## **Viewpoint**

# Bots and Misinformation Spread on Social Media: Implications for COVID-19

McKenzie Himelein-Wachowiak<sup>1</sup>, BA; Salvatore Giorgi<sup>1,2</sup>, MSc; Amanda Devoto<sup>1</sup>, PhD; Muhammad Rahman<sup>1</sup>, PhD; Lyle Ungar<sup>2</sup>, PhD; H Andrew Schwartz<sup>3</sup>, PhD; David H Epstein<sup>1</sup>, PhD; Lorenzo Leggio<sup>1</sup>, MD, PhD; Brenda Curtis<sup>1</sup>, MSPH, PhD

## **Corresponding Author:**

Brenda Curtis, MSPH, PhD Intramural Research Program National Institute on Drug Abuse 251 Bayview Blvd Suite 200 Baltimore, MD, 21224 United States

Phone: 1 443 740 2126 Email: <u>brenda.curtis@nih.gov</u>

## **Abstract**

As of March 2021, the SARS-CoV-2 virus has been responsible for over 115 million cases of COVID-19 worldwide, resulting in over 2.5 million deaths. As the virus spread exponentially, so did its media coverage, resulting in a proliferation of conflicting information on social media platforms—a so-called "infodemic." In this viewpoint, we survey past literature investigating the role of automated accounts, or "bots," in spreading such misinformation, drawing connections to the COVID-19 pandemic. We also review strategies used by bots to spread (mis)information and examine the potential origins of bots. We conclude by conducting and presenting a secondary analysis of data sets of known bots in which we find that up to 66% of bots are discussing COVID-19. The proliferation of COVID-19 (mis)information by bots, coupled with human susceptibility to believing and sharing misinformation, may well impact the course of the pandemic.

(J Med Internet Res 2021;23(5):e26933) doi: 10.2196/26933

## **KEYWORDS**

COVID-19; coronavirus; social media; bots; infodemiology; infoveillance; social listening; infodemic; spambots; misinformation; disinformation; fake news; online communities; Twitter; public health

## Introduction

Globally, 2020 has been characterized by COVID-19, the disease caused by the SARS-CoV-2 virus. As of March 2021, the COVID-19 pandemic has been responsible for over 115 million documented cases, resulting in over 2.5 million deaths. The United States accounts for 24.9% of the world's COVID-19 cases, more than any other country [1].

As the virus spread across the United States, media coverage and information from online sources grew along with it [2]. Among Americans, 72% report using an online news source for COVID-19 information in the last week, with 47% reporting that the source was social media [3]. The number of research

articles focusing on COVID-19 has also grown exponentially; more research articles about the disease were published in the first 4 months of the COVID-19 pandemic than throughout the entirety of the severe acute respiratory syndrome (SARS) and Middle East respiratory syndrome (MERS) pandemics combined [4]. Unfortunately, this breadth, and the speed with which information can travel, sets the stage for the rapid transmission of misinformation, conspiracy theories, and "fake news" about the pandemic [5]. One study found that 33% of people in the United States report having seen "a lot" or "a great deal" of false or misleading information about the virus on social media [3]. Dr Tedros Adhanom Ghebreyesus, the Director-General of the World Health Organization, referred to this accelerated flow



<sup>&</sup>lt;sup>1</sup>Intramural Research Program, National Institute on Drug Abuse, Baltimore, MD, United States

<sup>&</sup>lt;sup>2</sup>Department of Computer and Information Science, University of Pennsylvania, Philadelphia, PA, United States

<sup>&</sup>lt;sup>3</sup>Department of Computer Science, Stony Brook University, Stony Brook, NY, United States

of information about COVID-19, much of it inaccurate, as an "infodemic" [6].

Though the pandemic is ongoing, evidence is emerging regarding COVID-19 misinformation on social media. Rumors have spread about the origin of the virus, potential treatments or protections, and the severity and prevalence of the disease. In one sample of tweets related to COVID-19, 24.8% of tweets included misinformation and 17.4% included unverifiable information [7]. The authors found no difference in engagement patterns with misinformation and verified information, suggesting that myths about the virus reach as many people on Twitter as truths. A similar study demonstrated that fully false claims about the virus propagated more rapidly and were more frequently liked than partially false claims. Tweets containing false claims also had less tentative language than valid claims [8].

This trend of misinformation emerging during times of humanitarian crises and propagating via social media platforms is not new. Previous research has documented the spread of misinformation, rumors, and conspiracies on social media in the aftermath of the 2010 Haiti earthquake [9], the 2012 Sandy Hook Elementary School shooting [10], Hurricane Sandy in 2012 [11], the 2013 Boston Marathon bombings [12,13], and the 2013 Ebola outbreak [14].

Misinformation can be spread directly by humans, as well as by automated online accounts, colloquially called "bots." Social bots, which pose as real (human) users on platforms such as Twitter, use behaviors like excessive posting, early and frequent retweeting of emerging news, and tagging or mentioning influential figures in the hope they will spread the content to their thousands of followers [15]. Bots have been found to disproportionately contribute to Twitter conversations on controversial political and public health matters, although there is less evidence they are biased toward one "side" of these issues [16-18].

