

Trumping Hate on Twitter? Online Hate Speech in the 2016 U.S. Election Campaign and its Aftermath

Alexandra A. Siegel^{1,2}, Evgenii Nikitin^{2,3}, Pablo Barberá^{2,4},
Joanna Sterling^{2,5}, Bethany Pullen², Richard Bonneau^{2,6}, Jonathan Nagler^{2,3}
and Joshua A. Tucker^{2,3*}

¹*University of Colorado Boulder, Department of Political Science, Boulder, CO, USA; alexandra.siegel@colorado.edu*

²*New York University, Center for Social Media and Politics (CSMaP), New York, NY, USA; bethanyjpullen@gmail.com*

³*New York University, Department of Politics, New York, NY, USA; e.nikitin@nyu.edu; jonathan.nagler@nyu.edu; joshua.tucker@nyu.edu*

⁴*University of Southern California, Department of Political Science and International Relations, Los Angeles, CA, USA; pbarbera@usc.edu*

⁵*Princeton University, Department of Psychology, Princeton, NJ, USA, New York University, New York, NY, USA; joanna.sterling@princeton.edu*

⁶*New York University, Department of Biology, and Simons Foundation, New York, NY, USA; rbonneau@flatironinstitute.org*

*The authors gratefully acknowledge the financial support for the NYU Social Media and Political Participation (SMaPP) lab from the INSPIRE program of the National Science Foundation (Award SES-1248077), the William and Flora Hewlett Foundation, the Rita Allen Foundation, the Knight Foundation, the Bill and Melinda Gates Foundation, Craig Newmark Philanthropies, the Democracy Fund, the Intel Corporation, the New York University Global Institute for Advanced Study, and the Faculty of Arts and Sciences Research Investment Fund at New York University. We thank Sean Kates for his feedback in designing our coding scheme, NYU Undergraduate SMaPP Research Assistants for their coding work, and Yvan Scher and Leon Yin for programming support. A.S. and J.T. designed the research plan and outline for the paper. A.S. conducted the statistical analysis and wrote the first draft of the paper. A.S, J.S., and B.P. designed and implemented the dictionary-based coding method. E.N. designed and conducted the non-dictionary based analysis. P.B., R.B., and J.N. contributed to the data collection and design of the analytic tools, and strategy for data analysis and presentation. AS, JT, PB, RB, and JN contributed to revising the manuscript.

Online Appendix available from:

http://dx.doi.org/10.1561/100.00019045_app

Supplementary Material available from:

http://dx.doi.org/10.1561/100.00019045_supp

MS submitted on 20 March 2019; final version received 27 March 2020

ISSN 1554-0626; DOI 10.1561/100.00019045

© 2021 A. A. Siegel *et al.*

ABSTRACT

To what extent did online hate speech and white nationalist rhetoric on Twitter increase over the course of Donald Trump's 2016 presidential election campaign and its immediate aftermath? The prevailing narrative suggests that Trump's political rise — and his unexpected victory — lent legitimacy to and popularized bigoted rhetoric that was once relegated to the dark corners of the Internet. However, our analysis of over 750 million tweets related to the election, in addition to almost 400 million tweets from a random sample of American Twitter users, provides systematic evidence that hate speech did not increase on Twitter over this period. Using both machine-learning-augmented dictionary-based methods and a novel classification approach leveraging data from Reddit communities associated with the alt-right movement, we observe no persistent increase in hate speech or white nationalist language either over the course of the campaign or in the six months following Trump's election. While key campaign events and policy announcements produced brief spikes in hateful language, these bursts quickly dissipated. Overall we find no empirical support for the proposition that Trump's divisive campaign or election increased hate speech on Twitter.

Keywords: Hate speech; social media; Donald Trump; Twitter; text-as-data

From his calls for a “Muslim ban” to his retweets of white nationalist content and leaked tapes endorsing sexual assault, experts and casual observers alike warned that Donald Trump's divisive 2016 campaign and election fanned the flames of bigotry across the United States (Camp, 2016). Increased reports of hate crimes and vandalism targeting religious and racial minorities, as well as a notable rise in the number of active American hate groups, caused 2016 to be described as an “unprecedented year for hate” (SPLC, 2017). Citing a “massive rise” in online hate speech following Trump's election, reports by civil rights organizations suggest that Trump's campaign legitimized extremist ideologies, bringing hostile messages that were once relegated to the dark corners of the Internet into the mainstream (ADL, 2017a; Guynn, 2016). Fearing that Trump's election created a new “safe space for hate,” academics, journalists, policy makers, and everyday citizens have increasingly voiced concern about the consequences of Trump's actions and rhetoric both online and offline

(Milligan, 2017; Müller and Schwarz, 2018; Ott, 2017). As a result, proposals for regulation to control online hate speech, both in the United States and abroad, have become increasingly common.¹ Analyzing the impact of Trump’s Twitter rhetoric on public discourse, a 2017 study concluded that “Trump’s simple, impulsive, and uncivil Tweets do more than merely reflect sexism, racism, homophobia, and xenophobia; they spread those ideologies like a social cancer...His Tweets teach us to see others as less-than-human and they inspire hatred and violence”(Ott, 2017).

However, despite a wealth of anecdotal and small-scale empirical evidence of this “Trump effect,” little is known about how the quantity of online hate speech, or the number of individuals producing it, changed over the course of Trump’s 2016 campaign or in the aftermath of his election.² Here, we take a first step at addressing this gap in knowledge about the link between Trump’s political rise and the use of hate speech online. More specifically, using Twitter collections of over 150 million tweets referencing Hillary Clinton and related keywords, over 600 million tweets referencing Donald Trump and related keywords, and over 400 million tweets collected from a random sample of 500,000 American Twitter users, we systematically investigate the degree to which the quantity of hateful tweets and size of the population producing them on Twitter increased over the course of Trump’s campaign or following his election on November 8, 2016.

We identify hate speech and white nationalist language using two methods: a machine-learning-augmented dictionary-based approach; and a non-dictionary-based method harnessing large naturally annotated corpora of text containing hate speech and white nationalist language from alt-right subreddits.³ Using two different sources of data and these two different measurement strategies, we find — in contrast to the conventional wisdom — no persistent increase in hate speech or white nationalist rhetoric on Twitter, Trump’s preferred social media platform. While key events including terror attacks or Trump’s travel ban announcement produce temporary upticks in hateful rhetoric — and these bursts of hate speech are not inconsequential — they are not indicative of a systematic increase of hate speech in the American Twittersphere over the course of the 2016 US Presidential election campaign and its immediate aftermath, and appear similar to the bursts of such rhetoric regularly observed in the pre-Trump era.

¹See, for example, Financial Times (2017), Marwick (2017), and Rainie *et al.* (2017).

²Indeed, a parallel concern — about the rise of bots as a vehicle for spreading online hate — might in fact suggest that current estimates of the number of *people* producing hate speech on Twitter could be biased upwards.

³Subreddits are anonymous online forums dedicated to discussing specific topics on Reddit, a popular social news aggregation, web content rating, discussion forum, and social media platform. This method is explained in detail in the discussion of our analysis below, as well in Section A2 in the online appendix.

Our study therefore highlights the need to move beyond conventional narratives and small-scale empirical analysis to more systematically assess trends in online discourse and behavior over time. Trump's campaign and election undoubtedly drew a great deal of attention to online hate speech and white nationalist rhetoric, but, at least on the portions of Twitter we examined, such language did not become more common or popular. The extent to which this trend may have materialized in other online and offline forums — as well as in other parts of the Twittersphere not captured by our analysis — is beyond the scope of this paper and remains a subject for future research.

Beyond our empirical findings, this work also contributes to the development of methods for understanding the evolution of speech in the digital era. As the role of social media platforms in fostering extremism and offline violence has come under scrutiny, online hate speech has received increased attention from academics and policy makers alike. But despite a growing body of research devoted to defining and detecting online hate speech, the existing scientific literature lacks a systematic framework for assessing how the volume and content of these harmful messages change over time. We hope that the research design applied here — relying on multiple sources of data and utilizing multiple methods to analyze changes in online speech — might be used by other researchers to evaluate the real-time dynamics of online discourse in diverse contexts.

Motivation and Expectations

Beginning in the early days of his campaign, Donald Trump's rise as a mainstream political candidate was marked by unusually blunt bigoted and racist statements. From his assertion in July 2015 that Mexican immigrants were criminals and rapists (Walker, 2015) to his repeated derogatory remarks against Muslims including phrases like "Islam hates us", his campaign rhetoric directed at minority groups and immigrants was far from subtle (Filimon, 2016). His courting of support from anti-LGBT groups stoked fears in the LGBT community (Stack, 2016), and his use of misogynistic rhetoric, including repeated blaming of sexual assault victims and a leaked tape in which he advocated groping women against their will, highlighted his lack of respect for women (Cohen, 2016). His refusal to denounce the support of former Ku Klux Klan leader David Duke, not to mention his post-election appointment of Steve Bannon, the self-proclaimed creator of a "platform for the alt right," lent further credence to the view that the Trump administration tolerated, if not condoned, extremist white nationalist ideology (Huber, 2016).

During the campaign, round-the-clock traditional and social media coverage of Donald Trump were littered with examples of his prejudicial rhetoric and behavior. Trump's prolific Twitter use far surpassed that of all other candidates

and became “a tool of political promotion, distraction, score-settling and attack”(Barbaro, 2015). Many of Trump’s most inflammatory statements during the campaign were 140-character statements disseminated through his Twitter account. These were then repeatedly amplified and defended by his growing base of online supporters (Wells *et al.*, 2016). This endless stream of social media content meant that Trump was regularly trending on Twitter, providing constant fodder for journalists covering the campaign. As Wells *et al.* (2016) describe, Trump attracted a great deal of attention by taking advantage of this hybrid media system. He received media attention through conventional channels including rallies, press conferences, and interviews — not to mention uninvited call-ins to radio and television programs. But he also unleashed “tweetstorms” that galvanized his supporters and helped catapult him to unmatched media coverage in the 2016 campaign.

As Trump’s campaign gained momentum, reports began to emerge of increased hate speech, bias incidents, and hate crimes — with some perpetrators explicitly claiming that Trump had motivated their actions. Attracting more mainstream coverage with the release of Univision anchor Jorge Ramos’ widely publicized documentary, “Hate Rising,” stories of how the Trump campaign was emboldening hate groups and giving rise to a new wave of anti-minority hostility proliferated (Montagne, 2016). Graphic evidence of anti-Semitic harassment of Jewish journalists on Twitter by Trump supporters gained widespread attention as well, frequently accompanied by death threats and Holocaust imagery (Rapaport, 2016).

