

# Examining Factors Associated with Twitter Account Suspension Following the 2020 U.S. Presidential Election

Farhan Asif Chowdhury<sup>1</sup>, Dheeman Saha<sup>1</sup>, Md Rashidul Hasan<sup>1</sup>, Koustuv Saha<sup>2</sup>, Abdullah Mueen<sup>1</sup>

<sup>1</sup>University of New Mexico, <sup>2</sup>Georgia Tech  
{fasifchowdhury, dsaha, mdhasan, mueen}@unm.edu, koustuv.saha@gatech.edu

## Abstract

Online social media enables mass-level, transparent, and democratized discussion on numerous socio-political issues. Due to such openness, these platforms often endure manipulation and misinformation - leading to negative impacts. To prevent such harmful activities, platform moderators employ countermeasures to safeguard against actors violating their rules. However, the correlation between publicly outlined policies and employed action is less clear to general people.

In this work, we examine violations and subsequent moderations related to the 2020 U.S. President Election discussion on Twitter. We focus on quantifying plausible reasons for the suspension, drawing on Twitter’s rules and policies by identifying suspended users (*Case*) and comparing their activities and properties with (yet) non-suspended (*Control*) users. Using a dataset of 240M election-related tweets made by 21M unique users, we observe that *Suspended* users violate Twitter’s rules at a higher rate (statistically significant) than *Control* users across all the considered aspects - hate speech, offensiveness, spamming, and civic integrity. Moreover, through the lens of Twitter’s suspension mechanism, we qualitatively examine the targeted topics for manipulation.

## 1 Introduction

Social media platforms such as Facebook, Twitter, and Reddit have become vastly popular in public discussion related to societal, economic, and political issues (Gil de Zúñiga, Jung, and Valenzuela 2012). However, in the recent past, these platforms were heavily targeted for manipulation and spreading misinformation relating to numerous civic issues across the world (Bessi and Ferrara 2016), which is often referred to as “Computational Propaganda” (Woolley and Howard 2018). In particular, the coordinated misinformation and influence campaigns of foreign state-sponsored actors during the 2016 U.S. presidential election were highly scrutinized - which eventually led to the U.S. Congressional hearings and investigation by the U.S. Department of Justice (Congress-Hearing 2017; Mueller-Report 2019).

In the aftermath of the 2016 U.S. presidential election, social platforms announced strict and improved platform moderation policy (Facebook-Update 2017; Twitter-Update 2018). However, these platform’s moderation and suspension policies have been largely debated and have faced severe criticism from the political leaders and supporters about bias towards their opposition (Bias 2016, 2020). Although these platforms publicly outline their moderation policy, there is no third-party monitoring of their enacted moderation. Moreover, social platforms like Twitter and Facebook employ ex-

tensive safeguard mechanisms that consider various aspects of user activities (coordinated activities, impersonation, etc.) to identify malicious behavior (Twitter-Safety 2021). Therefore, analyzing these suspended users’ tweets and shared content might shed light on violators’ targeted topics.

In the context of inadequate countermeasures against manipulation during the 2016 U.S. presidential election, the 2020 election was of paramount importance for platform operators to provide a safe and democratized public discussion sphere (Twitter-Policy 2021). The impact and importance of these safeguard measures are not confined to this particular election; instead, they bear cardinal implications for future online political discussions exceeding all geopolitical boundaries. Therefore, it is requisite to assess these platforms’ moderation policy — to investigate the correlation between their policies and actions, examine for potential political biases, and make general people aware of the targeted malice topics.

In this respect, we focus on analyzing the moderation policy of popular micro-blogging site Twitter as a case-study by asking the following research questions:

- **RQ1:** What factors associate with Twitter account suspension following the 2020 U.S. President Election?
- **RQ2:** How do political ideologies associate with the suspended accounts?
- **RQ3:** What was the topic of discussion among suspended users? What type of content these users shared?

**This work.** To answer these research questions, we collect a large-scale dataset of 240M tweets made by 21M unique users over eight weeks centering the 2020 U.S. President Election. Afterward, we identify 355K suspended users who participated in this election discussion. We draw upon Twitter’s rules and policies to examine plausible suspension factors. To investigate the user activity that might lead to suspension, we adopt the “*case-control*” study design from epidemiology (Schulz and Grimes 2002). We consider the suspended users as *Case* group and sample similar number of non-suspended users as *Control* group. We devise several classification techniques to quantify suspension factors among these two groups. We infer these users’ political leaning by utilizing the political bias of the news media outlet they share. By employing a language differentiation technique, we contrast the conversational topics and shared content among *Case* and *Control* groups. Through the lens of Twitter’s suspension policy, we passively infer the targeted topics by platform violators and identify the online content platforms utilized to sway the discussion.