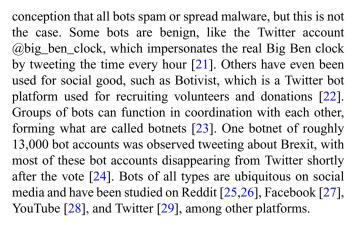
This paper combines a scoping review with an unpublished secondary analysis, similar in style to Leggio et al [19] and Zhu et al [20]. We begin with a high-level survey of the current bot literature: how bots are defined, what technical features distinguish bots, and the detection of bots using machine learning methods. We also examine how bots spread information, including misinformation, and explore the potential consequences with respect to the COVID-19 pandemic. Finally, we analyze and present the extent to which known bots are publishing COVID-19—related content.

## What Are Bots?

Before addressing issues surrounding the spread of misinformation, we provide a definition of bots, describe their typical features, and explain how detection algorithms identify bots

### **Definition and Identification**

Bots, shorthand for "software robots," come in a large variety of forms. Bots are typically automated in some fashion, either fully automated or human-in-the-loop. There is a common



Given their large variety, bots are often organized into subclasses, a selection of which we discuss here. Content polluters are one subclass; these are "accounts that disseminate malware and unsolicited content." Traditional spambots, another subclass, are "designed to be recognizable as bots" [30]. Social bots—a newer, more advanced type of bot [31-33]—use "a computer algorithm that automatically produces content and interacts with humans on social media, trying to emulate and possibly alter their behavior." There are also hybrid human-bot accounts (often called cyborgs) [34], which "exhibit human-like behavior and messages through loosely structured, generic, automated messages and from borrowed content copied from other sources" [35]. It is not always clear which category a bot may fall into (eg, if a given social bot is also a cyborg).

Various methods have been used to identify bots "in the wild," so as to build the data sets of known bots used to train bot-detection algorithms. One method, the "social honeypot" [36], mimics methods traditionally used by researchers to monitor hacker activity [37] and email harvesting [38]. Specifically, social honeypots are fake social media profiles set up with characteristics desirable to spammers, such as certain demographics, relationship statuses, and profile pictures [39]. When bots attempt to spam the honeypots (by linking malware-infested content or pushing product websites), researchers can easily identify them.

### **Technical Features of Bots**

#### **Overview**

Features that distinguish bots from humans roughly fall into three categories: (1) network properties, such as hashtags and friend/follower connections, (2) account activity and temporal patterns, and (3) profile and tweet content. These feature categories have the advantage of being applicable across different social media platforms [27].

#### Network Properties

Networks based on friend/follower connections, hashtag use, retweets, and mentions have been used in a number of studies that seek to identify social bots [40-43], exploiting network homophily (ie, humans tend to follow other humans and bots tend to follow other bots). As bots become more sophisticated, network properties become less indicative of them; studies have found groups of bots that were able to build social networks that mimic those of humans [44].



## Account Activity and Temporal Patterns

Patterns of content generation can be good markers of bots. Bots compose fewer original tweets than humans, but retweet others' tweets much more frequently, and have a shorter time interval between tweets [40]. Ferrara et al [31] found that humans are retweeted by others more than are bots, suggesting that bots may struggle to compose convincing or interesting tweets. However, many others have found this not to be the case [15,16,33]. Finally, humans typically modify their behavior during each online session; as the session progresses, the density of new tweets decreases. Bots do not engage in these "sessions" of social media usage, and accordingly do not modify their behavior [45].

## **Profile and Tweet Content**

Profile metadata such as account age and username can be used to identify social bots. Ferrara et al [31] showed that bots have shorter account age (ie, the accounts were created more

recently), as well as longer usernames. Automatic sentiment analysis of tweet content has also been studied as a means of distinguishing bots from humans. One study found humans expressed stronger positive sentiment than bots, and that humans more frequently "flip-flopped" in their sentiment [42].

#### **Detection of Bots**

Over the past decade, several teams have sought to develop algorithms that successfully identify bots online. Social media platforms use similar algorithms internally to remove accounts likely to be bots. These algorithms originated in early attempts to identify spam emails [46], social phishing [47], and other types of cybercrimes [37]. With the advent of online communities, cybercriminals turned their attention to these sites, eventually creating fake, automated accounts at scale [48]. Table 1 provides a summary of several prominent papers on bot identification. We note that the details of specific machine learning algorithms are beyond the scope of this paper and therefore are not included in this manuscript.

Table 1. Review of state-of-the-art detection of bots on Facebook and Twitter.