These reports continued and intensified in the aftermath of Trump’s election. Articles like the *New Yorker*’s “Hate Is on the Rise After Trump’s Election,” The *Guardian*’s “Trump’s Election led to Barrage of Hate,” and *Vox*’s “The Wave of Post-Election Hate Reportedly Sweeping the Nation, Explained,” became increasingly widespread. James King’s year-in-review column, “The Year in Hate: From Donald Trump to the Rise of the Alt-Right,” *Salon*’s “A Short History of Hate” which tracks the alt-right’s 2016 ascendance, and the *New York Times*’ hate-speech aggregator, “This Week in Hate,” are just a few examples of this trend (Duncan, 2017).⁴

⁴Links to the aforementioned news stories are listed here: Desmond-Harris, Jenna. “The Wave of Post-Election Hate Reportedly Sweeping the Nation, Explained.” *Vox* 17 Nov. 2016. <http://www.vox.com/2016/11/17/13639138/trump-hate-crimes-attacks-racism-xenophobia-islamophobia-schools>; King, James. “This Year in Hate.” *Vocativ* 12 Dec. 2016. <http://www.vocativ.com/383234/hate-crime-donald-trump-alt-right-2016/>; Okeowo, Alexis. “Hate on the Rise after Trump’s Election.” *New Yorker* 17 Nov. 2016. <http://www.newyorker.com/news/news-desk/hate-on-the-rise-after-trumps-election>; Sidahmed, Mazin. “Trump’s Election Led to ‘Barrage of Hate’, Report Finds.” *The Guardian* 29 Nov. 2016. <https://www.theguardian.com/society/2016/nov/29/trump-related-hate-crimes-report-southern-poverty-law-center>; Weisberg, Jacob. “The Alt-Right and a Deluge of Hate.” *Slate* 1 Nov. 2016. http://www.slate.com/articles/podcasts/trumpcast/2016/11/how_the_alt_right_harassed_david_french_on_twitter_and_at_home.html.

This anecdotal evidence of increased hate crimes and hate speech suggests that Trump's rise may have played a role in legitimizing and mainstreaming extremist rhetoric. It is also consistent with a long-standing political science literature on the effects of elite cuing on mass attitudes and behaviors. The bulk of this literature has tended to focus on the effects of elite cuing on public opinion formation and policy preferences (Brader *et al.*, 2013; Zaller, 1992, 1994). Perhaps of greater relevance to understanding the purported "Trump effect" on rising hate speech, the social movements literature has demonstrated that elite cues can create discursive opportunity structures that make it easier for far right movements or ideologies to gain traction (Giugni *et al.*, 2005; Koopmans and Muis, 2009; Koopmans and Olzak, 2004). Exploring the effect of elite cues on the tenor of mass rhetoric, political scientists have demonstrated that elites play an important role in changing social norms and spreading racist and intolerant discourse in a variety of cultural contexts (Siegel and Badaan, 2020; Van Dijk, 1992). This finding is also reflected in the communications literature on the spiral of silence, which suggests that people are alerted to social norms by elites, particularly the media. While this theory was first formulated to explain social desirability biases in survey responses (Glynn *et al.*, 1997; Noelle-Neumann, 1974; Scheufle and Moy, 2000), it has also been used to explain self-expression on social networking sites (Fox and Warber, 2014; Lee and Chun, 2016; Pang *et al.*, 2016; Sherrick and Hoewe, 2018). Moreover, recent work in the Arab context demonstrates that elites play an important role in instigating and spreading hate speech and intolerant rhetoric in the online sphere (Siegel, 2015; Siegel *et al.*, 2018, 2020). Academics have also explicitly theorized about Trump's role in legitimizing fringe groups and ideologies (Barkun, 2017; Ott, 2017). These findings from diverse bodies of literature all highlight mechanisms by which the spread of online hate speech during the 2016 election campaign and its aftermath may have been exacerbated by Trump's divisive campaign and election.

While both social science theory and popular narratives suggest that Trump's campaign and election may have increased the popularity of hate speech online, existing evidence is largely anecdotal. Here we provide the first large-scale empirical test of this relationship. More specifically, we test the extent to which hate speech and white nationalist rhetoric on Twitter increased over the course of the campaign and/or following the election. What exactly would such empirical evidence look like in practice? On the one hand we might expect that Trump's political rise throughout the 2016 campaign, and the period following his election, were accompanied by a steady increase in hate speech (a positive relationship between the time since Trump declared his candidacy and the prevalence of hate speech). On the other hand we might expect that Trump's unexpected election led to a sudden increase in

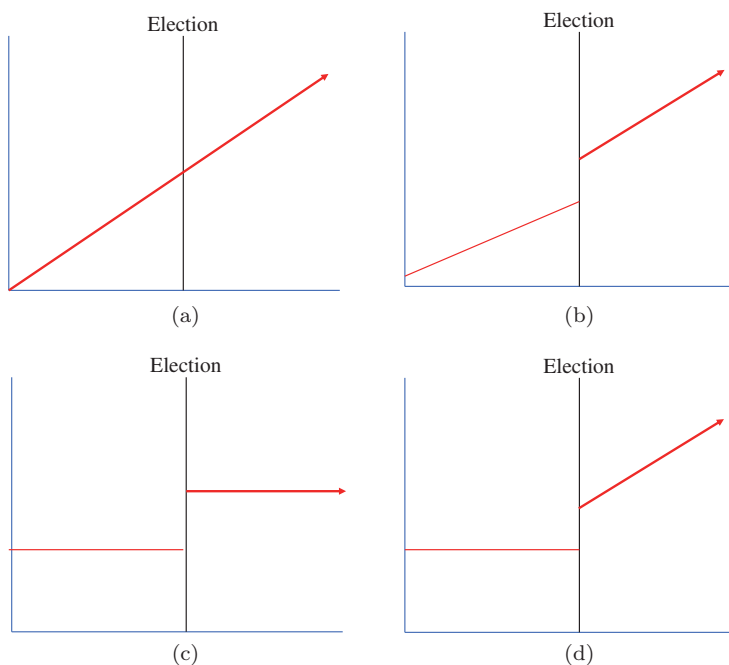


Figure 1: Hypothesized effects of 2016 election campaign and election on hate speech. Here the y -axis is the prevalence of hate speech and the x -axis represents time before and after November 8, 2016, which is marked by the vertical “Election” line.

hate speech (a discontinuity surrounding Trump’s victory). We can therefore imagine several different patterns that might be evidence of a lasting effect of the campaign and/or election on the popularity of online hate speech. These representations are displayed in Figure 1.

If Trump’s rise as a mainstream candidate and his divisive campaign increased the popularity of hateful language, then we might expect to see a positive upward slope in hate speech over the course of the campaign and through the election period (panel (a) of Figure 1).⁵ If the election itself — Trump’s unexpected victory — impacted the popularity of hate speech, then

⁵As we will explain in much greater detail in the following section, there are a variety of different ways to measure the prevalence of hate speech. We rely on four such measures for our dictionary-based analysis: the counts of tweets containing hate speech; the count of users producing hate speech; the proportion of tweets containing hate speech; and the proportion of users producing hate speech. For the non-dictionary based methods, we rely on changes in the semantic similarity of daily tweets over time to naturally annotated corpora of online hate speech.

we should see a discontinuity or jump after the election (panel (c) of Figure 1) and/or an upward trend in the use of hate speech following Trump's victory (panel (d) of Figure 1), relative to the pre-election period. We might also observe a combination of increasingly hateful rhetoric over the course of the campaign as well as a jump following the election (panel (b) of Figure 1). Other combinations, such as a positive increase over the course of the campaign (as in panel (a) of Figure 1 or panel (b) of Figure 1 in the pre-election period), followed by a sustained or flat level of online hate following the election (as in panel (c) of Figure 1), with or without a discontinuity at the election, would also be evidence of this relationship.

Data and Measurement

In this section, we first describe the primary data sources used in our analyses. We then present our machine-learning-augmented dictionary-based approach to measuring hate speech and white nationalist language. To analyze the extent to which Trump's campaign and election were associated with an increase in hate speech and white nationalist rhetoric on Twitter, we rely on two sources of data. Our *Political Twitter* data sets contain all tweets referencing Donald Trump and Hillary Clinton produced between June 17, 2015 and June 15, 2017. These include tweets directly mentioning the candidates using their Twitter handles, such as tweets that include @realDonaldTrump or @HillaryClinton, and any tweets that contain the names of the candidates.⁶ We began collecting tweets mentioning Donald Trump on June 17, 2015, the day after he declared his candidacy. Our analysis therefore covers the entire period of his campaign as well as more than half a year after his election. This Political Twitter data set contains over 150 million tweets related to Hillary Clinton and over 600 million tweets related to Donald Trump. This gives us a very large snapshot of political discourse throughout the 2016 campaign and election period.⁷

However, in order to test the degree to which hate speech and white nationalist language became more mainstream over the course of the 2016 presidential election campaign and its aftermath, we also move beyond the

⁶These tweets were collected using Twitter's streaming API and capture all tweets sent in this period mentioning the candidates, with the exception of infrequent missing data caused by Twitter's rate limits. Because the streaming API allows us to collect data in real time, tweets that were ultimately deleted or produced by accounts that were later suspended are still included in our counts and analyses, thus alleviating concerns that we might be under-counting incidents of hate speech for these reasons (deleted tweets or suspended accounts). Note that this would not have been the case had we used historical data from Twitter or GNIP provided at a later date.

⁷While our Trump and Clinton collections of course do not encompass all discussion of politics in this period, this is the subset of online political discourse where we might most expect to observe a Trump effect and constitutes a large sample of popular political discourse in this period.

political Twittersphere to assess these trends in the U.S. Twittersphere more broadly. To do so, we rely on a collection of tweets sent by a random sample of 500,000 American Twitter users, sampled by generating random user IDs and then checking that those random numbers correspond to accounts that were active and located in the United States.⁸ In the period under study, from June 17, 2015 to June 15, 2017, this collection contained approximately 400 million tweets. Together, these data sets enable us to test whether there was a relationship between Trump’s political rise and the prevalence of online hate speech — in any of the ways identified in the previous section and illustrated in Figure 1 — in either political or general discourse on Twitter.

Systematically measuring shifts in the popularity of hate speech and white nationalist language over the course of the 2016 campaign and election first requires defining and operationalizing these terms. There is no agreed upon definition of hate speech, and the topic has been hotly debated by academics and policymakers alike. In general, there are two primary tendencies in the literature. At one end of the spectrum are broad, more comprehensive, definitions that are designed to identify hate speech in a variety of incarnations. At the other end are narrower definitions, which characterize hate speech as “dangerous speech” that is explicitly intended to incite violence or to call for threatening action against an outgroup.⁹

Attempting to gain a more systematic understanding of the use of online hate speech in the 2016 election period, we define hate speech as *bias-motivated, hostile and malicious language targeted at a person or group because of their actual or perceived innate characteristics, especially when the group or individual are unnecessarily labeled* (Cohen-Almagor, 2011). By this definition, endorsements of groups associated with hate crimes or bias-motivated behavior, such as the Ku Klux Klan, or statements showing excessive pride in one’s own race or group, do not constitute hate speech. While such content is often offensive and frequently occurs alongside hateful language, following previous studies we define hate speech as requiring a disparagement of others (Warner and Hirschberg, 2012).¹⁰

⁸The random selection was achieved by sampling users based on their numeric ID: first, we generated random numbers between 1 and the highest numeric ID assigned at the time of data collection (3.3 billion); then, for each number we checked whether the user existed, and whether the ‘time_zone’ field in their profile was one of the time zones in the United States or whether their ‘location’ field mentioned the full name or abbreviation of a U.S. state or one of the top 1,000 most populated cities; if the user met one of these conditions, they were included in the sample.