**Summary findings.** We find that across all suspension factors, the *Suspended* users have higher (statistically significant) violation occurrences. Coherent with prior work, we

find that *Suspended* users are short-lived, have fewer followers, and show more tweeting activity. We observe a higher presence of right-leaning users than left-leaning users among *Suspended* users. We find that *Suspended* users use more curse and derogatory words and personally-attacking and propaganda related hashtags. We also notice that these users share news content from heavily biased right-leaning news-propaganda sites. We discuss the implications and limitations of our work in the Conclusion.

## 2 Related Work

There have been several works related to Twitter suspension, most of which focused primarily on spam-related suspensions (Thomas et al. 2011; Amleshwaram et al. 2021). More recently, Le et al. studied suspended users in the context of the 2016 U.S. presidential election (Le et al. 2019), and Chowdhury et al. examined a large group of suspended users related to a large-scale Twitter purge in 2018 (Chowdhury et al. 2020). We refer readers to (Le et al. 2019; Chowdhury et al. 2020) for a more comprehensive understanding of suspension and moderation on online platforms. However, none of these works quantify specific factors associated with Twitter suspension. Additionally, political discussions and related manipulation on online platforms have been thoroughly studied previously (Ferrara 2017; Im et al. 2020), mostly related to the 2016 U.S. presidential election (Badawy, Ferrara, and Lerman 2018; Zannettou et al. 2019). These works primarily focus on characterizing malicious users and inferring their motivation and impact. Im et al.; Badawy, Ferrara, and Lerman provide an extensive overview of this line of work (Im et al. 2020; Badawy, Ferrara, and Lerman 2018). In contrast, we focus on quantifying suspension factors and examining malice topics related to the 2020 U.S. presidential election.

## 3 Data

To collect tweets related to the 2020 U.S. President Election discussion, we deployed an uninterrupted data collection framework utilizing Twitter’s streaming API to filter real-time tweets based on given keywords. Similar to prior work that collected specific theme or event-related tweets (Olteanu et al. 2014), we initialize the keyword set with manually curated election-related words and hashtags. To cover the continuously evolving election discussions and topics, we update the keyword set daily with new trending hashtags and words from previous days collection, as employed in (Abu-El-Rub and Mueen 2019; Olteanu et al. 2015). We provide a list of all the keywords in the supporting website (Website 2021). Our data collection time-period spans over around ten weeks, centering the election date - from “September 28, 2020” to “December 04, 2020”. Approximately one month after our data collection ends, on “January 01, 2021” we start probing Twitter for each of the participating users from our dataset to identify the suspended users. Twitter returns the response code 63 when requested user information for a suspended user. In this process, we identify 355,573 suspended users from roughly 21M participating users. We provide a summary descriptive statistics of our dataset in Table 1 and plot per user tweet count in Figure 1. We discuss the limitation of our data collection framework in Section 6. **Ethical Concerns.** Throughout our data collection, experiment design, and analysis process, we maintain ethical re-

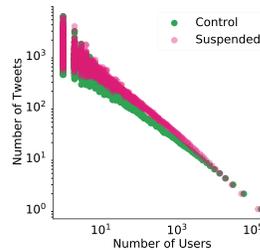


Figure 1: Distribution of # of users by # of tweets.

Statistic	Value
# Tweets	240M
# Unique users	21M
# Retweets	173M
# Quote tweets	60M
# $\mu$ tweets per user	11
# Suspended users	355K
# Tweets by sus. users	7.2M

Table 1: Descriptive statistics of Twitter Dataset.

search standards (Rivers and Lewis 2014). Hence, the accommodating academic institution’s *Institutional Review Board* exempted this project after a formal review. Following Twitter’s terms of services guideline: (1) we use the Twitter API key only for passive data collection purposes, (2) we do not publish user-specific information, (3) we do not redistribute the data, and (4) we only share aggregated statistics and derived information to facilitate future work.