Type and reference	Platform	Number of accounts	Features			Model	Metric	Predictive accuracy
			$N^{a}$	$T^{b}$	$C^c$			
Human judgment (manual a	nnotation)			•		-		
Cresci et al (2017) [33]	Twitter	928				Manual annotation	F1-score	0.57
Automatic methods								
Ahmed and Abulaish (2013) [27]	Facebook and Twitter	320 (Facebook), 305 (Twitter)			1	Naïve Bayes, decision trees, rule learners	Detection rate	0.96 (Facebook), 0.99 (Twitter)
Dickerson et al (2014) [42]	Twitter	897	✓	✓	1	Gradient boosting	Area under the curve	0.73
Cresci et al (2017) [33]	Twitter	928		✓		Digital DNA sequences	F1-score	0.92
Varol et al (2017) [41]	Twitter	21,000	✓	1	1	Random forests	Area under the curve	0.95
Kudugunta and Ferrara (2018) [49]	Twitter	8386			1	AdaBoost	Area under the curve	>0.99
Mazza et al (2019) [50]	Twitter	1000		✓		Long short-term memory networks	F1-score	0.87
Santia et al (2019) [51]	Facebook	1000			✓	Support vector machines, decision trees, Naïve Bayes	F1-score	0.72
Yang et al (2020) [52]	Twitter	137,520		✓	✓	Random forests	Area under the curve	0.60-0.99

<sup>&</sup>lt;sup>a</sup>N: network properties.

The first reference in Table 1 involved a manual annotation task in which raters were asked to label a Twitter account as human or bot. The fourth study listed in the table is the same as the first study; in this study, the same data set was evaluated by both human annotators and machine learning methods [33]. There was a large discrepancy in predictive accuracy (F1-score) between the two methods: 0.57 for the human annotators versus 0.92 for the automated method. Stated another way, human

participants correctly identified social bots less than 25% of the time, though they were quite good at identifying genuine (human) accounts (92%) and traditional spambots (91%). These results suggest that social bots have a very different online presence from traditional spambots, or "content polluters"—and that this presence is convincingly human. Even if the human annotators are compared to the lowest scoring automated method (which we note is in a different domain and, thus, not directly



<sup>&</sup>lt;sup>b</sup>T: account activity and temporal patterns.

<sup>&</sup>lt;sup>c</sup>C: profile and tweet content.

comparable), the machine learning algorithm still provides a considerable boost in F1-score (0.57 versus 0.72).

There is no good way to compare all automated methods directly, as data sets are typically built in a single domain (ie, a single social media platform) and rapid advances in machine learning techniques prevent comparisons between models published even a few years apart. Furthermore, results suggest that models trained on highly curated bot data sets (eg, groups of accounts promoting certain hashtags or spamming a particular honeypot) may not perform well at detecting bots in other contexts. Yang et al [52] used a large number of publicly available bot data sets, training machine learning models on each set and testing them on those remaining. The result was a wide range of predictive accuracies across different bot data sets.

## How Do Bots Amplify and Spread Misinformation?

We adopt the definition of misinformation used by Treen and colleagues: "misleading [or false] information that is created and spread, regardless of intent to deceive" [53]. For the purposes of this paper, we include fake news and false conspiracy theories under this umbrella term.

Many features of bots likely enable them to be "super-spreaders" of misinformation. Bots have been shown to retweet articles within seconds of their first being posted, contributing to the articles going viral [15]. Moreover, the authors of this study found that 33% of the top sharers of content from low-credibility sources were likely to be bots, significantly higher than the proportion of bots among top sharers of fact-checked content. Similarly, in a study of bots and "anti-vaxxer" tweets, Yuan et al [18] found that bots were "hyper-social," disproportionately contributing to content distribution. Bots also employ the strategy of mentioning influential users, such @realDonaldTrump, in tweets linking to false or misleading articles, and are more likely to do so than their human counterparts [15]. The hope is that these users will share the article with their many followers, contributing to its spread and boosting its credibility. "Verified" (blue check) Twitter users, often celebrities, have been shown to both author and propagate COVID-19-related misinformation [54]. Interestingly, the frequency of false claims about the 2020 election dropped dramatically in the week after former president Donald Trump was removed from the platform [55].

In light of findings that humans are largely unable to distinguish social bots from genuine (human) accounts [33], it is likely that humans unknowingly contribute to the spread of misinformation as well. Accordingly, one study found that in regard to low-credibility content, humans retweet bots and other humans at the same rate [15]. Similarly, Vosoughi et al [56] found that "fake news" articles spread faster on Twitter than true news articles because humans, not bots, were more likely to retweet fake articles. Given human susceptibility to both automated accounts and "fake news," some have warned that intelligent social bots could be leveraged for mass deception or even "political overthrow" [57].

There is reason to believe that bots have already infiltrated political conversations online. Leading up to the 2016 presidential election in the United States, 20% of all political tweets originated from accounts that were likely to be bots [16]. While it did not specifically implicate bots, one study found that a majority of "fake" or extremely biased news articles relating to the 2016 election were shared by unverified accounts—that is, accounts that were not confirmed to be human [58]. There is also evidence that bots spread misinformation in the 2017 French presidential election, though ultimately the bot campaign was unsuccessful, in part because the human users who engaged with the bots were mostly foreigners with no say in the election outcome [59]. Bot strategies specifically relevant to political campaigns include "hashtag hijacking," in which bots adopt an opponent's hashtags in order to spam or otherwise undermine them, as well as flagging their opponent's legitimate content in the hopes it gets removed from the platform [60].