⁹See Gagliardone *et al.* (2014), Kennedy *et al.* (2018), and Siegel (2020) for a detailed overview of the literature on defining hate speech on and offline.

¹⁰As will be discussed in more detail later, this “definitional” discussion guides our approach to dictionary based methods for identifying hate speech. Our non-dictionary based approach measure the extent to which speech on Twitter “resembles” (i.e., shares semantic similarity with) hate-speech as utilized in practice online. See the section “Robustness Check: Reference-Text Based Analysis” below.

However, because we are also interested in the extent to which white nationalist or extreme right wing rhetoric increased over the course of Trump’s campaign or following his election, we explicitly measure white nationalist rhetoric as well. Following the social science literature on white nationalism, we define this language as content praising or associated with “white nationalist” groups also known as “racist right-wing,” “extreme right,” “far right,” or “hate” groups.¹¹ What distinguishes white nationalist ideology is its concern with protecting white racial identity. White nationalist groups ascribe to the belief that people of white European backgrounds are a separate community based on myths of common ancestry and culture that transcend national boundaries (Kaplan, 2000). In particular, they oppose anything that they believe will dilute the purity of their exceptionalist white culture including immigration, interracial marriage, globalization, and multiculturalism (Fording, 2014). For the purposes of this paper, we define white nationalist rhetoric — in contrast to hate speech — as *any rhetoric or content that praises known white-nationalist groups, shows excessive pride in the white race, espouses white supremacist or white separatist ideologies, or focuses on the alleged inferiority of nonwhites.*

Past studies of online hate speech have frequently relied on dictionary-based methods, which require identifying words associated with hate speech in advance (Silva et al., 2016; Tuckwood, 2014). Other studies have incorporated sentiment analysis, natural language processing, neural networks, and other supervised and unsupervised machine learning approaches to classifying hate speech (Davidson et al., 2017; Fortuna and Nunes, 2018; Gitari et al., 2015; Kennedy et al., 2018; Olteanu et al., 2018; Siegel, 2020; Waseem and Hovy, 2016). We begin with a dictionary-based method, which enables us to classify tweets that contain the types of hate speech and white nationalist language that have received a great deal of coverage in reporting on how Trump’s rise has increased the popularity of online hate speech. In particular, we develop dictionaries of anti-Asian, anti-Black, anti-Immigrant, anti-Latinx, anti-Muslim, anti-Semitic, homophobic, and misogynistic slurs,¹² as well as a dictionary of white nationalist rhetoric using terms collected in pre-existing databases of hate speech. These include Hatebase and the Racial Slur Database, comprehensive online repositories of global hate speech (Belnik, 2017; Tuckwood, 2017), and the Anti-Defamation League’s database of slogans, terms, and symbols used by white-nationalist groups (ADL, 2017b). A list of the terms used in our

¹¹The specific criteria used to identify hate groups are debated in the literature (for a review of the literature, see Hainsworth (2000)). See also George and Wilcox (1996) who focus on political style and tactics over ideological dimensions. For an overview of defining the American “Right” and “extreme right”, see Eatwell and O’Sullivan (1989).

¹²While these categories are not exhaustive, they are the most common types of hate speech documented in the US context in pre-existing databases of online hate speech such as Hatebase and the Racial Slur Database described below.

Table 1: Hate speech and white nationalist tweet examples (warning: offensive language).

Hate speech tweets

Cant wait for donald trump to send all the monkey looking n*****s
back home to mexico

F***** Jews, calling Donald Trump “violent” while giving
n***** and Muslims a pass on their actual violence.

@realDonaldTrump yes get rid of all the b*****s, monkey looking c*****s
and muslims trump2016

White nationalist tweets

#DonaldTrump Peaceful White nationalists protect
beauty, family, and land, #AntiWhites want to destroy those things

Donald Trump means we don’t ever have to
apologize for being white ever again #NPI #AltRight

@HillaryClinton What is CHASING DOWN every last White person, assimilating
them with nonwhites and calling it 'Diversity?'. . .#whitegenocide

analysis can be found in the online appendix Section A4.¹³ It is worth noting that our hate speech categories are not mutually exclusive. One tweet can direct hate speech at several groups or individuals, as the examples of hate speech tweets from our political Twitter data set Table 1 illustrate.

However, as past research demonstrates, one of the main challenges in automatic hate-speech detection on social media is the ability to distinguish between the actual use of hate speech, posts denouncing or appropriating slurs, and speech using the terms associated with hate speech but conveying a different meaning. These methods often have low precision because they identify all messages containing particular slurs as hate speech, failing to recognize other uses and meanings of these terms (Davidson *et al.*, 2017). To evaluate the extent to which our dictionaries were accurately identifying hate speech, we used trained undergraduates and crowd-sourced coders to classify a random sample of about 25,000 tweets from our political Twitter data set

¹³We began with these pre-existing databases of hate speech and white nationalist terms. We excluded common English terms from these lists that are rarely used as hate speech and appear frequently in non-hateful contexts on Twitter. The removed terms are marked in the list in the online appendix, Section A4. We then looked at random samples of tweets from our Twitter data sets containing these terms and, where relevant, supplemented these dictionaries with additional terms that frequently co-occurred in our data. This yielded a combined dictionary of 3268 terms, which include both well-known slurs and derogatory terms often cited in reports of hate speech, as well as more obscure and rarely used terms. Because we started with this very large list of keywords from well known databases of online and offline hate speech and supplemented them with any other relevant terms, we are more concerned with false positives than recall, which prompted our use of classifiers to remove false positives from the data.

containing terms from each of our eight hate speech categories as well as white nationalist rhetoric.¹⁴ After each tweet was coded by three coders, we found that only a fraction of the tweets actually contained hate speech or white nationalist language. Many of the terms occurred in Twitter users' Twitter handles (@angry[bitch]), as part of other words ([spic]y), or homonyms (a "[chink] in his armor"). Moreover, detailed human coding on Figure Eight¹⁵ of a random sample of 5,400 tweets¹⁶ from our political Twitter data set revealed that 360 tweets, or about 7% of the 5,400 human coded tweets, were explicitly *condemning* the use of hate speech against particular groups. Examples of such tweets are displayed in Table 2.

This prevalence of anti-hate speech tweets, as well as tweets using slurs or white nationalist terms in an irrelevant manner, highlight the need to move beyond a purely dictionary-based approach to classify our tweets. Using the 25,000 total set of human coded tweets as a training data set, we trained two binary Naive Bayes classifiers, one to identify hate speech tweets and one to

Table 2: Tweets condemning hate speech and white nationalism (warning: offensive language).

RT @[HANDLE OBSCURED]: #DonaldTrump won the election & white people already don't know how to act This white boy told me I'm a N*****

Already been flicked off and called a w***** and it's only been 3 days... thanks Donald trump

Donald Trump is the type to call every east asian people as c**** c***** too. I'm not shocked <https://t.co/zFHwPLmsr9>

This just happened in Indiana. "F*** you n***** b****. Trump is going to deport you back to Africa." Day 1 of Donald

@realDonaldTrump Quite frankly Mr. Trump, you could have achieved this win without trying to get support from the # WhiteGenocide nutcases.!

¹⁴Our random sample included up to 3,000 tweets from each of our nine categories, for a total of approximately 25,000 tweets. After initially using trained undergraduate coders to ensure that tweets could be accurately classified, crowd-sourced coding was done using Figure Eight (formerly Crowdfunder), a data enrichment platform that allows a researcher to launch microtasks to a "crowd" of over five million contributors. For a recent overview of how to use Figure Eight in political science research see Benoit *et al.* (2016). Tweets were each coded by three coders. Test questions for quality control ensured that the contributors coding tweets were responding to tasks truthfully and conscientiously. If a contributor answered test questions incorrectly, that contributor was removed from the job and their data was erased.

¹⁵Coding instructions can be found in the online appendix Section A1.

¹⁶600 containing terms from each of our 8 hate speech dictionary and 600 containing white nationalist terms.

identify tweets containing white nationalist rhetoric.¹⁷ Our classifiers allowed us to identify which of the tweets containing terms from our dictionaries actually expressed hate speech and white nationalist sentiments, significantly improving the accuracy of our method.

With a method in hand for classifying tweets as containing hate speech or not, we can now return to the question of how to best measure the “prevalence” or “popularity” of hate speech on Twitter, in order to test the hypotheses laid out in Figure 1. To ensure that our findings are not driven by one particular approach, we develop four measures of the prevalence or popularity of hate speech in each of our collections on Twitter: (1) the number of tweets containing hate speech each day; (2) the number of unique users producing hate speech each day; (3) the proportion of tweets each day containing hate speech;¹⁸ and (4) the proportion of unique users producing hate speech each day.¹⁹

All four of these measures are substantively interesting and capture slightly different dimensions of the prevalence or popularity of online hate. We can imagine being concerned about the overall incidence hate speech, in which case we would want to measure total tweets as opposed to total users. However, we might also think that the more politically relevant outcome was whether more *people* started using hate speech online. Comparing counts to proportions, one could certainly argue that what is most politically relevant is the growth in quantity of hate speech in the public discourse. However, a documented increase in the raw number of hate speech tweets might not be an effective way to measure a change in the prevalence of such language if the overall volume of tweets was growing during the same period; examining the proportion of tweets containing hate speech over time insulates us against that critique.

With no *a priori* reason to favor one of these measures over the others, we run our analyses using all of them. For reasons of space, we will present analyses of the proportion of tweets containing hate speech in the main text of the paper; analyses using the remaining three measures can be found in the online appendix (Section A3). As it turns out, analyses from all four measures tell essentially the same story.

¹⁷This process is described in detail in the online appendix (Section A1). Our hatespeech classifier performed with 94% accuracy (true positives + true negatives/ total cases); 95% precision (true positives/true positives + false positives) and 90% recall (true positives/true positives + false negatives), while our white nationalist classifier performed with 97% accuracy (true positives + true negatives/ total cases); 71% precision (true positives/true positives + false positives) and 88% recall (true positives/true positives + false negatives).

¹⁸This is measured as (tweets containing hate speech each day/total tweets in each data set each day).

¹⁹This is measured as (unique users producing hate speech each day/unique users whose tweets are in our data sets).