## 4 Methods

### RQ1: Inferring Suspension Factors

**Twitter rules and policies.** To infer the plausible factors that explain suspension, we draw upon Twitter’s rules, and policies for free and safe public discussion (TwitterRules 2021). Twitter outlines three specific categories — (1) safety, (2) privacy, (3) authenticity, each of which entails finer sub-categories on specific violating activities. We specifically focus on five sub-categories that are more likely to be enacted upon on election discussion: three from safety — (1) hateful conduct, (2) abuse/harassment, (3) terrorism/violent extremism; two from authenticity — (4) spamming, (5) civic integrity. The rest are either not largely relevant (i.e., copyright, nudity) or inferred from our data (i.e., impersonation - as we do not have user and tweet information for all the tweet, sensitive media content - as we do not crawl media content).

**Hateful Conduct and Offensive Behavior.** Several recent works aim to identify hateful and abusive activities in online platforms, which produced publicly available datasets and trained models (Founta et al. 2018; Davidson et al. 2017). However, as the distinction between abusive and violent language is ill-defined, they unified these categories. Similarly, we combine both abusive and violent tweets into one category - *offensive*. We utilize an automatic hate speech and offensive language detection technique to detect hateful and offensive tweet, known as *HateSonar* (Davidson et al. 2017). HateSonar is a logistic regression classifier that uses several text features (i.e., TF-IDF of word-grams, sentiment, etc.), which has been trained on manually labeled tweet corpus. We use a pre-trained HateSonar model to classify each tweet into three categories: (1) hateful, (2) offensive, and (3) normal.

**Civic integrity.** Twitter has established strict rules to prevent users from “manipulating or interfering in elections or other civic processes”, including “posting and sharing misleading content” (Twitter-Integrity 2021). To infer such violation, we utilize the posted hashtags and shared news website URLs. We curate a list of hashtags related to misinformation, pro-

paganda, and conspiracy theories, borrowing from the work by (Ferrara et al. 2020), where they curated a list of conspiracy related hashtags. Additionally, we compile a list of biased and propaganda-spreader news websites based on a publicly available dataset from (FactCheck 2021; Politifact 2021). If a tweet contains a hashtag or news article from our curated list, we consider it a violation of civic integrity policy.

**Spam.** Several previous works have identified spamming on Twitter, most of which consider both tweet content and user attributes for user-level classification. Here, we primarily infer spamming violations at tweet-level based on tweet content, for which we utilize a collection of spam keywords (Benvenuto et al. 2010). However, to quantify spammers at the user-level, we also examine several account attributes (i.e., account age, tweet rate, etc.), which are most prominent for spammer detection (Thomas et al. 2011; Yang et al. 2020).

We note that the above-defined classification techniques are no match to Twitter’s actual countermeasure mechanism. Rather, we posit these methods as high-precision approaches — which utilize language models and keyword matching to avoid false-positives. The detected violations can be regarded as the lower-bound for actual ensued violations, only to increase with more comprehensive approaches.

### RQ2: Political Ideology of Suspended Users

We infer political leaning based on the political bias of the shared media outlets in the tweets. Similar to previous work on studying political ideology on Twitter (Badawy et al. 2019), we curate a list of “politically inclined” media outlets based on publicly available data from (AllSides 2021; MediaBias 2021). Additionally, if a user retweets one of the presidential candidates without adding a quote, we consider it as ideological endorsement (Ferrara et al. 2020).

### RQ3: Conversational topics and shared content

Twitter employs extensive countermeasure tools that consider a multitude of factors, features, and algorithms (Twitter-Measures 2018); which is beyond the scope of any third-party observer to reproduce. However, through the lens of Twitter’s suspension policy, we can identify platform violators’ targeted topics as a passive sensing mechanism for detecting online malice. Towards that, we contrast the conversational topics of *Suspended* and *Control* users. In particular, we consider (1) top uni-grams and bi-grams – to infer the commonality of discussion language; (2) hashtags – which are used for signaling and discoverability purposes (Bruns and Burgess 2011) and have been instrumental in several political and social movements (Arif, Stewart, and Starbird 2018). Additionally, Twitter is often used as an amplifier and fishing platform to disseminate news and multimedia content. Hence, we also examine the shared URL-domains to examine the online content platforms utilized by platform violators.

