0.71^{***}
	(0.01)
Orlando	0.31^{***}
	(0.02)
San Jose	0.27^{***}
	(0.05)
Seattle	0.72^{***}
	(0.03)
Tucson	0.65^{***}
	(0.05)
Closest Trauma (10s miles)	0.16^{***}
	(0.05)
Intercept	-1.65^{***}
	(0.07)
Num. obs.	715

***p < 0.001; **p < 0.01; *p < 0.05

Table 9: Poisson model regression results. Each model is run on a subset of the data with only Black and white civilians. Clustered-robust standard errors are shown.

	Distance (10s miles)	Los Angeles	Seattle
RR_{UY}	1.17	2.03	$2.05 \\ 1.33$
RR_{DU}	1.105^{*}	1.01	
BF	1.014	1.005	1.15
BR/BF	1.43	1.45	1.27
p	0.303	0.31	0.21

Table 10: Benchmark adjustment to relative risk

	Black/White	White/Hispanic	White/Asian
White	0.37^{*}	-0.08^{*}	-0.26^{*}
	(0.15)	(0.04)	(0.11)
Houston	-0.04	0.33***	0.11^{***}
	(0.02)	(0.02)	(0.03)
King County	0.10	0.72^{***}	0.63^{**}
	(0.09)	(0.06)	(0.21)
Los Angeles	0.71^{***}	1.12^{***}	1.04^{***}
	(0.01)	(0.03)	(0.04)
Orlando	0.31^{***}	0.14^{***}	0.04
	(0.02)	(0.01)	(0.02)
San Jose	0.27^{***}	0.73^{***}	0.73^{***}
	(0.05)	(0.02)	(0.02)
Seattle	0.72^{***}	1.44^{***}	1.35^{***}
	(0.03)	(0.01)	(0.04)
Tucson	0.65^{***}	1.50^{***}	1.34^{***}
	(0.05)	(0.02)	(0.05)
Closest Trauma (10s miles)	0.16^{***}	0.27^{***}	0.28
	(0.05)	(0.05)	(0.18)
Intercept	-1.65^{***}	-1.79^{***}	-1.51^{***}
	(0.07)	(0.05)	(0.14)
Num. obs.	715	732	279

 $^{***}p < 0.001; \ ^{**}p < 0.01; \ ^*p < 0.05$

Table 11: Poisson model regression results. Each model is run on a subset of the data to make either a Black/White, Hispanic/White or Asian/White comparison. Cluster-robust standard errors clustered on city. See Zou (2004) for a discussion of the robust standard error correction when Poisson models are used to calculate relative risks.

	Black/White	White/Hispanic	White/Asian
White	0.16^{*}	-0.04	-0.16
	(0.07)	(0.03)	(0.08)
Houston	-0.02	0.04^{*}	0.01
	(0.01)	(0.02)	(0.01)
King County	-0.01	0.17^{***}	0.16
	(0.03)	(0.02)	(0.09)
Los Angeles	0.24^{***}	0.39***	0.37^{***}
	(0.01)	(0.02)	(0.02)
Orlando	0.09***	0.05^{***}	0.02
	(0.01)	(0.00)	(0.01)
San Jose	0.06^{*}	0.20***	0.21^{***}
	(0.02)	(0.02)	(0.01)
Seattle	0.25^{***}	0.57^{***}	0.58^{***}
	(0.01)	(0.00)	(0.02)
Tucson	0.22^{***}	0.63***	0.55^{***}
	(0.02)	(0.01)	(0.02)
Closest Trauma (10s miles)	0.07^{*}	0.15^{***}	0.14
	(0.03)	(0.02)	(0.09)
Intercept	0.15^{***}	0.13^{**}	0.28^{*}
	(0.02)	(0.03)	(0.09)
RMSE	0.47	0.47	0.45
Num. obs.	715	732	279

***p < 0.001; **p < 0.01; *p < 0.05

Table 12: Linear probability model regression results. Each model is run on a subset of the data to make either a Black/White, Hispanic/White or Asian/White comparison. Cluster-robust standard errors clustered on city.

	Black/White	White/Hispanic	White/Asian
White	0.70^{*}	-0.20	-0.80
	(0.28)	(0.12)	(0.43)
Houston	-0.08	0.28***	0.05
	(0.06)	(0.08)	(0.10)
King County	0.00	0.83^{***}	0.70
	(0.17)	(0.11)	(0.46)
Los Angeles	1.12^{***}	1.78^{***}	1.66^{***}
	(0.03)	(0.08)	(0.07)
Orlando	0.47^{***}	0.23***	0.08
	(0.04)	(0.02)	(0.05)
San Jose	0.33^{***}	1.00^{***}	1.04^{***}
	(0.09)	(0.07)	(0.05)
Seattle	1.15^{***}	2.62^{***}	2.66^{***}
	(0.03)	(0.02)	(0.11)
Tucson	1.02^{***}	2.96^{***}	2.50^{***}
	(0.07)	(0.06)	(0.12)
Closest Trauma (10s miles)	0.31^{*}	0.70^{***}	0.66
	(0.13)	(0.11)	(0.46)
Intercept	-1.57^{***}	-1.73^{***}	-0.97^{*}
	(0.13)	(0.12)	(0.47)
Num. obs.	715	732	279
Deviance	896.32	911.29	329.27
Log Likelihood	-448.16	-455.65	-164.63

***p < 0.001; ** p < 0.01; *p < 0.05

Table 13: Logistic model regression results. Each model is run on a subset of the data to make either a Black/White, Hispanic/White or Asian/White comparison. Cluster-robust standard errors clustered on city.

References

- Alpert, Geoffrey P. and Roger G. Dunham (2004). Understanding Police Use of Force: Officers, Suspects, and Reciprocity. Cambridge University Press.
- Belvedere, Kimberly, John L Worrall, and Stephen G Tibbetts (2005). "Explaining suspect resistance in police-citizen encounters". In: *Criminal Justice Review* 30.1, pp. 30–44.
- Bierie, David M (2017). "Assault of police". In: Crime & Delinquency 63.8, pp. 899–925.
- Bierie, David M, Paul J Detar, and Sarah W Craun (2016). "Firearm violence directed at police". In: Crime & Delinquency 62.4, pp. 501–524.
- Cohen, Elisha (2021). "Sensitivity Analysis for Outcome Tests with Binary Data".
- Cohen, Elisha and Adam Glynn (2021). "Estimating Bounds on Selection Bias with Outcome Tests".
- Kahn, Kimberly Barsamian et al. (2017). "How suspect race affects police use of force in an interaction over time." In: *Law and human behavior* 41.2, p. 117.
- Knox, Dean, Will Lowe, and Jonathan Mummolo (2020). "Administrative records mask racially biased policing". In: American Political Science Review 114.3, pp. 619–637.
- VanderWeele, Tyler J and Peng Ding (2017). "Sensitivity analysis in observational research: introducing the E-value". In: Annals of internal medicine.
- Wheeler, Andrew P. et al. (2017). "What Factors Influence an Officer's Decision to Shoot? The Promise and Limitations of Using Public Data". In: Justice Research and Policy 18.1, pp. 48–76.
- Worrall, John L. et al. (2018). "Exploring Bias in Police Shooting Decisions With Real Shoot/Don't Shoot Cases". In: Crime & Delinquency, pp. 1171–1192.
- Zou, Guangyong (2004). "A modified poisson regression approach to prospective studies with binary data". In: *American journal of epidemiology* 159.7, pp. 702–706.