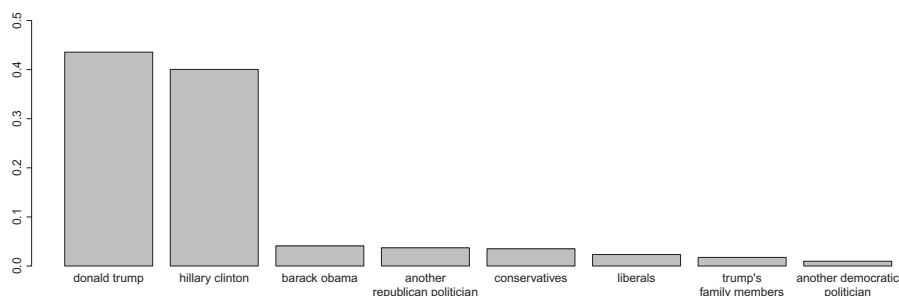


Figure 2: Proportion of human coded tweets directed at political actors. Of the 2,203 tweets classified by human coders as containing hate speech in our random sample of 5,400 tweets from the Political Twitter data sets containing dictionary terms for each type of speech, 512 (or almost one quarter of the hate speech tweets) were directed at political actors. This histogram shows the proportion of these 512 hate speech tweets directed at each type of political actor.

Before proceeding to our analysis, it is worth noting that almost one quarter of the tweets containing hate speech in our human-coded data²⁰ were directing the hate speech at political actors, especially at Hillary Clinton and Donald Trump. The breakdown of these tweets can be seen in Figure 2.

The widespread use of hate speech *against* Donald Trump suggests that although the rise in hate speech has mostly been characterized as a phenomenon that has emboldened Trump supporters, many of Trump’s opponents — on the right and left — produced online hate speech throughout the 2016 election cycle. Table 3 provides examples of these tweets containing hate speech directed at political candidates.

Empirical Strategy and Dictionary-Based Results

In order to assess whether there is empirical support for the any of the hypothesized relationships displayed in Figure 1, we employ Interrupted Time Series Analysis, a statistical method we describe in more detail below. Before presenting the results of these analyses, however, we start by simply displaying trends in the prevalence of hate speech in our datasets over time. Figure 3 shows the monthly proportion of hate speech tweets produced in the Clinton, Trump, and random sample data sets between June 17, 2015 and June 15,

²⁰While we asked human coders to code the first 5,400 hate speech tweets along several dimensions (displayed in the online appendix) in order to get a better sense of their content, we did not train a classifier to examine targets of hate speech and therefore only conducted this exploratory analysis on human coded tweets.

Table 3: Examples of hate speech tweets directed at political actors (warning: offensive language).

Donald Trump Ahead In Internal Polls & that c*** Hillary is same as Obama
n***** s***! She's a loser.

Donald trump uses self tanner... what sort of f***** s*** is that. I thought we
elected a man not a fairy!

@realDonaldTrump blow putin brown noser loser # dumpdrumpf
ivankarusianbride # tiffanytrumpussians***princess # deport them all

F*** Donald Trump let's all impeach that pathetic neo-nazi f***** b***** a***
fat f*** whale and kick his a** back to Europe feeling lucky

@realDonaldTrump called on ALL GREAT AMERICANS to UNITE to
TrumpThatB***** # Lockherup # BuildtheWall # DrainTheSwamp # MAGA!!!

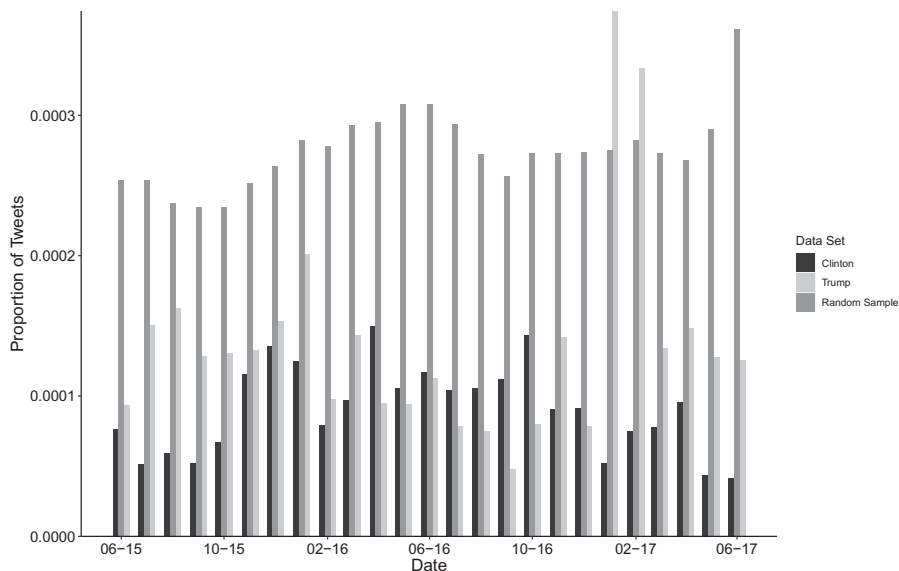


Figure 3: Monthly proportion of classified hate speech tweets in the Clinton, Trump, and random sample data sets. This figure shows the monthly proportion of classified hate-speech tweets in the Clinton, Trump, and random sample data sets. We classified tweets as hate speech (or not) using a Naive Bayes Classifier to remove false positives from tweets containing hate speech dictionary terms. Similar plots for white nationalist language, and plots displaying raw counts of the data rather than proportions, are available in the supplementary materials (Figures A13–A19).

2017.²¹ Plots showing the raw counts of tweets as well as the number and proportion of unique users producing those tweets look quite similar.²² We chose to include retweets in our primary analysis in order to capture both the number of original utterances of hate speech, as well as the popularity or spread of that hate speech, though we also disaggregate the analysis by tweets and retweets in Figures A7–A9. When we exclude retweets we see that the same general patterns persist though the spikes in our data largely disappear. We find that in any given month between 0.1% and 0.3% of tweets (including retweets) contain hate speech, a tiny fraction of both political language and general content produced by American Twitter users.²³ More importantly, Figure 3 also reveals that Trump’s victory (in November 2016) does not appear to have increased the proportion of hate speech. The Clinton data set contains less hate speech after the election, while the random sample data remains about the same for several months. The largest spike in monthly hate speech in the Trump data set occurs in late January 2017. Interestingly, while still only a fraction of a percentage point, we observe a higher proportion of hate speech in the random sample data than we do in the political data sets.

Analysis of the Trump data by hate speech type reveals that this spike is largely explained by a substantial uptick in misogynistic hate speech following the announcement of Trump’s “travel ban” executive order. This increased misogynistic language appears to be a reaction to Clinton’s decision to break her post-election silence to criticize the ban, as well as language directed at then acting Attorney General Sally Yates, who declined to defend the ban and was subsequently fired by Trump. There are also spikes in anti-Asian, anti-Muslim, and anti-Black language in this period, though their volume is much lower. These patterns disaggregated by hate speech type are displayed in Figure 4.²⁴ These plots suggest that although the data is quite bursty — for example, anti-Muslim tweets spike around terror attacks and

²¹For our statistical analysis, as well as our non-dictionary based robustness tests in the following section, we rely on the day as the primary unit of analysis for time. As Figure 3 contains data from all three collections, we pool the data by month in order to make the visualization feasible.

²²See Figures A2, A4, A6, A10–12 in the online appendix for plots displaying these other measures.

²³To be very clear, this is an empirical statement about the size of the proportion tweets containing hate speech, and *not* a normative claim that such a level of hate speech should be interpreted as problematic or not.

²⁴Comparable plots for the Clinton and Random Sample data sets, broken down by different measures of popularity, are provided in the online appendix Figures A1–A12. We also provide tables of descriptive statistics including tweets broken down by dictionary type, classification, unique users, and retweets, which generally follow very similar patterns (see Tables A1–A3).

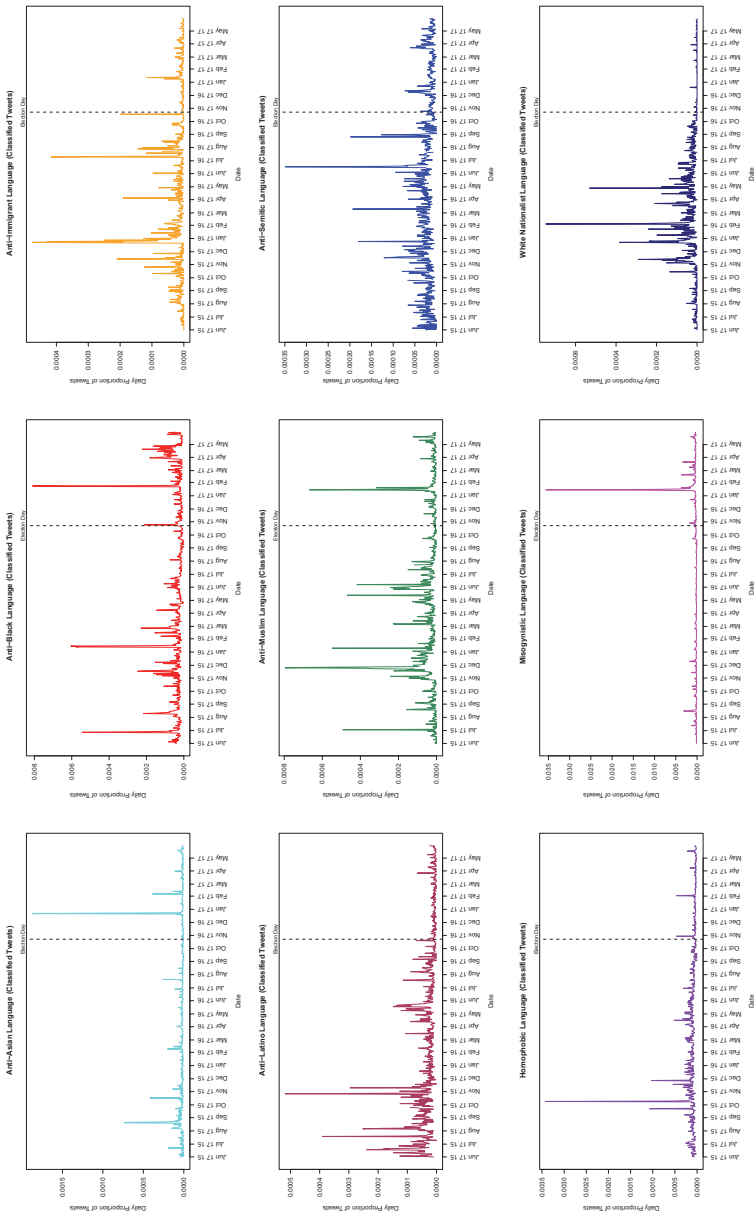


Figure 4: Daily proportion of tweets containing hatespeech and white nationalist language (Trump data set). This plot shows the daily proportion of tweets containing hate speech and white nationalist rhetoric in a data set of over 600 million tweets mentioning Donald Trump collected using Twitter's Streaming API between June 17, 2015 and June 15, 2017. Hate speech tweets were identified both using dictionaries of slurs and a Naive Bayes Classifier trained to remove false positives from our data.

anti-Semitic tweets spike following Trump’s retweet of anti-Semitic content — there is no evidence of a persistent increase in hate speech over the course of the campaign or in the aftermath of Trump’s election for any type of hate speech.