To examine the uniqueness across these dimensions, we use a generative text modeling method known as Sparse Additive Generative Models of Text, or SAGE (Eisenstein, Ahmed, and Xing 2011). We use SAGE as a text differentiation technique where each class label or latent topic is endowed with a model of the deviation in log-frequency from a constant background distribution. We utilize SAGE to identify the highly used distinctive word-grams, hashtags, and URL-domains across *Suspended* and *Control* user’s tweet corpus.

Measure	<i>Suspended</i>	<i>Control</i>	<i>d</i>	<i>t</i>	<i>KS</i>
<b>Suspension rule (Safety)</b>					
Hateful (%)	0.78	0.43	0.05	77.73***	0.003***
Offensive (%)	6.45	5.21	0.05	91.62***	0.01***
<b>Suspension rule (Authenticity)</b>					
Civic Integrity (%)	0.40	0.30	0.02	31.46***	0.001***
Spam (%)	0.56	0.38	0.03	46.83***	0.002***
<b>Account Properties</b>					
Active days	964.0	1972.9	-0.74	-151.77***	0.24***
Tweets per Day	33.5	20.82	0.25	50.20***	0.14***
Followers Count	1491.5	3337.2	-0.04	-8.54***	0.11***
Friends Count	1112.5	1263.7	-0.03	-6.74***	0.11***
<b>Political Ideology (% Tweets)</b>					
Left Leaning	3.97	5.65	-0.08	-137.9***	0.02***
Right Leaning	7.25	6.00	0.05	87.21***	0.01***

Table 2: Summary of differences in quantitative measures across *Suspended* and *Control* users. We report average occurrences across matched clusters, effect size (Cohen’s *d*), independent sample *t*-statistic, and *KS*-statistic. *p*-values are reported after Bonferroni correction (\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ ).

## 5 Results

### RQ1: Inferring Suspension Reason

In this subsection, we quantify the differences in the suspension rules between *Suspended* and *Control* users. We calculate effect size (Cohen’s *d*) and use independent sample *t*-tests to evaluate statistical significance in the differences. We perform Koglomorov-Smirnov (*KS*) test to test against the null hypothesis that the distribution of suspension rules for the *Suspended* and *Control* users are drawn from the same distribution (Saha et al. 2021). We summarize these differences in Table 2.

**Hateful Conduct and Offensive Behavior.** *Suspended* users are twice more likely to post hateful tweets than *Control* users ( $t=77.73$ ,  $p < 0.05$ ). Also, *Suspended* users post more offensive tweet than *Control* users ( $t=91.62$ ,  $p < 0.05$ ). These findings are in coherence with Twitter’s suspension policy.

**Civic Integrity and Spam.** We find that *Suspended* users are more likely to use hashtags related to conspiracy theories and share news from media sites with questionable authenticity ( $t=31.46$ ,  $p < 0.05$ ). Also, *Suspended* users are 50% more likely to post spam-tweets ( $t=46.83$ ,  $p < 0.05$ ).

**Account Properties.** We find significant difference in the active days between *Suspended* and *Control* users ( $t=-151.77$ ,  $p < 0.05$ ). The *Control* users are, on average, roughly three years older than *Suspended* users. Contrastingly, these short-lived *Suspended* users post 50% more tweets compared to *Control* users ( $t=50.20$ ,  $p < 0.05$ ). The *Suspended* users have less follower count, on average 100% less than *Control* users ( $t=-8.54$ ,  $p < 0.05$ ). However, both *Suspended* and *Control* users have similar friends count ( $t=-6.74$ ,  $p < 0.05$ ). These findings resonate with previous works studying spamming and suspension on Twitter (Thomas et al. 2011; Chowdhury et al. 2020), which identified that the rules violating users are generally short-lived, posts more tweet, and have a smaller follower base.

### RQ2: Political Ideology of Suspended Users

In Table 2, we observe 40% higher left-leaning tweets among *Control* users than *Suspended* users ( $t=-137.9$ ,  $p < 0.05$ ). In

Category	Words
Control	transition, health care, amy coney barrett, flynn, senate, sidney powell, absentee, graham, judge, legal, rigged, mail, ballot, rudy giuliani, fraudulent
Suspended	traitor, dumb, communist, biden family, idiot, seanhannity, fu*k, liar, stupid, breitbartnews, ukraine, treason, terrorist, evil, leftist

Table 3: Highly used fifteen distinctive words per user-group obtained using SAGE.

