



# Language from police body camera footage shows racial disparities in officer respect

Rob Voigt<sup>a,1</sup>, Nicholas P. Camp<sup>b</sup>, Vinodkumar Prabhakaran<sup>c</sup>, William L. Hamilton<sup>c</sup>, Rebecca C. Hetey<sup>b</sup>, Camilla M. Griffiths<sup>b</sup>, David Jurgens<sup>c</sup>, Dan Jurafsky<sup>a,c</sup>, and Jennifer L. Eberhardt<sup>b,1</sup>

<sup>a</sup>Department of Linguistics, Stanford University, Stanford, CA 94305; <sup>b</sup>Department of Psychology, Stanford University, Stanford, CA 94305; and <sup>c</sup>Department of Computer Science, Stanford University, Stanford, CA 94305

Contributed by Jennifer L. Eberhardt, March 26, 2017 (sent for review February 14, 2017; reviewed by James Pennebaker and Tom Tyler)

**Using footage from body-worn cameras, we analyze the respectfulness of police officer language toward white and black community members during routine traffic stops. We develop computational linguistic methods that extract levels of respect automatically from transcripts, informed by a thin-slicing study of participant ratings of officer utterances. We find that officers speak with consistently less respect toward black versus white community members, even after controlling for the race of the officer, the severity of the infraction, the location of the stop, and the outcome of the stop. Such disparities in common, everyday interactions between police and the communities they serve have important implications for procedural justice and the building of police–community trust.**

racial disparities | natural language processing | procedural justice | traffic stops | policing

Over the last several years, our nation has been rocked by an onslaught of incidents captured on video involving police officers' use of force with black suspects. The images from these cases are disturbing, both exposing and igniting police–community conflict all over the country: in New York, Missouri, Ohio, South Carolina, Maryland, Illinois, Wisconsin, Louisiana, Oklahoma, and North Carolina. These images have renewed conversations about modern-day race relations and have led many to question how far we have come (1). In an effort to increase accountability and transparency, law enforcement agencies are adopting body-worn cameras at an extremely rapid pace (2, 3).

Despite the rapid proliferation of body-worn cameras, no law enforcement agency has systematically analyzed the massive amounts of footage these cameras produce. Instead, the public and agencies alike tend to focus on the fraction of videos involving high-profile incidents, using footage as evidence of innocence or guilt in individual encounters.

Left unexamined are the common, everyday interactions between the police and the communities they serve. By best estimates, more than one quarter of the public (ages 16 y and over) comes into contact with the police during the course of a year, most frequently as the result of a police-initiated traffic stop (4, 5). Here, we examine body-worn camera footage of routine traffic stops in the large, racially diverse city of Oakland, CA.

Routine traffic stops are not only common, they are consequential, each an opportunity to build or erode public trust in the police. Being treated with respect builds trust in the fairness of an officer's behavior, whereas rude or disrespectful treatment can erode trust (6, 7). Moreover, a person's experiences of respect or disrespect in personal interactions with police officers play a central role in their judgments of how procedurally fair the police are as an institution, as well as their willingness to support or cooperate with the police (8, 9).

Blacks report more negative experiences in their interactions with the police than other groups (10). Across numerous studies, for example, blacks report being treated less fairly and respectfully in their contacts with the police than whites (6, 11). Indeed,

some have argued that racial disparities in perceived treatment during routine encounters help fuel the mistrust of police in the controversial officer-involved shootings that have received such great attention. However, do officers treat white community members with a greater degree of respect than they afford to blacks?

We address this question by analyzing officers' language during vehicle stops of white and black community members. Although many factors may shape these interactions, an officer's words are undoubtedly critical: Through them, the officer can communicate respect and understanding of a citizen's perspective, or contempt and disregard for their voice. Furthermore, the language of those in positions of institutional power (police officers, judges, work superiors) has greater influence over the course of the interaction than the language used by those with less power (12–16). Measuring officer language thus provides a quantitative lens on one key aspect of the quality or tone of police–community interactions, and offers new opportunities for advancing police training.

Previous research on police–community interactions has relied on citizens' recollection of past interactions (10) or researcher observation of officer behavior (17–20) to assess procedural fairness. Although these methods are invaluable, they offer an indirect view of officer behavior and are limited to a small number of interactions. Furthermore, the very presence of researchers may influence the police behavior those researchers seek to measure (21).

## Significance

**Police officers speak significantly less respectfully to black than to white community members in everyday traffic stops, even after controlling for officer race, infraction severity, stop location, and stop outcome. This paper presents a systematic analysis of officer body-worn camera footage, using computational linguistic techniques to automatically measure the respect level that officers display to community members. This work demonstrates that body camera footage can be used as a rich source of data rather than merely archival evidence, and paves the way for developing powerful language-based tools for studying and potentially improving police–community relations.**

Author contributions: R.V., N.P.C., D. Jurafsky, and J.L.E. designed research; R.V. and N.P.C. performed research; V.P., W.L.H., R.C.H., C.M.G., and D. Jurgens contributed new reagents/analytic tools; R.V. and N.P.C. analyzed data; R.V., N.P.C., D. Jurafsky, and J.L.E. wrote the paper; and D. Jurafsky and J.L.E. served as PI on this project.