Beyond these initial insights from the raw data, our daily measures of the number and proportion of tweets containing hate speech or white nationalist rhetoric and the number and proportion of unique users tweeting them across our data sets enable us to systematically test the extent to which this language became more popular both over the course of the 2016 campaign and in the aftermath of Trump’s election. In particular, we can utilize Interrupted Time Series Analysis (ITSA) to model the popularity of hate speech on Twitter over the course of Trump’s campaign and the effect of his election as follows:

$$Y_t = \beta_0 + \beta_1(T) + \beta_2(X_t) + \beta_3(X_tT) \quad (1)$$

In Equation (1), Y_t is the proportion of tweets containing hate speech or white nationalist rhetoric at time t , T is the time since Trump announced his candidacy, X_t is a dummy variable representing Trump’s 2016 election (here the pre-election period is coded as 0 and the post-election period is coded as 1), and XT_t is an interaction term. β_0 represents the baseline proportion of hate speech tweets in each data set at $t = 0$, β_1 shows the change in the proportion of hate speech tweets associated with a one unit time increase, representing the underlying daily pre-election trend. β_2 captures the immediate effect of the election on the proportion of hate speech tweets, or an intercept change, and β_3 captures the slope change in the popularity of hate tweets following the election, relative to the pre-election trend. In other words, ITSA is a segmented regression model.²⁵ Here, we use ITSA to measure the pre-election trend, the immediate changes in the proportion of hate tweets following the election, as well as the change in the slope of the daily proportion of hate tweets after the election.²⁶ If we observe a rise in hate speech over the course of the campaign, then we should see a positive statistically significant coefficient for the pre-election trend, (β_1). If Trump’s election resulted in a lasting increase in the use of this rhetoric, then we should either see a positive shift immediately after the election (β_2) followed by a non-negative post-event slope change (β_3), or a non-negative immediate effect of the election (β_2) followed by a positive slope change in the relative popularity of hate speech in the post-election period (β_3).

²⁵Segmented regression simply refers to a model with different intercept and slope coefficients for the pre- and post-intervention time periods.

²⁶In order to address serial autocorrelation in our data, we use a first order autoregressive (AR1) model in our analysis instead of the standard OLS ITSA model (Bernal *et al.*, 2016).

We also model this relationship including quadratic terms, in case the pre- and post-election trends are not well captured by a linear model.²⁷ In these models, the squared terms tell us whether the pre- and post-election trends are concave (on the whole, slowly decreasing over time) in the case of a negative coefficient or convex (on the whole, slowly increasing over time) in the case of a positive coefficient.

$$Y_t = \beta_0 + \beta_1(T) + \beta_2(X_t) + \beta_3(X_t T) + \beta_4(T^2) + \beta_5(X_t^2) \quad (2)$$

Put another way, this allows us to test all of our hypotheses laid out in Figure 1, which involve hate speech increasing over the course of the campaign (a positive slope over time), immediately after the election (a positive discontinuity), longer term following the election (a slope change from the pre-election to the post-election period), or some combination of those three patterns.

Conducting ITSA using our political Twitter and random sample data sets, we find no evidence of a lasting increase in hate speech or white nationalist rhetoric either over the course of the campaign or in the aftermath of Trump’s election. In Figure 5, we plot the pre and post-election trends over the observed daily proportion of hate speech tweets and white nationalist language tweets in our data sets. Beginning with the Trump data set (panels (a) and (b)) — by far the largest collection — we see no significant increase in the proportion of hate speech or white nationalist language in either period. As panel (a) demonstrates, the largest spike in the proportion of hate speech in the Trump data set occurred in late January and early February of 2017, in the period surrounding the aforementioned travel ban. By contrast the largest spike in white nationalist rhetoric occurs following Trump’s retweet of a white supremacist account in February 2016.²⁸

Turning to our Clinton hate speech data (panel (c) in Figure 5), we again observe no change in the proportion of hate speech over the course of the 2016 campaign or in the aftermath of Trump’s election. In fact, we actually observe

²⁷We add quadratic terms as a robustness check since they allow us to test a broader range of null hypotheses. Rather than only testing whether there was a linear increase in hate speech, these allow us to also test more possibilities. For example, we could see an increase in hate speech over the course of the campaign that manifested both as an initial boost in hate speech after Trump declared his candidacy and then a later boost as it became clear that he would be the candidate or an initial increase followed by a decrease in the post-election period. Both of these might be evidence of a sustained Trump effect from the middle of the campaign to the end of the campaign as well as in the initial period after the election, but might not be captured with our linear trend. As in our analysis of linear trends, our models with quadratic terms do not suggest that there is any kind of sustained “Trump effect” either over the course of the campaign or following his election.

²⁸Similarly, there are no persistent increases in the number of unique users producing this content and these results hold using both linear and quadratic models. Plots and regression tables displaying these results for all data sets and outcome variables are available in the online appendix (Figures A20–A37 and Tables A7–A24).

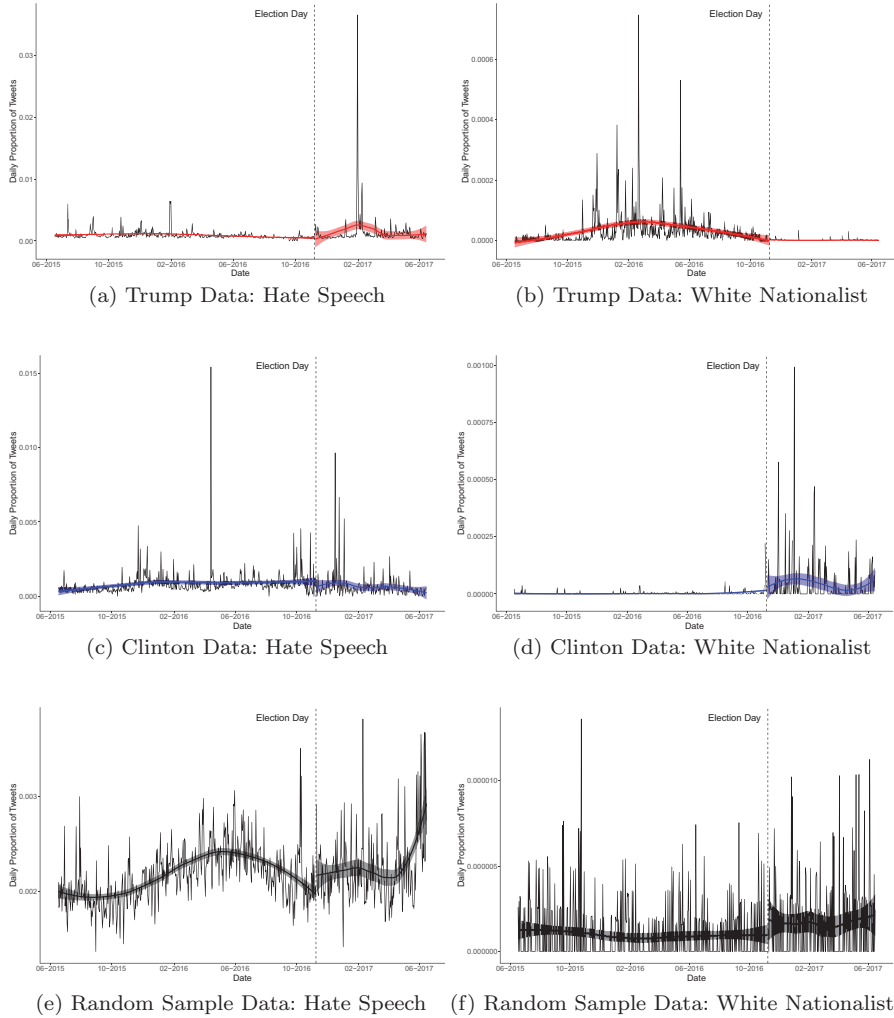


Figure 5: Effect of 2016 election on daily proportion of hate speech and white nationalist language tweets (Trump, Clinton, and random sample data sets). These plots show the pre- and post-election trends plotted with loess smoothing and 95% confidence intervals. These trend lines are plotted against the observed daily proportion of hate speech tweets and white nationalist language tweets in our data sets of over 600 million tweets referencing Donald Trump (a and b), 150 million tweets referencing Hillary Clinton (c and d), and 400 million tweets sent by a random sample of American twitter users collected using Twitter’s Streaming API between June 17, 2015 and June 15, 2017 (e and f). Hate speech and white nationalist language tweets were identified both using dictionaries of slurs and Naive Bayes classifiers trained to remove false positives from our data.

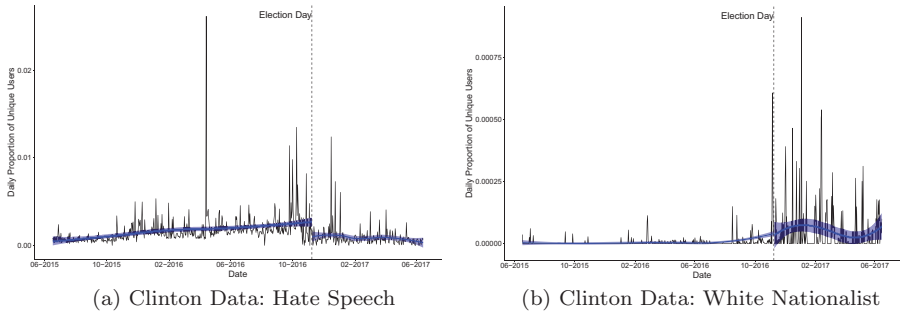


Figure 6: Effect of 2016 election on daily proportion of unique users producing hate speech and white nationalist language tweets (Clinton data set). These plots show the pre- and post-election trends plotted with loess smoothing and 95% confidence intervals. These trend lines are plotted against the observed daily proportion of unique users producing hate speech tweets and white nationalist language tweets in our data set referencing Hillary Clinton between June 17, 2015 and June 15, 2017. Hate speech and white nationalist language tweets were identified both using dictionaries of slurs and Naive Bayes classifiers trained to remove false positives from our data.

a statistically significant post-election *decrease* in the number of unique users producing hate speech, displayed in Figure 6. We do, however, observe a statistically significant increase in the proportion of unique users tweeting hate speech over the pre-election period in the Clinton data.²⁹ Figure 6 displays this trend in the daily proportion of unique users in the Clinton data. This uptick is primarily driven by the increase in misogynistic rhetoric over the course of the Clinton campaign, particularly a spike following her April 2016 debate against Bernie Sanders and a general increase as the election approached.

While we see no evidence of increasing white nationalist rhetoric over the course of the campaign in the Clinton data set (panel (d) in Figure 5), we do see an increase in the proportion of tweets containing white nationalist rhetoric — and the proportion of unique users tweeting them — after Trump’s election.³⁰ However, this coincides with a much larger decrease in the volume of white nationalist rhetoric in the Trump data set. Thus if we look at changes in the proportion of white nationalist tweets in political discourse referencing both major candidates, we see an overall decrease despite the increase we observe in the Clinton data. This pattern is even clearer examining the monthly volume of white nationalist tweets displayed in Figure A14 in the appendix, where we see a monthly average of between 1000 and 2000 white nationalist tweets

²⁹The regression table displaying this result and plots of the ITSA trend lines both aggregated and disaggregated by hate speech type are provided in Figures A20–A25 and Tables A7–A12 in the online appendix.