Category	Hashtags
Control	pentagon, bigtech, corruptkelly, quidproquo, doj, justicematters, climatechange, flipthesenate, michiganhearing, cnttapes
Suspended	stevebannon, bidenfamilycorruption, war-roompandemic, russia, hunterbidenemails, hunterbidenlaptop, democratsaredestroyingamerica, bidencrimesyndicate, chinajoe, chinabitchbiden

Table 4: Highly used ten distinctive hashtags per user-group obtained using SAGE.

contrast, we observe higher right-leaning tweets among *Suspended* users than *Control* users ( $t=87.21, p<0.05$ ). Our finding shows that both left-leaning and right-leaning users engaged in violating Twitter’s rules and policies, with 100% higher presence of right-leaning tweets among *Suspended* users than *Control* users.

### RQ3: Conversational topics and shared content

**Word Usage.** In Table 3, we present the top 15 most distinctive words used in per group’s tweet as obtained from the SAGE technique. We observe that *Control* users distinctly used words relating to different event-driven election-related topics (i.e., mail, ballot, rigged, fraudulent). However, we observe a large presence of swear words (i.e., idiot, dumb) and sensitive words (traitor, treason, terrorist) among *Suspended* users’ unique words, which supports the higher hate-speech detection among suspended users.

**Hashtag Usage.** Table 4 presents the top 10 most distinctively used hashtags by *Suspended* and *Control* users. Similar to word usage, the unique hashtags among *Control* users were related to specific events (i.e., cnttapes, corruptkelly, etc) and general election issues (i.e., bigtech, climatechange). In contrast, distinct hashtags from *Suspended* users were mostly related to the defamation of Democratic presidential candidate Joe Biden (i.e., bidencrimesyndicate, chinajoe, chinabitchbiden) and issues related to his son Hunter Biden (i.e., hunterbidenemails, hunterbidenlaptops).

**Shared Content.** In Table 5, we show the top 10 distinct shared domain names per user group. Among the *Control* users, we observe the presence of few moderately neutral news outlets (i.e., nytimes.com, npr.org), independent political monitoring organization (i.e., democracydocekt.com, citizenforethics.org), and few left-leaning news outlets (theatlantic.com, motherjones.com). However, among

Category	Domain
Control	democracydocekt.com, buildbackbetter.gov, www.infobae.com, latimes.com, texastribune.com, citizensforethics.org, theatlantic.com, motherjones.com, nytimes.com, npr.org
Suspended	usfuturenews.com, trumpsports.org, techfinguy.com, mostafadahroug.com, ovalofficeradio.com, wuestdevelopment.de, queenofsixteens.com, truenewshub.com, einnews.com, thefreeliberty.com

Table 5: Highly shared ten distinctive URL-domain per user-group obtained using SAGE.

*Suspended* users, we notice several heavily right-leaning non-mainstream news-propaganda sites (i.e., usfuturenews.com, ovalofficeradio.com, truenewshub.com, thefreeliberty.com).

## 6 Discussion and Conclusion

**Implications.** Our study bears an implication in shedding light on the transparency about Twitter’s content moderation policy. Although we cannot ascertain any quantitative estimation towards how far or through what means Twitter’s rules followed, our study makes insightful findings of the statistically significant occurrences of hateful, offensive, and misinformative content among the users whose accounts were suspended after a while. These findings support theoretical, empirical, and anecdotal evidence about Twitter’s moderation policies (TwitterRules 2021), which had only gained significant attention since January 2021 when Twitter suspended the U.S. President Donald Trump’s Twitter account owing to inciteful and unrest-provocative content (Trump-Ban 2020).

**Limitations and Future Work.** Our Twitter data collection has potential biases as we initialize our seed keywords manually. While investigating plausible suspension reasons, we use simple, interpretable, and high-precision approaches - which are no match to Twitter’s complex and multi-faceted safeguard mechanisms. We do not infer the exact reason for suspension for individual users; rather, we quantify violations at tweet level. Future research can use causal inference methods like matching (Saha and Sharma 2020) to minimize confounds and draw causal claims about why certain accounts were suspended. Moreover, we utilize several publicly available datasets that might suffer from biases.