## Where Do Bots Come From?

The origin of social bots is a challenging question to answer. Given the aforementioned concerns of political disruption by social bots, one may assume that foreign actors create social bots to interfere with political processes. Indeed, the Mueller report found evidence of Russian interference in the 2016 US election via social media platforms [61], and Twitter reports removing over 50,000 automated accounts with connections to Russia [62]. However, locating the origin of a social media account is difficult, as tweets from these accounts are rarely geotagged. Rheault and Musulan [63] proposed a methodology to identify clusters of foreign bots used during the 2019 Canadian election using uniform manifold approximation and projection combined with user-level document embeddings. Simply put, the authors constructed communities of users via linguistic similarities, and identified members significantly outside these communities as foreign bots.

Of note, studies have shown that a majority of social bots focusing on election-related content originate domestically [63]. Reasons for a candidate or their supporters to employ social bots may be relatively benign, such as boosting follower counts or sharing news stories, or they may involve smear campaigns [64].

While the ability to investigate the origin and motive of social bots is difficult, the means to create a social bot are fairly easy to access. Social bots are available for purchase on the dark web, and there are tens of thousands of codes for building social bots on free repositories like GitHub [65]. Of note, the top contributors of bot-development tools for mainstream social media sites are the United States, the United Kingdom, and Japan. The authors of this paper also note the intelligence and capabilities of these freely available bots may be overstated.

## Are Bots Tweeting About COVID-19?

In light of the COVID-19 "infodemic" and findings that social bots have contributed to misinformation spread in critical times, we sought to assess the number of known Twitter bots producing COVID-19—related content. To this end, we gathered a number



of publicly available bot data sets from the Bot Repository [66]. These data sets include both traditional spambots and social bots that were first identified through a number of different methods (see the original papers for more details).

Using the open-source Python package TwitterMySQL [67], which interfaces with the Twitter application programming interface (API), we were able to pull all tweets from 2020 for each bot in the combined data set. Of note, Twitter's API limits access to tweets and account information available at the time of collection. Tweets and accounts that have been deleted or made private since originally appearing in one of the above papers are not made available, meaning we had less data than what was reported in the original papers. Our final data set consisted of 3.03 million tweets from 3953 bots, with an average

of 768.9 (SD 1145.4) tweets per bot, spanning January 1, 2020, to August 21, 2020.

From these data, we pulled tweets using a set of 15 COVID-19–related keywords, which have previously been used to identify COVID-19 tweets in a study tracking mental health and psychiatric symptoms over time [68]. Sample keywords include #coronavirus, #covid19, and #socialdistancing. We then counted the number of accounts that mentioned these keywords in tweets since January 2020. Table 2 shows the percentage of bot accounts in each data set currently tweeting about COVID-19. Original sample size refers to the number of bots identified in this data set, while current sample size is the number of currently active bots (ie, tweeting in 2020). Between 53% (96/182) and 66% (515/780) of these bots are actively tweeting about COVID-19.

**Table 2.** Open-source data sets of bots discussing COVID-19<sup>a</sup>.

Reference	Year	Original sample size, n	Current sample size, n	Bots discussing COVID-19, n (%)
Lee et al [36]	2011	22,223	2623	1427 (54)
Varol et al [41]	2017	826	292	164 (56)
Gilani et al [69]	2017	1130	780	515 (66)
Cresci et al [33]	2017	4912	77	48 (62)
Mazza et al [50]	2019	391	182	96 (53)

<sup>&</sup>lt;sup>a</sup>Original sample size is the number of bot IDs publicly released on the Bot Repository, while current sample size is the number of active accounts tweeting in 2020. Percentage discussing COVID-19 is the percentage of bots with at least one tweet containing a COVID-19 keyword out of those active in 2020.

## Implications for the COVID-19 Pandemic

Here we have shown that a majority of known bots are tweeting about COVID-19, a finding that corroborates similar studies [68,70]. Early in the pandemic, one study found that 45% of COVID-19—related tweets originate from bots [71], although Twitter has pushed back on this claim, citing false-positive detection algorithms [72]. Another study showed that COVID-19 misinformation on Twitter was more likely to come from unverified accounts—that is, accounts not confirmed to be human [7]. In an analysis of 43 million COVID-19—related tweets, bots were found to be pushing a number of conspiracy theories, such as QAnon, in addition to retweeting links from partisan news sites [73]. Headlines from these links often suggested that the virus was made in Wuhan laboratories or was a biological weapon.

One limitation of our study is that we did not investigate the validity of COVID-19–related claims endorsed by bots in our analyses. It may be that bots are largely retweeting mainstream news sources, as was the case in a recent study of bots using #COVID19 or #COVID-19 hashtags [68]. However, previous research has connected bots to the spread of misinformation in other public health domains, such as vaccines [30] and e-cigarettes [74], and unsubstantiated medical claims surrounding the use of marijuana [75].