Reviewers: J.P., University of Texas at Austin; and T.T., Yale Law School.

Conflict of interest statement: J.L.E. was invited by a federal judge and monitor to serve as a Subject Matter Expert to assist with the Oakland Police Department's reform efforts. The assignment began prior to the studies reported here.

Freely available online through the PNAS open access option.

<sup>1</sup>To whom correspondence may be addressed. Email: robvoigt@stanford.edu or jleberhardt@stanford.edu.

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1702413114/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1702413114/-DCSupplemental).

In study 1, human participants rated officer utterances on several overlapping dimensions of respect. With a high degree of agreement, participants inferred these dimensions from officer language. Even though they were not told the race of the stopped driver, participants judged officer language directed toward black motorists to be less respectful than language directed toward whites. In study 2, we build statistical models capable of predicting aspects of respect based on linguistic features derived from theories of politeness, power, and social distance. We discuss the linguistic features that contribute to each model, finding that particular forms of politeness are implicated in perceptions of respect. In study 3, we apply these models to all vehicle stop interactions between officers of the Oakland Police Department and black/white community members during the month of April 2014. We find strong evidence that utterances spoken to white community members are consistently more respectful, even after controlling for contextual factors such as the severity of the offense or the outcome of the stop.

## Data

Our dataset consists of transcribed body camera footage from vehicle stops of white and black community members conducted by the Oakland Police Department during the month of April 2014. We examined 981 stops of black ( $N = 682$ ) and white ( $N = 299$ ) drivers from this period, 68.1% of the 1,440 stops of white and black drivers in this period. These 981 stops were conducted by 245 different officers (see *SI Appendix, Data Sampling Process* for inclusion criteria). Per Oakland Police Department policy, officers turn on their cameras before making contact with the driver and record for the duration of the stop. From the 183 h of footage in these interactions, we obtain 36,738 usable officer utterances for our analysis.

**Study 1: Perceptions of Officer Treatment from Language.** We first test whether human raters can reliably judge respect from officers' language, and whether these judgments reveal differences in officer respect toward black versus white community members.

Respect is a complex and gradient perception, incorporating elements of a number of correlated constructs like friendliness and formality. Therefore, in this study, we ask participants to rate transcribed utterances spoken by officers along five conceptually overlapping folk notions related to respect and officer treatment. We randomly sampled 414 unique officer utterances (1.1% of all usable utterances in the dataset) directed toward black ( $N = 312$ ) or white ( $N = 102$ ) community members. On each trial, participants viewed the text of an officer utterance, along with the driver's utterance that immediately preceded it. All proper names and places were anonymized, and participants were not told the race or gender of the driver. Participants indicated on four-point Likert scales how respectful, polite, friendly, formal, and impartial the officer was in each exchange. Each utterance was rated by at least 10 participants.

Could participants reliably glean these qualities from such brief exchanges? Previous work has demonstrated that different perceivers can arrive at similar judgments from "thin slices" of behavior (22). In a similar vein, participants showed consistency in their perceptions of officer language, with reliability for each item ranging from moderate (Cronbach's  $\alpha = 0.73$ ) to high ( $\alpha = 0.91$ ) agreement (see *SI Appendix, Annotator Agreement*). These results demonstrate that transcribed language provides a sufficient and consensual signal of officer communication, enough to gain a picture of the dynamics of an interaction at a given point in time.

To test whether participant ratings uncovered racial group differences, we averaged scores across raters to calculate a single rating on each dimension for each utterance, then built a linear mixed-effects regression model to estimate the fixed

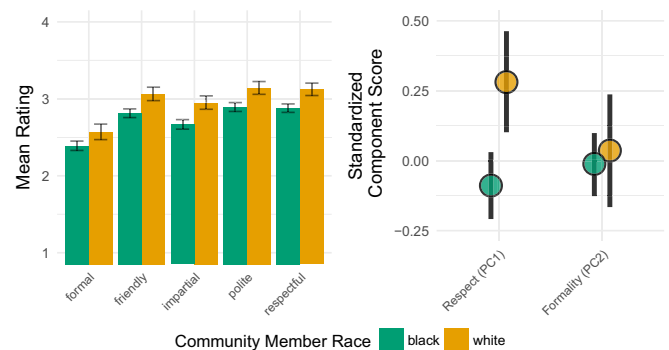
effect of community member race across interactions, controlling for variance of a random effect at the interaction level. Officer utterances directed toward black drivers were perceived as less respectful [ $b = -0.23$ , 95% confidence interval ( $-0.34$ ,  $-0.11$ )], polite [ $b = -0.23$  ( $-0.35$ ,  $-0.12$ )], friendly [ $b = -0.24$  ( $-0.36$ ,  $-0.12$ )], formal [ $b = -0.16$  ( $-0.30$ ,  $-0.03$ )], and impartial [ $b = -0.26$  ( $-0.39$ ,  $-0.12$ )] than language directed toward white drivers (Fig. 1). These differences persisted even when controlling for the age and sex of the driver (see *SI Appendix, Model Outputs for Each Rated Dimension*).