We argue to situate our work as an initial step towards understanding malice, misinformation, and subsequent moderation related to the 2020 U.S. presidential election on online platforms. Our presented insights and the derived information can instigate further in-depth examination. For example, the shared news articles’ content can be analyzed to understand the nature of propaganda news. To facilitate such research, we make these news URLs publicly available and other summary statistics (Website 2021). Similarly, future works can investigate the dynamics of propaganda hashtags and news articles unique to suspended users to understand their impact and influence. Additionally, interactions among suspended users can be explored to identify potential coordination.

**Conclusion.** In this work, we perform a computational study to analyze Twitter’s suspension policy situated in the context

of the 2020 U.S. presidential election. We facilitate our work by collecting large-scale tweet dataset during the election period and subsequently identifying the suspended users. By designing a *Case-Control* experimental study and devising high-precision classification approaches, we quantify associated factors related to the suspension. Additionally, we explore the political ideology and targeted topics of suspended users. We aim to motivate more rigorous and in-depth future works through our presented insights and shared datasets.

## References

- Abu-El-Rub, N.; and Mueen, A. 2019. Botcamp: Bot-driven interactions in social campaigns. In *The World Wide Web Conference*, 2529–2535.
- AllSides. 2021. <https://www.allsides.com/media-bias/media-bias-ratings>.
- Amleshwaram, A. A.; Reddy, A. N.; Yadav, S.; Gu, G.; and Yang, C. 2021. CATS: Characterizing automation of Twitter spammers.
- Arif, A.; Stewart, L. G.; and Starbird, K. 2018. Acting the part: Examining information operations within# BlackLivesMatter discourse. *Proceedings of the ACM on Human-Computer Interaction* 2(CSCW): 1–27.
- Badawy, A.; Addawood, A.; Lerman, K.; and Ferrara, E. 2019. Characterizing the 2016 Russian IRA influence campaign. *Social Network Analysis and Mining* 9(1): 31.
- Badawy, A.; Ferrara, E.; and Lerman, K. 2018. Analyzing the digital traces of political manipulation: The 2016 Russian interference Twitter campaign. In *2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, 258–265. IEEE.
- Benevenuto, F.; Magno, G.; Rodrigues, T.; and Almeida, V. 2010. Detecting spammers on twitter.
- Bessi, A.; and Ferrara, E. 2016. Social bots distort the 2016 US Presidential election online discussion. *First Monday* 21(11-7).
- Bias. 2016. <https://www.usatoday.com/story/tech/news/2016/11/18/conservatives-accuse-twitter-of-liberal-bias/94037802/>.
- Bias. 2020. <https://www.wired.co.uk/article/twitter-political-account-ban-us-mid-term-elections>.
- Bruns, A.; and Burgess, J. E. 2011. The use of Twitter hashtags in the formation of ad hoc publics. In *Proceedings of the 6th European consortium for political research (ECPR) general conference 2011*.
- Chowdhury, F. A.; Allen, L.; Yousuf, M.; and Mueen, A. 2020. On Twitter Purge: A Retrospective Analysis of Suspended Users. In *Companion Proceedings of the Web Conference 2020*, 371–378.
- Congress-Hearing. 2017. <https://www.govinfo.gov/content/pkg/CHRG-115shrg27398/pdf/CHRG-115shrg27398.pdf>.
- Davidson, T.; Warmusley, D.; Macy, M.; and Weber, I. 2017. Automated hate speech detection and the problem of offensive language. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 11.
- Eisenstein, J.; Ahmed, A.; and Xing, E. P. 2011. Sparse additive generative models of text .
- Facebook-Update. 2017. <https://about.fb.com/news/2017/09/information-operations-update/>.
- FactCheck. 2021. <https://www.factcheck.org/2017/07/websites-post-fake-satirical-stories/>.
- Ferrara, E. 2017. Disinformation and social bot operations in the run up to the 2017 French presidential election. *arXiv preprint arXiv:1707.00086* .
- Ferrara, E.; Chang, H.; Chen, E.; Muric, G.; and Patel, J. 2020. Characterizing social media manipulation in the 2020 US presidential election. *First Monday* .