Such misinformation can have detrimental consequences for the course of the COVID-19 pandemic. Examples of these real-world consequences include shortages of hydroxychloroquine (a drug that is crucial for treating lupus and malaria) due to increased demand from people who believe it will protect them from COVID-19 [76,77]. This drug has been promoted as a preventative against COVID-19 on social media, even though several randomized controlled trials have found it ineffective, [78,79], and the National Institutes of Health recently halted its own trial due to lack of effectiveness [80]. Moreover, belief in conspiracy theories about COVID-19 is associated with a decreased likelihood of engaging in protective measures such as frequent handwashing and social distancing, suggesting that misinformation may even contribute to the severity of the pandemic [81]. In addition, exposure to misinformation has been negatively correlated with intention to take a COVID-19 vaccine [82].

We are certainly not the first to express concern with viral misinformation; in May 2020, Twitter began labeling fake or misleading news related to COVID-19 in an effort to ensure the integrity of information shared on the platform [83]. Facebook introduced even more controls, such as organizing the most vetted articles at the top of the news feed, banning antimask groups, and sending antimisinformation messages to users who have shared fake news [84]. However, these measures are designed to target humans. In light of the numerous viral rumors relating to COVID-19 and the US response to the pandemic, we believe that bots likely contributed to their spread.

Major social media platforms like Twitter and Facebook do have methods to curtail suspected bots. In 2018, Twitter banned close to 70 million suspicious accounts in a matter of months



[85]. Facebook banned 1.3 billion suspicious accounts in the third quarter of 2020. The platform estimates these accounts represent 5% of its worldwide monthly active users. The vast majority of suspicious accounts were identified using automated detection methods, but 0.7% were first flagged by human users, suggesting that everyday Facebook users concerned about malicious activity on the platform can contribute to efforts to ban these accounts [86].

Mitigation of the harmful effects of social bots can also occur at the policy level. In 2018, California became the first and only state to pass a law requiring social bots to identify themselves as such [87]. In 2019, Senator Dianne Feinstein proposed a similar bill federally; the bill would allow the Federal Trade Commission to enforce bot transparency and would prohibit political candidates from incorporating social bots in their campaign strategy [88]. The United States Congress has brought top executives from Facebook, Twitter, and Google to testify before Congress about Russian influence on their platforms in advance of the 2016 election [89]. Scholars have interpreted these actions as a sign that the government wishes to maintain the right to regulate content on social media—a prospect that brings concerns of its own [90]. Presently, content problems on social media platforms are almost exclusively dealt with by the owners of those platforms, usually in response to user complaints, but in the coming years we may see an increase in government oversight on these platforms, fueling concerns about state-sponsored censorship [91,92]. More fundamentally, some have argued that, before any actionable policy or automatic interventions can be enabled, ambiguities in both bot definitions and jurisdiction and authority need to be addressed [90].

Even as citizens, social media platforms, and policy makers converge on the notion that bots and misinformation are urgent problems, the methods used to address the issue have had mixed results. When social media platforms crack down on bots and

misinformation, either through automated techniques or manual content moderation, they run the risk of censoring online speech and further disenfranchising minority populations. Content promotion and moderation can lead to arbitrary policy decisions that may be inconsistent across or even within platforms [93]. In one example, Facebook ignited a controversy when their moderators flagged a breastfeeding photo as obscene, leading to a large number of protests on both sides of the debate [94]. Automated methods suffer from similar drawbacks, with multiple studies showing that biases in machine learning models can have unintended downstream consequences [95]. For example, algorithms designed to detect hate speech were more likely to label a post as "toxic" when it showed markers of African American English [96].

Finally, there is a continued arms race between bot-detection algorithms and bot creators [21,33]. As bots inevitably become more intelligent and convincingly human, the means for identifying them will have to become more precise. We observed that the majority of known bots in a sample of publicly available data sets are now tweeting about COVID-19. These bots, identified between 2011 and 2019, were discovered before the pandemic and were originally designed for non-COVID-19 purposes: promoting product hashtags, retweeting political candidates, and spreading links to malicious content. The COVID-19 pandemic will eventually end, but we have reason to believe social bots, perhaps even the same accounts, will latch on to future global issues. Additionally, we can expect bot generation techniques to advance, especially as deep learning methods continue to improve on tasks such as text or image generation [97,98]. Bot creators will continue to deploy such techniques, possibly fooling detection algorithms and humans alike. In the end, we should not expect current detection techniques, self-policing of social media platforms, or public officials alone to fully recognize, or adequately address, the current landscape of bots and misinformation.

## Acknowledgments

This work was supported in part by the Intramural Research Program of the National Institutes of Health, National Institute on Drug Abuse. MHW wrote the manuscript, with help from SG and AD. SG conceptualized the paper, with input from BC, HAS, and LU. SG performed the analyses on Twitter bots and created Tables 1 and 2. MHW, SG, AD, and MR all contributed to the literature review. LL and DHE provided crucial edits. All authors reviewed and approved the final version of the manuscript.

## **Conflicts of Interest**

None declared.