Given the expected conceptual overlap in the five perceptual categories we presented to the participants, we used principal component analysis to decompose the ratings into their underlying components. Two principal components explained 93.2% of the variance in the data (see *SI Appendix, Principal Component Analysis (PCA) Loadings* for loadings). The first component, explaining 71.3% of the variance and composed of positive loadings on the impartial, respectful, friendly, and polite dimensions with some loading on the formal dimension, we characterize as Respect, broadly construed. The second, explaining 21.9% of the variance and composed primarily of a very high positive loading on the formal dimension and a weak negative loading on the friendly dimension, we characterize as Formality. This component captures formality as distinct from respect more generally, and is likely related to social distance.

Standardizing these factor scores as outcome variables in mixed-effects models, we find that officers were equal in Formality with white and black drivers [ $\beta = -0.01$  ( $-0.19$ ,  $0.16$ )], but higher in Respect with white drivers [ $\beta = 0.17$  ( $0.00$ ,  $0.33$ )] (Fig. 1).

Study 1 demonstrates that key features of police treatment can be reliably gleaned from officer speech. Participant ratings from thin slices of police-community interactions reveal racial disparities in how respectful, impartial, polite, friendly, and formal officers' language to community members was perceived. Such differences were driven by differences in the Respect officers communicated toward drivers rather than the Formality with which officers addressed them.

**Study 2: Linguistic Correlates of Respect.** The methods of study 1 (human coding of 414 individual utterances), although effective at discovering racial disparities in officer respect toward community members in our dataset, cannot offer a general solution to the analysis of body camera data. One problem is scale: Each year, on the order of 26 million vehicle stops are made (5). Furthermore, using only a small sample of individual utterances makes it impossible to study how police treatment varies over officers, or how the interaction progresses across time in each stop.

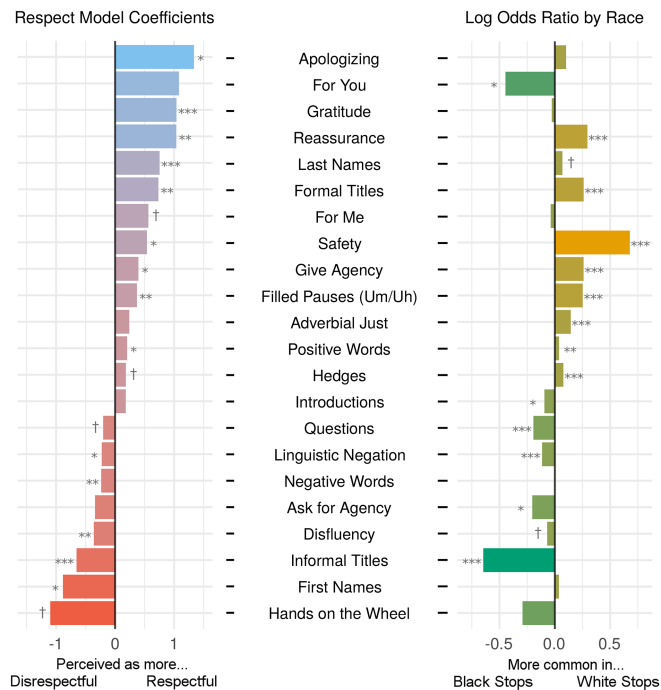


**Fig. 1.** (Left) Differences in raw participant ratings between interactions with black and white community members. (Right) When collapsed to two uncorrelated components, Respect and Formality, we find a significant difference for Respect but none for Formality. Error bars represent 95% confidence intervals. PC, principal component.

In this study, we therefore develop computational linguistic models of respect and formality and tune them on the 414 individual utterances; in study 3, we apply these models to our full dataset of 36,738 utterances. Our method is based on linguistic theories of respect that model how speakers use respectful language (apologizing, giving agency, softening of commands, etc.) to mitigate “face-threatening acts.” We use computational linguistic methods (e.g., refs. 23–26) to extract features of the language of each officer utterance. The log-transformed counts of these features are then used as independent variables in two linear regression models predicting the perceptual ratings of Respect and Formality from study 1.

Our model-assigned ratings agree with the average human from study 1 about as well as humans agree with each other. Our model for Respect obtains an adjusted  $R^2$  of 0.258 on the perceptual ratings obtained in study 1, and a root-mean-square error (RMSE) of 0.840, compared with an RMSE of 0.842 for the average rater relative to other raters. Our model for Formality obtains an adjusted  $R^2$  of 0.190, and an RMSE of 0.882 compared with 0.764 for the average rater (see *SI Appendix, Model Comparison to Annotators* for more details on how these values were calculated). These results indicate that, despite the sophisticated social and psychological cues participants are likely drawing upon in rating officers’ utterances, a constrained set of objectively measurable linguistic features can explain a meaningful portion of the variance in these ratings.