³⁰See Figures A21–A24 and Tables A7–A12 in the online appendix.

in the Trump data set in the pre-election period and a monthly average that hovers around 100 tweets in the Clinton data set in the post-election period.

In order to address the possibility that our *political* Twitter data differ systematically from the US Twittersphere as a whole, we replicate our analysis using a data set of over 400 million tweets produced by a random sample of 500,000 American Twitter users. Consistent with our results on political Twitter, panel (e) of Figure 5 shows no lasting increase in online hate speech over the course of Trump’s campaign. Furthermore, there is no statistically significant increase in the proportion of tweets containing hate speech following the election — or the number of unique users tweeting them — although we do see a brief one-day spike in the proportion of unique users tweeting hate speech on election day. Once again, this spike in hate speech is largely driven by misogynistic language, though we see a one day spike in anti-black language as well in the random sample data.

Examining trends in white nationalist rhetoric in the random sample data set, which we plot in panel (f) of Figure 5, we see no increase in the proportion of tweets containing white nationalist rhetoric over the course of the campaign. While there is some evidence of a slight increase in the proportion of white nationalist rhetoric after Trump’s election, this effect is not statistically significant and appears to be driven by a slight decrease in the total volume of tweets in the collection post-election. When we examine the raw volume of tweets produced in panel (f) of Figure 5 (reported in the appendix in Figure A12) we see that the average daily number of white nationalist tweets ranges from 0 to 6, averaging about 2 tweets per day across the entire period and there is clearly no visible jump in the raw data. Taken together, while we do observe an increase in white nationalist rhetoric in the Clinton data set post-election, it coincides with a much larger drop in white nationalist rhetoric in the Trump data set, resulting in a decrease in white nationalist rhetoric overall in our political collections. Taken together, these results therefore do not provide support for the proposition that Trump’s election prompted an immediate increased prevalence or popularity of online hate on Twitter along any of the ways proposed in Figure 1.³¹

Robustness Check: Reference-Text Based Analysis

One of the potential pitfalls of relying on dictionary-based methods for identifying hate speech — no matter how sophisticated the application of these approaches — is that they force the analyst to rely on a corpus of words used

³¹This conclusion also holds when we use any of the alternative measures of prevalence describing previously and reported upon in the online appendix (Section A3).

in the past to code speech in the present. In most contexts, this is unlikely to be problematic, due to the long life span of slurs and derogatory language. However, given our surprising finding in the previous section that hate speech and white nationalist rhetoric did not increase consistently either over the course of the 2016 election campaign or in its aftermath, we must seriously consider the fact that we have somehow failed to identify a significant subset of hate speech on Twitter.

To give an example, Nikhil Sonnad, writing at *Quartz*, detailed the existence of alt-right code words online, whereby “‘googles’ is an anti-Black slur; ‘skypes’ is an anti-Semitic slur; and ‘yahoos’ is an anti-Latinx slur (Sonnad, 2016). A dictionary-based method that did not contain these code words would therefore be missing the occurrence of hate speech. Perhaps even more problematically, if an event both led to an increase in hate speech *and* to the use of new code words for hate speech, we would completely miss the impact of this event.”³²

With this concern in mind, we repeat our analyses using a semi-supervised text-embedding method of measuring online hate speech and white nationalist rhetoric. The idea underlying this alternative method is to find an example of “hate speech in the wild,” or a large collection of text that contains the types of hate speech people actually use online. For this task, we rely on publicly available comments posted on Reddit.com for our reference corpus. Reddit.com is a popular news aggregation and discussion website, and Reddit entries are organized into forums with specific topics of interest (“subreddits”). Some of these subreddits are infamous for their explicitly racist, hateful and extreme alt-right content. Examples of these subreddits include /r/Coontown, /r/WhiteRights, /r/AntiPOZi, /r/european. Many of them were eventually banned or quarantined by the Reddit administration, but the comments that were posted in these subreddits are still available for analysis.

To the extent that these subreddits contain “real world” hate speech, we can measure how much hate speech is present on Twitter by examining the degree to which our tweets resemble the text in these subreddits. Now there will of course be differences in how speech is used across different platforms, so any absolute measure of similarity is going to be difficult to interpret. However, *relative* measures — such as whether tweets after Trump’s election

³²There is also a larger question as to whether hate speech that uses sanitized words ought to be considered hate speech at all. Addressing this issue is beyond the scope of the current paper, although it seems that if everyone knows what is being referred to by “kill all the skypes” then it is difficult to see why this would not be considered hate speech. For now, though, we simply note that the alternative approach described in the remainder of this section is flexible enough to pick up any instances of hate speech that might be missed by dictionary methods, regardless of whether they rely on hateful words that are simply not in the dictionaries we are using or if they represent some sort of new vocabulary — in code or not — for depicting hateful speech.

were more similar to racist subreddit content than tweets at the beginning of the campaign — can help us assess whether or not hate speech increased on Twitter over the course of the 2016 campaign or after Trump’s election. Indeed, we can examine tweets produced each day, just as we did using our machine-learning-augmented dictionary-based methods, but this time comparing not the proportion of classified tweets containing slurs, but rather how much the tweets on a given day resemble the text from hateful or white nationalist subreddits. We can then see if that similarity increases over time and/or shifts significantly following the election.³³

Not only does Reddit provide us with naturally annotated hateful text because the site is organized into “subreddits” that are explicitly devoted to particular topics — in this case communities that are infamous for their explicitly racist, hateful, and extreme alt-right content — but Reddit users frequently up-vote or down-vote posts. By excluding posts that have net negative votes from the data used to train our classifier, we subject our text to two forms of annotation: whether it is posted on a particular subreddit in the first place; and whether users of that subreddit think it belongs there. Because we are not interested in whether the tweets in our collections have a greater resemblance to one type of alt-right subreddit or another, we use a two-step classification process. We first train our classifier to group similar subreddits, and then measure the daily probability that the language in our Twitter collections shares features with the language on a collection of alt-right subreddits.³⁴ The advantage here is that unlike in our first method, we do not need to explicitly provide a dictionary of alt-right terms and phrases. Instead, our model can automatically learn relevant terms from the corpus of subreddit comments.

More specifically, we classify our tweets based on the probability that their text might have been generated by the same generative process as the text found in the hateful (or non-hateful) subreddits. Generally speaking, any kind of supervised classifier can be used to apply this method. Here we use a supervised version of the fastText model (Joulin *et al.*, 2016), which is conceptually similar to the skip-gram version of word2vec. Unlike word2vec though, instead of learning word representations in an unsupervised fashion, fastText updates these embeddings in a way that optimizes for a particular text classification task. In the current case, it is possible to train the model to predict in which subreddit (or group of subreddits) each of the comments in the corpus was posted. As a result, the model will learn semantic similarities of comments in each subreddit. After training the model to predict which subreddit (or

³³Moreover, we could also compare whether speech on Twitter in a particular time period is more similar to speech in a anti-minority subreddit than, for example, a pro-minority subreddit, a generic political discussion subreddit, or alternatively a subreddit having nothing to do with politics or identity at all (e.g., a baseball subreddit).

³⁴A detailed description of this method is available in the online appendix (Section A2).

group of subreddits) comments were posted in, it can be used to calculate class probabilities for each tweet in our collection. Finally, average daily probabilities can be calculated for each class, and changes in these probabilities can show us the dynamics of the relative popularity of a particular type of language — in this case hate speech or white nationalist rhetoric — on Twitter.

Consistent with our dictionary-based analysis, we do not observe the language in any of our three Twitter collections becoming more similar to content produced on alt-right subreddits over the course of the campaign. Trump’s election also has no effect on these probabilities. These over time findings are displayed in Figure 7.³⁵

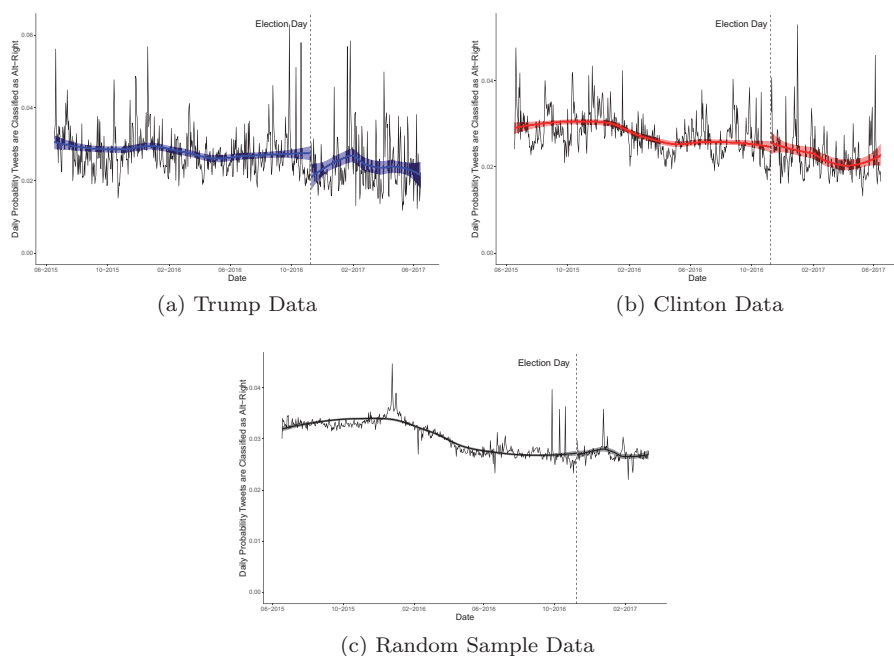


Figure 7: Effect of 2016 election on probability that tweets are classified as alt-right (Trump, Clinton, and random sample data sets). This plot shows the pre- and post-election trend lines with loess smoothing and 95% confidence intervals. These are plotted against the average daily predicted probabilities that tweets in the Trump (a), Clinton (b), and random sample (c) data sets collected using Twitter’s Streaming API between June 17, 2015 and June 15, 2017 are classified as belonging to alt-right subreddits.

³⁵Regression tables displaying results of our Interrupted Time Series Analysis (ITSA) are provided in the online appendix, Tables A28, A29, and A30.

Conclusions and Steps for Future Research

Contrary to the prevailing narrative that Trump's divisive 2016 campaign and election drove a rise in the popularity of online hate speech, we find little empirical evidence of a persistent "Trump effect" in political tweets or in a large random sample of American Twitter users during the 2016 US presidential campaign or in the six months following Trump's election. Both our dictionary and non-dictionary-based methods reveal no evidence of an increase in hate speech before or after the election across our data sets. By highlighting the shortcomings of the conventional wisdom regarding the rise of hate speech on Twitter over this period, our paper demonstrates the importance of moving beyond short term or small scale data sets when studying online speech.