## References

- 1. Coronavirus Resource Center. COVID-19 Dashboard by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (JHU). URL: <a href="https://coronavirus.jhu.edu/map.html">https://coronavirus.jhu.edu/map.html</a> [accessed 2020-12-21]
- 2. Gozzi N, Tizzani M, Starnini M, Ciulla F, Paolotti D, Panisson A, et al. Collective Response to Media Coverage of the COVID-19 Pandemic on Reddit and Wikipedia: Mixed-Methods Analysis. J Med Internet Res 2020 Oct 12;22(10):e21597 [FREE Full text] [doi: 10.2196/21597] [Medline: 32960775]
- 3. Nielsen RK, Fletcher R, Newman N, Brennen JS, Howard PN. Navigating the 'infodemic': How people in six countries access and rate news and information about coronavirus.: Reuters Institute; 2020 Apr. URL: <a href="http://www.fundacionindex.com/fi/wp-content/uploads/2020/04/Navigating-the-Coronavirus-Infodemic-FINAL.pdf">http://www.fundacionindex.com/fi/wp-content/uploads/2020/04/Navigating-the-Coronavirus-Infodemic-FINAL.pdf</a> [accessed 2021-05-09]
- 4. Valika TS, Maurrasse SE, Reichert L. A Second Pandemic? Perspective on Information Overload in the COVID-19 Era. Otolaryngol Head Neck Surg 2020 Nov;163(5):931-933. [doi: 10.1177/0194599820935850] [Medline: 32513072]



- 5. Ball P, Maxmen A. The Epic Battle Against Coronavirus Misinformation and Conspiracy Theories. Nature 2020 May;581(7809):371-374. [doi: 10.1038/d41586-020-01452-z] [Medline: 32461658]
- 6. Tangcharoensathien V, Calleja N, Nguyen T, Purnat T, D'Agostino M, Garcia Saiso S, et al. Framework for Managing the COVID-19 Infodemic: Methods and Results of an Online, Crowdsourced WHO Technical Consultation. J Med Internet Res 2020 Jun 26;22(6):e19659 [FREE Full text] [doi: 10.2196/19659] [Medline: 32558655]
- 7. Kouzy R, Abi Jaoude J, Kraitem A, El Alam MB, Karam B, Adib E, et al. Coronavirus Goes Viral: Quantifying the COVID-19 Misinformation Epidemic on Twitter. Cureus 2020 Mar 13;12(3):e7255 [FREE Full text] [doi: 10.7759/cureus.7255] [Medline: 32292669]
- 8. Shahi G, Dirkson A, Majchrzak T. An exploratory study of COVID-19 misinformation on Twitter. Online Soc Netw Media 2021 Mar;22:100104 [FREE Full text] [doi: 10.1016/j.osnem.2020.100104] [Medline: 33623836]
- 9. Oh O, Kwon KH, Rao HR. An exploration of social media in extreme events: Rumor theory and Twitter during the Haiti earthquake 2010. 2010 Presented at: International Conference on Information Systems 2020; 2010; St. Louis, MO URL: <a href="https://aisel.aisnet.org/icis2010">https://aisel.aisnet.org/icis2010</a> submissions/231
- 10. Williamson E. How Alex Jones and Infowars helped a Florida man torment Sandy Hook families. The New York Times. 2019 Mar 29. URL: <a href="https://www.nytimes.com/2019/03/29/us/politics/alex-jones-infowars-sandy-hook.html">https://www.nytimes.com/2019/03/29/us/politics/alex-jones-infowars-sandy-hook.html</a> [accessed 2021-05-09]
- 11. Wang B, Zhuang J. Rumor response, debunking response, and decision makings of misinformed Twitter users during disasters. Nat Hazards 2018 May 11;93(3):1145-1162 [FREE Full text] [doi: 10.1007/s11069-018-3344-6]
- 12. Gupta A, Lamba H, Kumaraguru P. \$1.00 per RT #BostonMarathon #PrayForBoston: Analyzing fake content on Twitter. 2013 Presented at: 2013 APWG eCrime Researchers Summit; September 17-18, 2013; San Francisco, CA p. 1-12 URL: <a href="https://doi.org/10.1109/eCRS.2013.6805772">https://doi.org/10.1109/eCRS.2013.6805772</a> [doi: 10.1109/ecrs.2013.6805772]