Fig. 2 lists the linguistic features that received significant weights in our model of Respect (arranged by their model coefficients). For example, apologizing, gratitude, and expressions of concern for citizen safety are all associated with respect. The bars on the right show the log-odds of the relative proportion of interactions in our dataset taken up by each feature, where negative numbers mean that a feature comprised a larger proportion of officers’ speech in interactions with black community members and positive numbers mean the same for interactions



**Fig. 2.** (Left) Respect weights assigned by final model to linguistic features and (Right) the corresponding log-odds of those features occurring in officer speech directed toward black versus white community members, calculated using Fisher’s exact test. † $P < 0.1$ ; \* $P < 0.05$ ; \*\* $P < 0.01$ ; \*\*\* $P < 0.001$ .

EXAMPLE	RESPECT SCORE
<p>FIRST NAME    ASK FOR AGENCY    QUESTIONS</p> <p>[name], can I see that driver's license again?</p> <p>It- it's showing suspended. Is that- that's you?</p> <p>DISFLUENCY    NEGATIVE WORD    DISFLUENCY</p>	-1.07
<p>INFORMAL TITLE    ASK FOR AGENCY    ADVERBIAL "JUST"</p> <p>All right, my man. Do me a favor. Just keep your hands on the steering wheel real quick.</p> <p>"HANDS ON THE WHEEL"</p>	-0.51
<p>APOLOGY    INTRODUCTION    LAST NAME</p> <p>Sorry to stop you. My name's Officer [name] with the Police Department.</p>	0.84
<p>FORMAL TITLE    SAFETY    PLEASE</p> <p>There you go, ma'am. Drive safe, please.</p>	1.21
<p>ADVERBIAL "JUST"    FILLED PAUSE    REASSURANCE</p> <p>It just says that, uh, you've fixed it. No problem. Thank you very much, sir.</p> <p>GRATITUDE    FORMAL TITLE</p>	2.07

**Fig. 3.** Sample sentences with automatically generated Respect scores. Features in blue have positive coefficients in the model and connote respect, such as offering reassurance (“no problem”) or mentioning community member well-being (“drive safe”). Features in red have negative coefficients in the model and connote disrespect, like informal titles (“my man”), or disfluencies (“that- that’s”).

with white community members. Example utterances containing instances of the highest-weighted features for the Respect model are shown in Fig. 3. See *SI Appendix, Study 2* for full regression outputs and more detailed discussion of particular linguistic findings.

**Study 3: Racial Disparities in Respect.** Having demonstrated that people can reliably infer features of procedural justice from officer speech (study 1), and that these ratings can be reliably predicted from statistical models of linguistic features (study 2), we are now able to address our central question: Controlling for contextual factors of the interaction, is officers’ language more respectful when speaking to white as opposed to black community members?

We apply our models from study 2 to the entire corpus of transcribed interactions to generate predicted scores for Respect and Formality for each of the 36,738 utterances in our dataset. We then build linear mixed-effects models for Respect and Formality over these utterances. We include, as covariates in our primary model, community member race, age, and gender; officer race; whether a search was conducted; and the result of the stop (warning, citation, or arrest). We include random intercepts for interactions nested within officers.

Controlling for these contextual factors, utterances spoken by officers to white community members score higher in Respect [ $\beta = 0.05$  (0.03, 0.08)]. Officer utterances were also higher in

Respect when spoken to older [ $\beta = 0.07$  (0.05, 0.09)] community members and when a citation was issued [ $\beta = 0.04$  (0.02, 0.06)]; Respect was lower in stops where a search was conducted [ $\beta = -0.08$  (-0.11, -0.05)]. Officer race did not contribute a significant effect. Furthermore, in an additional model on 965 stops for which geographic information was available, neither the crime rate nor density of businesses in the area of the stop were significant, although a higher crime rate was indicative of increased Formality [ $\beta = 0.03$  (0.01, 0.05)].

One might consider the hypothesis that officers were less respectful when pulling over community members for more severe offenses. We tested this by running another model on a subset of 869 interactions for which we obtained ratings of offense severity on a four-point Likert scale from Oakland Police Department officers, including these ratings as a covariate in addition to those mentioned above. We found that the offense severity was not predictive of officer respect levels, and did not substantially change the results described above.

To consider whether this disparity persists in the most “everyday” interactions, we also reran our analyses on the subset of interactions that did not involve arrests or searches ( $N = 781$ ), and found the results from our earlier models were fundamentally unchanged. Full regression tables for all models described above are given in *SI Appendix, Study 3*.

Another hypothesis is that the racial disparities might have been caused by officers being more formal to white community members, and more informal or colloquial to black community members. However, we found that race was not associated with the formality of officers’ utterances. Instead, utterances were higher in Formality in interactions with older [ $\beta = 0.05$  (0.03, 0.07)] and female [ $\beta = 0.02$  (0.00, 0.04)] community members.