Precisely because social media platforms like Twitter are so large and diverse, it is easy to find examples — even on a relatively large-scale — of almost any conceivable attitude or behavior. However, observing a particular kind of discourse online in a given moment does not necessarily mean it has increased or changed over time. This is particularly true given the bursty nature of Twitter data, where topics can trend briefly in response to events and then re-equilibrate shortly afterwards.³⁶ There is certainly evidence of tens of thousands of tweets containing hate speech and white nationalist rhetoric on Twitter over the course of the 2016 US presidential election campaign and in its aftermath. The existence of these hateful tweets thus makes it possible to use snapshots of this data as "evidence" of a "Trump effect." Further, the fact that there are indeed spikes in these data, as we have illustrated above, even makes it possible to find increases in hate speech over short time periods. Nevertheless, when we zoom out and examine the overall and relative popularity of this language at scale over a significant time period, we see that such content did not become more common in political comments on Twitter or among American Twitter users in general.

Our analysis also offers a number of important innovations for the study of online hate speech and online behavior more broadly that we hope will be adapted by scholars and practitioners alike. First, we employ two different but equally informative data sets: a collection of all tweets referencing the two candidates in the 2016 election, and a random sample of 500,000 American Twitter users. This allows us to study online hate speech both in an explicitly Trump-related political context — where we might most likely expect to see

³⁶Tables of dates with the highest volume of hate speech and white nationalist rhetoric are displayed in the online appendix, Tables A4–A6.

a “Trump effect”³⁷ — and in a representative sample of American Twitter users. Second, by using both a machine-learning augmented dictionary-based analysis and a non-dictionary approach leveraging data from subreddits to classify hate speech, we avoid drawing conclusions that are biased by one particular classification strategy or definition of speech.³⁸ Finally, by exploring changes in the volume and proportion of both hateful tweets and unique users producing hateful tweets (both including and excluding retweets), we ensure that our findings are robust to a broad set of measures of behavioral trends.

To be clear about the scope of our findings, the application of this method in the present manuscript is limited to Twitter data. While Twitter is of course not the only platform on which hate speech may have proliferated during the election period, our approach enables us to test whether people on a large, popular social media platform — indeed the president’s preferred platform — were likely to be incidentally exposed to hate speech, rather than seeking it out on specialized platforms such as Gab, Voat, Parler, or particular communities on Reddit. While recent studies have begun to investigate the spread of this language on such alternative platforms (Nithyanand *et al.*, 2017; Zannettou *et al.*, 2018), this is beyond the scope of our research on the popularization of online hate on Twitter. Indeed, we hope future research will explore the extent to which banning accounts may have pushed hateful language off of

³⁷ Given how vitriolic the campaign was, the fact that Trump explicitly retweeted white nationalists and produced derogatory tweets, and that white nationalists and other extremists who seek to broadcast their messages tend to use popular hashtags and tag well known individuals with large followings in their tweets to increase visibility, we believe that the political data sets we examine are a particularly likely place to observe a Trump effect if one existed. Moreover the 2016 election was undoubtedly a significant portion of political discourse among American Twitter users in the 18 months surrounding the 2016 election and therefore represents a large snapshot of political discourse overall, though of course not the full universe of political tweets. Perhaps the more conservative test — though the approach that best tests the popular narrative that Trump’s rise “mainstreamed” online hate in the American Twittersphere — is our use of the random sample of American Twitter users, where we again find no increase in this rhetoric either over the course of the campaign or in the aftermath of Trump’s election.

³⁸ While our definition of hate speech used to identify tweets with our dictionary-based approach is conservative in that it requires tweets to contain a slur or derogatory term, this is precisely the type of discourse that was purportedly increasing on Twitter during the time period under study. Moreover, by using our reddit-based approach to compare tweets in our Twitter data sets to language used broadly on alt-right platforms, this analysis captures a wide range of language beyond slurs and derogatory terms and, we argue, provides a fairly expansive measure of hate speech and white nationalist language.

mainstream platforms and onto more specialized ones, as well as how everyday Internet users encounter and interact with this content on other platforms.³⁹

Trolling and harassment of journalists on Twitter were frequently reported over the course of the election campaign and may have contributed to the perception of increased widespread online hate in this period. Our approach to measuring broad trends in online hate speech over time does not necessarily allow us to capture these specific incidents if they did not include references to Trump or Clinton or were not perpetrated by users in our random sample of American Twitter users. Thus it is possible that hateful attacks on individuals could have increased over the time period we analyzed, even while hate speech was not increasing generally on Twitter or in discussions of the elections. However, this too would need to be carefully studied, as hateful attacks on individuals on Twitter were taking place before the summer of 2015 as well (Parkin, 2014). Moreover, our analysis of Twitter data tells us nothing about trends in hate crimes, bias incidents, or other offline events that have also contributed to the popular narrative of a “Trump effect” and deserve further study (Müller and Schwarz, 2018; Rushin and Edwards, 2018).

The perception of the *volume* of hate speech may also have changed because the *effect* of hate speech changed over this period. When the space of possible political options is expanded to include deportation of an ethnic group, or the ban on immigration of a religious group, then people may feel the sting of hate speech more than they would in more tolerant times. And in a period of increased reporting of actual hate-crimes, hate speech may have a bigger impact on people who see it as it is associated with a greater threat. But such *impacts* of hate speech are distinct from the quantity of hate speech. And while we may ultimately care more about the impact of hate speech than the production of hate speech, in order to understand the impact of hate speech we need to accurately measure the quantity of it being circulated. Confusing the two is not likely to lead to optimal policy interventions or informed civil discourse.

Additionally, we can of course say nothing about the potentially chilling effects of Trump’s political rise on speech of either journalists or ordinary citizens, who might have believed that forays into political discussions were likely to be met with vicious and hateful personal attacks. Our manuscript therefore should not be read as evidence that there were not negative consequences from hateful speech — or hateful acts — during and after the 2016 US election campaign. Nevertheless, the fact that some of these potential negative consequences might have been driven by the *perception* that hateful speech was on the rise itself points to the importance of moving beyond anecdotal reports on hateful speech to rigorous empirical studies such as those presented here.

³⁹Of course, it would have been extremely useful to be able to replicate our study using data from Facebook — a platform used by many more Americans than Twitter — but such data are not currently available for scholarly analysis.

We hope further research will systematically explore when elites inspire changes in language, as well as when this manifests in changes in offline behavior.

Finally, our analysis only runs through June of 2017, so no conclusions should be drawn from these data concerning developments after that time. It is of course possible that events such as the Unite the Right rally, which took place in August 2017, altered the dynamics we observed over the two years prior to that date, although at the very least our study should caution against assuming this to be the case without a rigorous analysis of the relevant data.

Despite a growing body of research defining and detecting online hate, the existing scientific literature lacks a systematic approach for tracking the prevalence of this harmful speech over time (Fortuna and Nunes, 2018; Gagliardone *et al.*, 2016; Olteanu *et al.*, 2018). Although almost no empirical work has explicitly measured the overall prevalence or temporal dynamics of harmful speech on popular social media sites (Olteanu *et al.*, 2018), governments and online platforms have increasingly proposed and adopted policy interventions to combat online hate speech (Gagliardone *et al.*, 2016; House of Commons Digital, Culture, Media and Sport Committee, 2019; Marwick, 2017; Rainie *et al.*, 2017). By introducing a new systematic approach to studying the over-time dynamics of hate speech on widely used platforms like Twitter, our work offers a valuable contribution to the study of online hate speech and, in so doing, can hopefully inform policy debates.

Finally, our approach could be applied to the study of trends in many other types of online discussion and behavior beyond hate speech and white nationalist rhetoric. Finding consistent results across two different data sets, employing two different means of measuring hate speech, and using several different measures of popularity substantially increases our confidence that we are drawing meaningful inferences about behavior on Twitter over time. Our hope is that by bringing new tools and data sources to the study of online hate speech — and other online discourse — such work will enable academics, policymakers, and everyday citizens alike to better understand and address divisive social and political forces currently at play in the United States and countries around the world.

References

- ADL. 2017a. “Anti-Semitic Targeting of Journalists During the 2016 Presidential Campaign”. Available at: <https://www.adl.org/sites/default/files/documents/>.
- ADL. 2017b. “Database of Hate Symbols and Terms”. *Anti-Defamation League*. Available at: <https://www.adl.org/education/references/hate-symbols>.

- Barbaro, M. 2015. "Pithy, Mean and Powerful: How Donald Trump Mastered Twitter for 2016". *The New York Times*. Available at: <https://www.nytimes.com/2015/10/06/us/politics/donald-trump-twitter-use-campaign-2016.html>.
- Barkun, M. 2017. "President Trump and the 'Fringe'". *Terrorism and Political Violence* 29(3): 437–43.
- Belnik, S. 2017. "Racial Slur Database". Available at: <http://www.rsdb.org/>.
- Benoit, K., D. Conway, B. E. Lauderdale, M. Laver, and S. Mikhaylov. 2016. "Crowd-sourced Text Analysis: Reproducible and Agile Production of Political Data". *American Political Science Review* 110(2): 278–95.
- Bernal, J. L., S. Cummins, and A. Gasparrini. 2016. "Interrupted Time Series Regression for the Evaluation of Public Health Interventions: A Tutorial". *International Journal of Epidemiology*.
- Brader, T., J. A. Tucker, and D. Duell. 2013. "Which Parties can Lead Opinion? Experimental Evidence on Partisan Cue Taking in Multiparty Democracies". *Comparative Political Studies* 46(11): 1485–517.
- Camp, K. 2016. "Donald Trump and the Escalation of Hate". Available at: <http://billmoyers.com/story/donald-trump-escalation-hate/>.
- Cohen, C. 2016. "Donald Trump Sexism Tracker: Every Offensive Comment in One Place". Available at: <http://www.telegraph.co.uk/women/politics/donald-trump-sexism-tracker-every-offensive-comment-in-one-place/>.
- Cohen-Almagor, R. 2011. "Fighting Hate and Bigotry on the Internet". *Policy & Internet* 3(3): 1–26.
- Davidson, T., D. Warmesley, M. Macy, and I. Weber. 2017. "Automated Hate Speech Detection and the Problem of Offensive Language". Available at: <https://arxiv.org/pdf/1703.04009.pdf>.
- Duncan, P. K. 2017. "The Uses of Hate: On Hate as a Political Category". *M/C Journal* 20(1).
- Eatwell, R. and N. O'Sullivan. 1989. *The Nature of the Right: European and American Politics and Political Thought Since 1789*. Burns & Oates.
- Filimon, L. M. 2016. "Politics and Magical Thinking: How Falsehoods, Showmanship and Hawkishness Became Trademarks of Republican Presidential Electoral Campaign". *British and American Studies* (22): 211–26.
- Financial Times. 2017. "Germany Applies the Brake to Online Hate Speech". *Financial Times*. Available at: <https://www.ft.com/content/75a8c190-08c1-11e7-97d1-5e720a26771b>.
- Fording, R. C. 2014. "The Political Origins of Extremism: Minority Descriptive Representation and the Mobilization of American Hate Groups". SSRN Scholarly Paper ID 3116303, Social Science Research Network, Rochester, NY.
- Fortuna, P. and S. Nunes. 2018. "A Survey on Automatic Detection of Hate Speech in Text". *ACM Computing Surveys (CSUR)* 51(4): 85.