- 13. Starbird K, Maddock J, Orand M, Achterman P, Mason RM. Rumors, False Flags, and Digital Vigilantes: Misinformation on Twitter after the 2013 Boston Marathon Bombing. In: iConference 2014 Proceedings. 2014 Presented at: iConference 2014; March 4-7, 2014; Berlin, Germany p. 654-662. [doi: 10.9776/14308]
- 14. Jin F, Wang W, Zhao L, Dougherty E, Cao Y, Lu C, et al. Misinformation Propagation in the Age of Twitter. IEEE Annals of the History of Computing 2014 Dec;47(12):90-94. [doi: 10.1109/MC.2014.361]
- 15. Shao C, Ciampaglia GL, Varol O, Yang K, Flammini A, Menczer F. The spread of low-credibility content by social bots. Nat Commun 2018 Nov 20;9(1):4787 [FREE Full text] [doi: 10.1038/s41467-018-06930-7] [Medline: 30459415]
- 16. Bessi A, Ferrara E. Social bots distort the 2016 U.S. Presidential election online discussion. First Monday 2016 Nov 03;21(11):1 [FREE Full text] [doi: 10.5210/fm.v21i11.7090]
- 17. Badawy A, Ferrara E, Lerman K. Analyzing the digital traces of political manipulation: The 2016 Russian interference Twitter campaign. 2018 Presented at: IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM); August 28-31, 2018; Barcelona, Spain p. 258-265. [doi: 10.1109/asonam.2018.8508646]
- 18. Yuan X, Schuchard RJ, Crooks AT. Examining Emergent Communities and Social Bots Within the Polarized Online Vaccination Debate in Twitter. Social Media + Society 2019 Sep 04;5(3):205630511986546 [FREE Full text] [doi: 10.1177/2056305119865465]
- 19. Leggio L, Garbutt JC, Addolorato G. Effectiveness and safety of baclofen in the treatment of alcohol dependent patients. CNS Neurol Disord Drug Targets 2010 Mar;9(1):33-44. [doi: 10.2174/187152710790966614] [Medline: 20201813]
- 20. Zhu Y, Guntuku SC, Lin W, Ghinea G, Redi JA. Measuring Individual Video QoE: A survey, and proposal for future directions using social media. ACM Trans Multimedia Comput Commun Appl 2018 May 22;14(2s):1-24. [doi: 10.1145/3183512]
- 21. Yang KC, Varol O, Davis CA, Ferrara E, Flammini A, Menczer F. Arming the public with artificial intelligence to counter social bots. Hum Behav & Emerg Tech 2019 Feb 06;1(1):48-61. [doi: 10.1002/hbe2.115]
- 22. Savage S, Monroy-Hernandez A, Hollerer T. Botivist: Calling volunteers to action using online bots. 2016 Presented at: 16th ACM Conference on Computer-Supported Cooperative Work & Social Computing; February 27-March 2, 2016; San Francisco, CA p. 813-822. [doi: 10.1145/2818048.2819985]
- 23. Abokhodair N, Yoo D, McDonald DW. Dissecting a Social Botnet: Growth, Content and Influence in Twitter. 2015 Presented at: 18th ACM Conference on Computer-Supported Cooperative Work & Social Computing; March 14-18, 2015; Vancouver, Canada p. 839-851. [doi: 10.1145/2675133.2675208]
- 24. Bastos MT, Mercea D. The Brexit Botnet and User-Generated Hyperpartisan News. Social Science Computer Review 2017 Oct 10;37(1):38-54. [doi: 10.1177/0894439317734157]
- 25. Jhaver S, Bruckman A, Gilbert E. Does Transparency in Moderation Really Matter?: User Behavior After Content Removal Explanations on Reddit. Proceedings of the ACM on Human-Computer Interaction 2019;3(CSCW):1-27. [doi: 10.1145/3359252]
- 26. Ma MC, Lalor JP. An Empirical Analysis of Human-Bot Interaction on Reddit.: Association for Computational Linguistics; 2020 Presented at: Sixth Workshop on Noisy User-generated Text (W-NUT 2020); November 2020; Virtual p. 101-106. [doi: 10.18653/v1/2020.wnut-1.14]
- 27. Ahmed F, Abulaish M. A generic statistical approach for spam detection in Online Social Networks. Computer Communications 2013 Jun;36(10-11):1120-1129. [doi: 10.1016/j.comcom.2013.04.004]