Are the racial disparities in the respectfulness of officer speech we observe driven by a small number of officers? We calculated the officer-level difference between white and black stops for every officer ( $N = 90$ ) in the dataset who had interactions with both blacks and whites (Fig. 4). We find a roughly normal distribution of these deltas for officers of all races. This contrasts with the case of stop-and-frisk, where individual outlier officers account for a substantial proportion of racial disparities (27); the disparities we observe here cannot be explained by a small number of extreme officers.

Because our model is able to generate scores across all utterances in our dataset, we can also consider aspects of the trajectory of interactions beyond the mean level of respect (Fig. 5). Growth-curve analyses revealed that officers spoke with greater Respect [ $b = 0.35$  (0.29, 0.40)] and reduced Formality [ $b = -0.57$  (-0.62, -0.53)] as interactions progressed. However, these trajectories varied by community member race: Although stops of white and black drivers converged in the Formality expressed during the interaction [ $b = -0.09$  (-0.13, -0.05)], the gap in Respect increased over time [ $b = 0.10$  (0.05, 0.15)]. That is, offi-

cer Respect increased more quickly in interactions with white drivers [ $b = 0.45$  (0.38, 0.54)] than in interactions with black drivers [ $b = 0.24$  (0.19, 0.29)].

**Discussion.** Despite the formative role officer respect plays in establishing or eroding police legitimacy (7), it has been impossible to measure how police officers communicate with the public, let alone gauge racial disparities in officer respect. However, body-worn cameras capture such interactions every day. Computational linguistic techniques let us examine police–community contacts in a manner powerful enough to scale to any number of interactions, but sensitive enough to capture the interpersonal qualities that matter to the police and public alike.

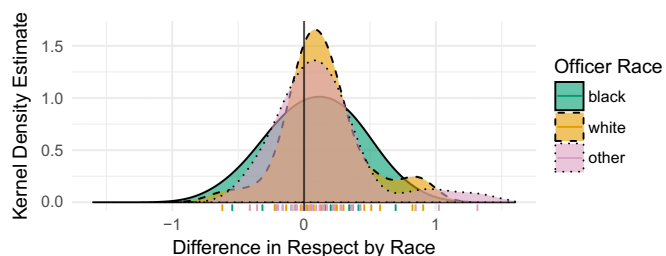
In doing so, we first showed that people make consistent judgments about such interactions from officers’ language, and we identified two underlying, uncorrelated constructs perceived by participants: Respect and Formality. We then built computational linguistic models of these constructs, identifying crucial positive and negative politeness strategies in the police–community interactional context. Applying these models to an entire month of vehicle stops, we showed strong evidence for racial disparities in Respect, but not in Formality: Officers’ language is less respectful when speaking to black community members.

Indeed, we find that white community members are 57% more likely to hear an officer say one of the most respectful utterances in our dataset, whereas black community members are 61% more likely to hear an officer say one of the least respectful utterances in our dataset. (Here we define the top 10% of utterances to be most respectful and the bottom 10% to be least respectful.)

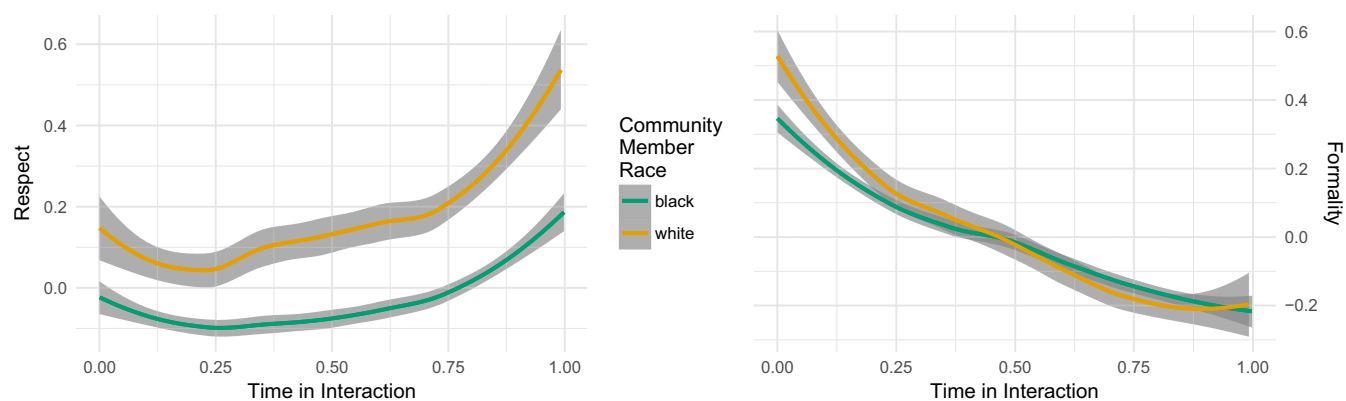
This work demonstrates the power of body camera footage as an important source of data, not just as evidence, addressing limitations with methodologies that rely on citizens’ recollection of past interactions (10) or direct researcher observation of police behavior (17–20). However, studying body camera footage presents numerous hurdles, including privacy concerns and the raw scale of the data. The computational linguistic models presented here offer a path toward addressing both these concerns, allowing for the analysis of transcribed datasets of any size, and generating reliable ratings of respect automatically. These models have the potential to allow for useful information about an interaction to be extracted while maintaining officer and community member privacy.