- Fox, J. and K. M. Warber. 2014. "Queer Identity Management and Political Self-expression on Social Networking Sites: A Co-cultural Approach to the Spiral of Silence". *Journal of Communication* 65(1): 79–100.
- Gagliardone, I., M. Pohjonen, Z. Beyene, A. Zerai, G. Aynekulu, M. Bekalu, J. Bright, M. Moges, M. Seifu, N. Stremlau, *et al.* 2016. "Mechachal: Online Debates and Elections in Ethiopia-From Hate Speech to Engagement in Social Media".
- Gagliardone, I., A. Patel, and M. Pohjonen. 2014. "Mapping and Analysing Hate Speech Online: Opportunities and Challenges for Ethiopia".
- George, J. and L. M. Wilcox. 1996. *American Extremists: Militias, Supremacists, Klansmen, Communists & Others*. Prometheus Books Amherst, NY.
- Gitari, N. D., Z. Zuping, H. Damien, and J. Long. 2015. "A Lexicon-based Approach for Hate Speech Detection". *International Journal of Multimedia and Ubiquitous Engineering* 10(4): 215–30.
- Giugni, M., R. Koopmans, F. Passy, and P. Statham. 2005. "Institutional and Discursive Opportunities for Extreme-right Mobilization in Five Countries". *Mobilization: An International Quarterly* 10(1): 145–62.
- Glynn, C. J., A. F. Hayes, and J. Shanahan. 1997. "Perceived Support for One's Opinions and Willingness to Speak Out: A Meta-Analysis of Survey Studies on the "Spiral of Silence"". *Public Opinion Quarterly*: 452–63.
- Guynn, J. 2016. "'Massive Rise' in Hate Speech on Twitter During Presidential Election". *USA Today*. Available at: <http://www.usatoday.com/story/tech/news/2016/10/21/massive-rise-in-hate-speech-twitter-during-presidential-election-donald-trump/92486210/>.
- Hainsworth, P. 2000. "The Front National: From Ascendancy to Fragmentation on the French Extreme Right". *The Politics of the Extreme Right: From the Margins to the Mainstream*: 18–32.
- House of Commons Digital, Culture, Media and Sport Committee. 2019. "Disinformation and Fake News: Final Report". *Eighth Report of Session 2017–19*. Available at: <https://publications.parliament.uk/pa/cm201719/cmselect/cmcumeds/1791/1791.pdf>.
- Huber, L. P. 2016. "Make America Great again: Donald Trump, Racist Nativism and the Virulent Adherence to White Supremacy Amid US Demographic Change". *Charleston Law Review* 10: 215.
- Joulin, A., E. Grave, P. Bojanowski, and T. Mikolov. 2016. "Bag of Tricks for Efficient Text Classification". Available at: <http://arxiv.org/abs/1607.01759>.
- Kaplan, J. 2000. *Encyclopedia of White Power: A Sourcebook on the Radical Racist Right*. Rowman & Littlefield.
- Kennedy, B., D. Kogon, K. Coombs, J. Hoover, C. Park, G. Portillo-Wightman, A. Mostafazadeh, M. Atari, and M. Dehghani. 2018. "A Typology and Coding Manual for the Study of Hate-based Rhetoric". Available at: <https://psyarxiv.com/hqjxn/>.

- Koopmans, R. and J. Muis. 2009. "The Rise of Right-wing Populist Pim Fortuyn in the Netherlands: A Discursive Opportunity Approach". *European Journal of Political Research* 48(5): 642–64.
- Koopmans, R. and S. Olzak. 2004. "Discursive Opportunities and the Evolution of Right-wing Violence in Germany". *American Journal of Sociology* 110(1): 198–230.
- Lee, M. J. and J. W. Chun. 2016. "Reading Others Comments and Public Opinion Poll Results on Social Media: Social Judgment and Spiral of Empowerment". *Computers in Human Behavior* 65: 479–87.
- Marwick, A. 2017. "Are There Limits to Online Free Speech?" Available at: <https://points.datasociety.net/are-there-limits-to-online-free-speech-14dbb7069aec>.
- Milligan, S. 2017. "A Safe Space for Hate". Available at: <https://www.usnews.com/news/the-report/articles/2017-03-24/donald-trump-and-the-politics-of-hate>.
- Montagne, R. 2016. "To Make 'Hate Rising,' Jorge Ramos Spent Time With Hate Groups". Available at: <http://www.npr.org/2016/10/21/498804694/to-make-hate-rising-jorge-ramos-spent-time-with-hate-groups>.
- Müller, K. and C. Schwarz. 2018. "Making America Hate Again? Twitter and Hate Crime under Trump". Available at: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3149103.
- Nithyanand, R., B. Schaffner, and P. Gill. 2017. "Online Political Discourse in the Trump Era". Available at: <https://arxiv.org/abs/1711.05303>.
- Noelle-Neumann, E. 1974. "The Spiral of Silence a Theory of Public Opinion". *Journal of Communication* 24(2): 43–51.
- Olteanu, A., C. Castillo, J. Boy, and K. R. Varshney. 2018. "The Effect of Extremist Violence on Hateful Speech Online". *arXiv preprint arXiv:1804.05704*.
- Ott, B. L. 2017. "The Age of Twitter: Donald J. Trump and the Politics of Debasement". *Critical Studies in Media Communication* 34(1): 59–68.
- Pang, N., S. S. Ho, A. M. R. Zhang, J. S. W. Ko, W. X. Low, and K. S. Y. Tan. 2016. "Can Spiral of Silence and Civility Predict Click Speech on Facebook?" *Computers in Human Behavior* 64: 898–905.
- Parkin, S. 2014. "Gamergate: A Scandal Erupts in the Video-game Community". *The New Yorker* 17.
- Rainie, L., J. Anderson, and J. Albright. 2017. "The Future of Free Speech, Trolls, Anonymity and Fake News Online". Available at: <http://www.pewinternet.org/2017/03/29/the-future-of-free-speech-trolls-anonymity-and-fake-news-online/>.
- Rapaport, H. 2016. "Big Data: Communicating Outside the Medium of Meaning". *Symploke* 24(1): 447–57.

- Rushin, S. and G. S. Edwards. 2018. "The Effect of President Trump's Election on Hate Crimes". Available at: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3102652.
- Scheufele, D. A. and P. Moy. 2000. "Twenty-five Years of the Spiral of Silence: A Conceptual Review and Empirical Outlook". *International Journal of Public Opinion Research* 12(1): 3–28.
- Sherrick, B. and J. Hoewe. 2018. "The Effect of Explicit Online Comment Moderation on Three Spiral of Silence Outcomes". *New Media & Society* 20(2): 453–74.
- Siegel, A. and V. Badaan. 2020. "#No2Sectarianism: Experimental Approaches to Reducing Sectarian Hate Speech Online". *American Political Science Review* 114(3): 837–55.
- Siegel, A. 2015. *Sectarian Twitter Wars: Sunni-Shia Conflict and Cooperation in the Digital Age*. Vol. 20. Carnegie Endowment for International Peace.
- Siegel, A. 2020. "Online Hate Speech: An Overview of the Literature". In: *Social Media and Democracy*. Ed. N. Persily and J. A. Tucker. New York, NY: Cambridge University Press.
- Siegel, A., J. Tucker, J. Nagler, and R. Bonneau. 2018. "Socially Mediated Sectarianism". *Unpublished Manuscript*.
- Siegel, A., J. Tucker, J. Nagler, and R. Bonneau. 2020. "Tweeting Beyond Tahrir: Ideological Diversity and Political Intolerance in Egyptian Twitter Networks". *Forthcoming in World Politics*.
- Silva, L., M. Mondal, D. Correa, F. Benevenuto, and I. Weber. 2016. "Analyzing the Targets of Hate in Online Social Media". Available at: <https://arxiv.org/abs/1603.07709v1>.
- Sonnad, N. 2016. "Alt-right Trolls are Using These Code Words for Racial Slurs Online". *Quartz*. Available at: <https://qz.com/798305/alt-right-trolls-are-using-googles-yahoos-skittles-and-skypes-as-code-words-for-racial-slurs-on-twitter/>.
- SPLC. 2017. "The Year in Hate and Extremism". *Southern Poverty Law Center*. Available at: <https://www.splcenter.org/fighting-hate/intelligence-report/2017/year-hate-and-extremism>.
- Stack, L. 2016. "Trump's Victory Alarms Gay and Transgender Groups". *New York Times*. Available at: <https://www.nytimes.com/2016/11/11/us/politics/trump-victory-alarms-gay-and-transgender-groups.html>.
- Tuckwood, C. 2014. "The State of the Field: Technology for Atrocity Response". *Genocide Studies and Prevention: An International Journal* 8(3): 9.
- Tuckwood, C. 2017. "Hatebase: Online Database of Hate Speech". *The Sentinel Project*. Available at: <https://www.hatebase.org/>.
- Van Dijk, T. A. 1992. "Denying Racism: Elite Discourse and Racism". *Discourse and Society* 3(1): 87–118.

- Walker, H. 2015. "Donald Trump Just Released an Epic Statement Raging Against Mexican Immigrants and 'Disease'". Available at: <http://www.businessinsider.com/donald-trumps-epic-statement-on-mexico-2015-7>.
- Warner, W. and J. Hirschberg. 2012. "Detecting Hate Speech on the World Wide Web". In: *Proceedings of the Second Workshop on Language in Social Media*. Association for Computational Linguistics, 19–26.
- Waseem, Z. and D. Hovy. 2016. "Hateful Symbols or Hateful People? Predictive Features for Hate Speech Detection on Twitter". *Proceedings of NAACL-HLT*: 88–93.
- Wells, C., D. V. Shah, J. C. Pevehouse, J. Yang, A. Pelled, F. Boehm, J. Lukito, S. Ghosh, and J. L. Schmidt. 2016. "How Trump Drove Coverage to the Nomination: Hybrid Media Campaigning". *Political Communication* 33(4): 669–76.
- Zaller, J. 1992. *The Nature and Origins of Mass Opinion*. New York, NY: Cambridge University Press.
- Zaller, J. 1994. "Elite Leadership of Mass Opinion". In: *Taken by Storm: The Media, Public Opinion, and US Foreign Policy in the Gulf War*. Ed. W. L. Bennett and D. L. Paletz. Chicago, IL: University of Chicago Press, 186–209.
- Zannettou, S., B. Bradlyn, E. De Cristofaro, M. Sirivianos, G. Stringhini, H. Kwak, and J. Blackburn. 2018. "What is Gab? A Bastion of Free Speech or an Alt-Right Echo Chamber?" Available at: <https://arxiv.org/abs/1802.05287>.