- Founta, A.; Djouvas, C.; Chatzakou, D.; Leontiadis, I.; Blackburn, J.; Stringhini, G.; Vakali, A.; Sirivianos, M.; and Kourtellis, N. 2018. Large scale crowdsourcing and characterization of twitter abusive behavior. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 12.
- Gil de Zúñiga, H.; Jung, N.; and Valenzuela, S. 2012. Social media use for news and individuals’ social capital, civic engagement and political participation. *Journal of computer-mediated communication* 17(3): 319–336.
- Im, J.; Chandrasekharan, E.; Sargent, J.; Lighthammer, P.; Denby, T.; Bhargava, A.; Hemphill, L.; Jurgens, D.; and Gilbert, E. 2020. Still out there: Modeling and identifying Russian troll accounts on twitter. In *12th ACM Conference on Web Science*, 1–10.
- Le, H.; Boynton, G.; Shafiq, Z.; and Srinivasan, P. 2019. A postmortem of suspended Twitter accounts in the 2016 US presidential election. In *2019 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, 258–265. IEEE.
- MediaBias. 2021. <https://mediabiasfactcheck.com/>.
- Mueller-Report. 2019. <https://www.justice.gov/storage/report.pdf>.
- Olteanu, A.; Castillo, C.; Diakopoulos, N.; and Aberer, K. 2015. Comparing events coverage in online news and social media: The case of climate change. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 9.
- Olteanu, A.; Castillo, C.; Diaz, F.; and Vieweg, S. 2014. Crisislex: A lexicon for collecting and filtering microblogged communications in crises. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 8.
- Politifact. 2021. <https://www.politifact.com/article/2017/apr/20/politifact-guide-fake-news-websites-and-what-they/>.
- Rivers, C. M.; and Lewis, B. L. 2014. Ethical research standards in a world of big data. *F1000Research* 3.
- Saha, K.; Liu, Y.; Vincent, N.; Chowdhury, F. A.; Neves, L.; Shah, N.; and Bos, M. W. 2021. AdverTiming Matters: Examining User Ad Consumption for Effective Ad Allocations on Social Media. In *Proc. CHI*.
- Saha, K.; and Sharma, A. 2020. Causal Factors of Effective Psychosocial Outcomes in Online Mental Health Communities. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 14, 590–601.
- Schulz, K. F.; and Grimes, D. A. 2002. Case-control studies: research in reverse. *The Lancet* 359(9304): 431–434.
- Thomas, K.; Grier, C.; Song, D.; and Paxson, V. 2011. Suspended accounts in retrospect: an analysis of twitter spam. In *Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference*, 243–258.
- Trump-Ban. 2020. [https://blog.twitter.com/en\\_us/topics/company/2020/suspension.html](https://blog.twitter.com/en_us/topics/company/2020/suspension.html).
- Twitter-Integrity. 2021. <https://help.twitter.com/en/rules-and-policies/election-integrity-policy>.
- Twitter-Measures. 2018. [https://blog.twitter.com/en\\_us/topics/company/2018/how-twitter-is-fighting-spam-and-malicious-automation.html](https://blog.twitter.com/en_us/topics/company/2018/how-twitter-is-fighting-spam-and-malicious-automation.html).
- Twitter-Policy. 2021. [https://blog.twitter.com/en\\_us/topics/company/2020/2020-election-changes.html](https://blog.twitter.com/en_us/topics/company/2020/2020-election-changes.html).
- Twitter-Safety. 2021. [https://blog.twitter.com/en\\_us/topics/company/2018/how-twitter-is-fighting-spam-and-malicious-automation.html](https://blog.twitter.com/en_us/topics/company/2018/how-twitter-is-fighting-spam-and-malicious-automation.html).
- Twitter-Update. 2018. [https://blog.twitter.com/en\\_us/topics/company/2018/2016-election-update.html](https://blog.twitter.com/en_us/topics/company/2018/2016-election-update.html).
- TwitterRules. 2021. <https://help.twitter.com/en/rules-and-policies/twitter-rules>.
- Website. 2021. <https://sites.google.com/view/us-election20-twitter-suspend>.
- Woolley, S. C.; and Howard, P. N. 2018. *Computational propaganda: political parties, politicians, and political manipulation on social media*. Oxford University Press.
- Yang, K.-C.; Varol, O.; Hui, P.-M.; and Menczer, F. 2020. Scalable and generalizable social bot detection through data selection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 1096–1103.
- Zannettou, S.; Caulfield, T.; De Cristofaro, E.; Sirivianos, M.; Stringhini, G.; and Blackburn, J. 2019. Disinformation warfare: Understanding state-sponsored trolls on Twitter and their influence on the web. In *Companion proceedings of the 2019 world wide web conference*, 218–226.