The racial disparities in officer respect are clear and consistent, yet the causes of these disparities are less clear. It is certainly possible that some of these disparities are prompted by the language and behavior of the community members themselves, particularly as historical tensions in Oakland and preexisting beliefs about the legitimacy of the police may induce fear, anger, or stereotype threat. However, community member speech cannot be the sole cause of these disparities. Study 1 found racial disparities in police language even when annotators judged that language in the context of the community member’s utterances. We observe racial disparities in officer respect even in police utterances from the initial 5% of an interaction, suggesting that officers speak differently to community members of different races even before the driver has had the opportunity to say much at all.

Regardless of cause, we have found that police officers’ interactions with blacks tend to be more fraught, not only in terms of disproportionate outcomes (as previous work has shown) but also interpersonally, even when no arrest is made and no use of force occurs. These disparities could have adverse downstream effects, as experiences of respect or disrespect in personal interactions with police officers play a central role in community members’ judgments of how procedurally fair the police are as an institution, as well as the community’s willingness to support or cooperate with the police (8, 9).



**Fig. 4.** Kernel density estimate of individual officer-level differences in Respect when talking to white as opposed to black community members, for the 90 officers in our dataset who have interactions with both blacks and whites. More positive numbers on the x axis represent a greater positive shift in Respect toward white community members.



**Fig. 5.** Loess-smoothed estimates of the (Left) Respect and (Right) Formality of officers' utterances relative to the point in an interaction at which they occur. Respect tends to start low and increase over an interaction, whereas the opposite is true for Formality. The race discrepancy in Respect is consistent throughout the interactions in our dataset.

We now have a method for quantifying these troubled interactions. Although the circumstances of any particular stop can vary dramatically, our approach allows us to measure aggregate department-level trends, revealing disparities across hundreds of interactions. These disparities are part of a constellation of differences in officer language spoken toward black versus white community members; a simple classifier trained on only the words used by officers is able to correctly predict the race of the community member in over two thirds of the interactions (see *SI Appendix, Linguistic Classification Accuracy of Race*).

Future research could expand body camera analysis beyond text to include information from the audio such as speech intonation and emotional prosody, and video, such as the citizen's facial expressions and body movement, offering even more insight into how interactions progress and can sometimes go awry. In addition, footage analysis could help us better understand what linguistic acts lead interactions to go well, which can inform police training and quantify its impacts over time.

The studies presented here open a path toward these future opportunities and represent an important area of research for the study of policing: Computational, large-scale analyses of language give us a way to examine and improve police–community interaction that we have never had before.

## Materials and Methods

**Data and Processing.** The video for each traffic stop was transcribed into text by professional transcribers, who transcribed while listening to audio and watching the video. Extensive measures were taken to preserve privacy; data were kept on a central server, and transcribers (as well as all researchers) underwent background checks with the Oakland Police Department. Transcribers also “diarized” the text (labeling who was speaking at each time point). We used the diarization to automatically remove all officer speech to the dispatcher or to other officers, leaving only speech from the officer directed toward the community member. After transcription, transcripts were manually cleaned up, heuristically fixing transcriber diarization errors, and correcting typographical errors involving utterance timing so that all transcripts were automatically readable. Every utterance in the dataset was processed with Stanford CoreNLP 3.4.1 (28) to generate sentence and word segmentation, part-of-speech tags, and dependency parses used for feature extraction and analysis.

The raw video footage associated with this paper was available for our research purposes with the cooperation of the Oakland Police Department, and naturally cannot be publicly distributed. However, we make available deidentified data frames for each study described here, so that other researchers can replicate our results. We also release all of the code for the computational linguistic models, as well as pretrained models that can be run on arbitrary text.

**Human Annotation of Utterances.** A subset of 420 exchanges, consisting of one officer utterance (defined as a “turn” of one or more sentences by transcribers) and, if applicable, the immediately preceding community member utterance were sampled from the corpus for annotation. Utterances were sampled with the constraint that at least 15 words were spoken between the two speakers, and that at least five words were spoken by the officer. These utterances were grouped into seven “batches” of 60 utterances apiece. Due to a data error, six duplicate utterances were annotated, but were excluded from subsequent analyses, resulting in 414 unique utterances toward black ( $N = 312$ ) and white ( $N = 102$ ) community members.

Each of 70 participants (39 female,  $M_{age} = 25.3$ ) rated a batch of 60 of these utterances, such that each utterance was rated by at least 10 participants. On each trial, participants viewed the text of an exchange between a police officer and a community member: the text of the officer utterance, as well as the text of the community member utterance that immediately preceded it, if there was one. They then indicated, on four-point bipolar Likert scales, how respectful, polite, friendly, formal, and impartial the officer was in each exchange. Participants were allowed to indicate that they could not rate an utterance on a particular dimension, but were encouraged to nonetheless indicate their best guess. Participants had no other information about the interaction besides the officer's utterance and the immediately preceding community member utterance.

All research was approved by the Stanford University Institutional Review Board, and written informed consent was obtained from all raters before their participation.

**Computational Annotation of Utterances.** Our model draws on linguistic theories of politeness; the technical term “politeness” refers to how concepts like respect, formality, and social distance take shape in language. These theories suggest that speakers use polite or respectful language to mitigate face-threatening acts (29–31).

Negative politeness is used to mitigate direct commands or other impositions that limit the freedom of action of the listener, for example, by minimizing the imposition or emphasizing the agency of the interlocutor. Such strategies are central to police–community interactions because of the inherently coercive nature of a traffic stop. For instance, the use of the word “please” can soften requests and provide a sense of agency or choice; apologizing (“sorry,” “excuse me”) can admit regret on the part of the officer that some request is necessary; the use of hedges (“may,” “kinda,” “probably”) may reduce the perception of imposition.

Positive politeness is used to show that the speaker values the interlocutor and their interests, or to minimize the impact of actions that could damage such a perception. Positive politeness strategies are also crucial for police–community interactions, where the inherently unequal social roles at play may necessitate a particular sensitivity to the community member's positive face. For instance, greetings and introductions can establish a friendly context at the beginning of an interaction and convey openness. Expressions of reassurance (“no big deal,” “don't worry”) seek to assuage the community member's potential concerns in tense circumstances, and expressions of gratitude (“thank you”) serve to reduce the perceived power differential by deferring to the actions of the community member. Mentions of safety (“Drive safely now”) explicitly acknowledge concern for the community member's personal well-being. Referring expressions are another important component of positive politeness;

formal titles (“sir,” “ma’am,” “Mr.,” “Ms.”) and surnames may convey a contrast with informal titles (“dude,” “bro,” “bud”) and first names (31–33).

We also include features we expect to capture officer anxiety, such as speech disfluencies (“w- well”) and commands to keep “hands on the wheel,” which may contribute to a community member’s perception of disrespect. These are of a different character than the politeness strategies discussed above, but we found that all analyses presented here hold true even if these features are not included.

We use standard techniques to automatically extract features from the text of each utterance (23–26). These features include lexicons (lists of words). For example, to detect informal titles, we used an augmented version of a word list from ref. 34. We also used regular expressions, such as for detecting tag questions (“do that for me, will you?”), and syntactic parse

features, such as a feature that detects when “just” is used in constructions as an adverbial modifier.

Features were modeled as log-transformed counts in each utterance, and were used as independent variables in two linear regression models predicting the human perceptual ratings of respect and formality obtained in study 1. They were introduced into the regression using stepwise forward selection by  $R^2$  to remove features that don’t substantially contribute to the model’s accuracy.

**ACKNOWLEDGMENTS.** This research was supported by the John D. and Catherine T. MacArthur Foundation, with additional support from the Stanford Institute for Research in the Social Sciences, the Stanford School of Humanities and Sciences, and the Stanford Data Science Initiative. We also thank the City of Oakland and the Oakland Police Department for their support and cooperation.

1. President’s Task Force on 21st Century Policing (2015) *Final Report of the President’s Task Force on 21st Century Policing* (Off Commun Oriented Policing Serv, Washington, DC).
2. The White House (December 1, 2014) Fact sheet: Strengthening community policing. Press release (Washington, DC). Available at <https://obamawhitehouse.archives.gov/the-press-office/2014/12/01/fact-sheet-strengthening-community-policing>. Accessed February 1, 2017.
3. Reaves B (2015) *Local Police Departments, 2013: Personnel, Policies, and Practices* (US Dep Justice, Washington, DC), NCJ 248677.
4. Eith C, Durose M (2011) *Contacts Between Police and the Public, 2008* (Bur Justice Stat, Washington, DC).
5. Langton L, Durose M (2013) *Special Report: Police Behavior During Traffic and Street Stops, 2011* (Bur Justice Stat, Washington, DC).
6. Tyler TR, Huo Y (2002) *Trust in the Law: Encouraging Public Cooperation with the Police and Courts* (Russell Sage Found, New York).
7. Tyler TR, Blader SL (2003) The group engagement model: Procedural justice, social identity, and cooperative behavior. *Pers Soc Psychol Rev* 7:349–361.
8. Tyler TR, Bies RJ (1990) Beyond formal procedures: The interpersonal context of procedural justice. *Applied Social Psychology and Organizational Settings*, ed Carroll JS (Lawrence Erlbaum, Hillsdale, NJ), pp 77–98.
9. Mazerolle L, Antrobus E, Bennett S, Tyler TR (2013) Shaping citizen perceptions of police legitimacy: A randomized field trial of procedural justice. *Criminology* 51:33–63.
10. Epp CR, Maynard-Moody S, Haider-Markel DP (2014) *Pulled Over: How Police Stops Define Race and Citizenship* (Univ Chicago Press, Chicago).
11. Peffley M, Hurwitz J (2010) *Justice in America: The Separate Realities of Blacks and Whites* (Cambridge Univ Press, New York).
12. Giles H, Coupland J, Coupland N (1991) Accommodation theory: Communication, context and consequences. *Contexts of Accommodation: Developments in Applied Sociolinguistics*, eds Giles H, Coupland J, Coupland N (Cambridge Univ Press, New York), pp 1–68.
13. Gnisci A (2005) Sequential strategies of accommodation: A new method in courtroom. *Br J Soc Psychol* 44:621–643.
14. Ng SH, Bell D, Brooke M (1993) Gaining turns and achieving high influence ranking in small conversational groups. *Br J Soc Psychol* 32:265–275.
15. Nguyen VA, et al. (2014) Modeling topic control to detect influence in conversations using nonparametric topic models. *Mach Learn* 95:381–421.
16. Prabhakaran V, Rambow O (2014) Predicting power relations between participants in written dialog from a single thread. *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics* (Assoc Comput Linguist, Stroudsburg, PA), pp 339–344.
17. Mastrofski SD, Parks RB, McCluskey JD (2010) Systematic social observation in criminology. *Handbook of Quantitative Criminology*, eds Piquero AR, Weisburd D (Springer, New York), pp 225–247.
18. Dai M, Frank J, Sun I (2011) Procedural justice during police-citizen encounters: The effects of process-based policing on citizen compliance and demeanor. *J Crim Justice* 39:159–168.
19. Jonathan-Zamir T, Mastrofski SD, Moyal S (2015) Measuring procedural justice in police-citizen encounters. *Justice Q* 32:845–871.
20. Mastrofski SD, Jonathan-Zamir T, Moyal S, Willis JJ (2016) Predicting procedural justice in police-citizen encounters. *Crim Justice Behav* 43:119–139.
21. Mastrofski S, Parks RB (1990) Improving observational studies of police. *Criminology* 28:475–496.
22. Ambady N, Bernieri FJ, Richeson JA (2000) Toward a histology of social behavior: Judgmental accuracy from thin slices of the behavioral stream. *Adv Exp Soc Psychol* 32:201–271.
23. Tausczik YR, Pennebaker JW (2010) The psychological meaning of words: LIWC and computerized text analysis methods. *J Lang Soc Psychol* 29:24–54.
24. Prabhakaran V, Rambow O, Diab M (2012) Predicting overt display of power in written dialogs. *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics* (Assoc Comput Linguist, Stroudsburg, PA), pp 518–522.
25. Danescu-Niculescu-Mizil C, Lee L, Pang B, Kleinberg J (2012) Echoes of power: Language effects and power differences in social interaction. *Proceedings of the 21st International Conference on World Wide Web* (Assoc Comput Mach, New York), pp 699–708.
26. Danescu-Niculescu-Mizil C, Sudhof M, Jurafsky D, Leskovec J, Potts C (2013) A computational approach to politeness with application to social factors. *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics* (Assoc Comput Linguist, Stroudsburg, PA), pp 250–259.
27. Goel S, Rao JM, Shroff R (2016) Precinct or prejudice? Understanding racial disparities in New York City’s stop-and-frisk policy. *Ann Appl Stat* 10:365–394.
28. Manning CD, et al. (2014) The Stanford CoreNLP natural language processing toolkit. *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics* (Assoc Comput Linguist, Stroudsburg, PA), pp 55–60.
29. Goffman E (1967) On face-work. *Interaction Ritual: Essays on Face-to-Face Behavior* (Anchor, Garden City, NY), pp 5–45.
30. Lakoff RT (1973) The logic of politeness: Minding your p’s and q’s. *Papers from the 9th Regional Meeting of the Chicago Linguistic Society*, eds Corum C, Smith-Stark T, Weiser A (Chicago Linguist Soc, Chicago), pp 292–305.
31. Brown P, Levinson SC (1987) *Politeness: Some Universals in Language Usage* (Cambridge Univ Press, Cambridge, UK).
32. Wood LA, Kroger RO (1991) Politeness and forms of address. *J Lang Soc Psychol* 10:145–168.
33. Boxer D (1993) Social distance and speech behavior: The case of indirect complaints. *J Pragmat* 19:103–125.
34. Krishnan V, Eisenstein J (2015) “You’re Mr. Lebowsky, I’m the Dude”: Inducing address term formality in signed social networks. *Proceedings of the North American Chapter of the Association for Computational Linguistics*, eds Elangovan V, Eisenstein J (Assoc Comput Linguist, Stroudsburg, PA), pp 1616–1626.