- 28. Hussain MN, Tokdemir S, Agarwal N, Al-Khateeb S. Analyzing disinformation and crowd manipulation tactics on YouTube. : IEEE; 2018 Presented at: IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM); August 28-31, 2018; Barcelona, Spain p. 1092-1095. [doi: 10.1109/asonam.2018.8508766]
- 29. Subrahmanian V, Azaria A, Durst S, Kagan V, Galstyan A, Lerman K, et al. The DARPA Twitter Bot Challenge. Computer 2016 Jun;49(6):38-46. [doi: 10.1109/MC.2016.183] [Medline: 27295638]
- 30. Broniatowski DA, Jamison AM, Qi S, AlKulaib L, Chen T, Benton A, et al. Weaponized Health Communication: Twitter Bots and Russian Trolls Amplify the Vaccine Debate. Am J Public Health 2018 Oct;108(10):1378-1384. [doi: 10.2105/AJPH.2018.304567] [Medline: 30138075]
- 31. Ferrara E, Varol O, Davis C, Menczer F, Flammini A. The rise of social bots. Commun ACM 2016 Jun 24;59(7):96-104. [doi: 10.1145/2818717]
- 32. Zhang J, Zhang R, Zhang Y, Yan G. The Rise of Social Botnets: Attacks and Countermeasures. IEEE Trans Dependable and Secure Comput 2018 Nov 1;15(6):1068-1082. [doi: 10.1109/tdsc.2016.2641441]
- 33. Cresci S, Di Pietro R, Petrocchi M, Spognardi A, Tesconi M. The paradigm-shift of social spambots: Evidence, theories, and tools for the arms race. 2017 Presented at: The 26th International Conference on World Wide Web Companion; April 3-7, 2017; Perth, Australia p. 963-972. [doi: 10.1145/3041021.3055135]
- 34. Chu Z, Gianvecchio S, Wang H, Jajodia S. Detecting Automation of Twitter Accounts: Are You a Human, Bot, or Cyborg? IEEE Trans Dependable and Secure Comput 2012 Nov;9(6):811-824. [doi: 10.1109/Tdsc.2012.75]
- 35. Clark EM, Williams JR, Jones CA, Galbraith RA, Danforth CM, Dodds PS. Sifting robotic from organic text: A natural language approach for detecting automation on Twitter. Journal of Computational Science 2016 Sep;16:1-7 [FREE Full text] [doi: 10.1016/j.jocs.2015.11.002]
- 36. Lee K, Eoff B, Caverlee J. Seven months with the devils: A long-term study of content polluters on Twitter. 2011 Presented at: Fifth International AAAI Conference on Weblogs and Social Media (ICWSM); July 17-21, 2011; Barcelona, Spain URL: <a href="https://www.aaai.org/ocs/index.php/ICWSM/ICWSM/1CWSM/1/paper/viewFile/2780/3296">https://www.aaai.org/ocs/index.php/ICWSM/ICWSM/1CWSM/1/paper/viewFile/2780/3296</a>
- 37. Spitzner L. The Honeynet Project: Trapping the hackers. IEEE Secur Privacy 2003 Mar;1(2):15-23. [doi: 10.1109/msecp.2003.1193207]
- 38. Prince MB, Dahl BM, Holloway L, Keller AM, Langheinrich E. Understanding how spammers steal your email address: An analysis of the first six months of data from Project Honey Pot. 2005 Presented at: Collaboration, Electronic messaging, Anti-Abuse and Spam Conference (CEAS); July 21-22, 2005; Stanford, CA URL: <a href="https://www.cse.scu.edu/~tschwarz/coen252">https://www.cse.scu.edu/~tschwarz/coen252</a> 07/Resources/spammer.pdf
- 39. Webb S, Caverlee J, Pu C. Social honeypots: Making friends with a spammer near you. 2008 Presented at: Collaboration, Electronic messaging, Anti-Abuse and Spam Conference (CEAS); August 21-22, 2008; Mountain View, CA p. 1-10 URL: <a href="https://people.engr.tamu.edu/caverlee/pubs/webb08socialhoneypots.pdf">https://people.engr.tamu.edu/caverlee/pubs/webb08socialhoneypots.pdf</a>
- 40. Davis CA, Varol O, Ferrara E, Flammini A, Menczer F. Botornot: A system to evaluate social bots. 2016 Presented at: 25th International Conference Companion on World Wide Web; April 11-15, 2016; Montreal, Canada. [doi: 10.1145/2872518.2889302]
- 41. Varol O, Ferrara E, Davis CA, Menczer F, Flammini A. Online human-bot interactions: Detection, estimation, and characterization. 2017 Presented at: International AAAI Conference on Web and Social Media; May 15-18, 2017; Montreal, Canada URL: <a href="https://arxiv.org/abs/1703.03107">https://arxiv.org/abs/1703.03107</a>
- 42. Dickerson JP, Kagan V, Subrahmanian V. Using sentiment to detect bots on twitter: Are humans more opinionated than bots? 2014 Presented at: IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM); Aug 17-20, 2014; Beijing, China. [doi: 10.1109/asonam.2014.6921650]
- 43. Stephens M, Poorthuis A. Follow thy neighbor: Connecting the social and the spatial networks on Twitter. Computers, Environment and Urban Systems 2015 Sep;53:87-95. [doi: <a href="https://doi.org/10.1016/j.compenvurbsys.2014.07.002">10.1016/j.compenvurbsys.2014.07.002</a>]
- 44. Yang Z, Wilson C, Wang X, Gao T, Zhao BY, Dai Y. Uncovering social network Sybils in the wild. ACM Trans Knowl Discov Data 2014 Feb;8(1):1-29. [doi: 10.1145/2556609]
- 45. Pozzana I, Ferrara E. Measuring Bot and Human Behavioral Dynamics. Front Phys 2020 Apr 22;8:125. [doi: 10.3389/fphy.2020.00125]
- 46. Sahami M, Dumais S, Heckerman D, Horvitz E. A Bayesian approach to filtering junk email. Learning for Text Categorization: Papers from the AAAI workshop. 1998. URL: <a href="https://www.aaai.org/Papers/Workshops/1998/WS-98-05/WS98-05-009.pdf">https://www.aaai.org/Papers/Workshops/1998/WS-98-05/WS98-05-009.pdf</a> [accessed 2021-05-11]
- 47. Jagatic TN, Johnson NA, Jakobsson M, Menczer F. Social phishing. Commun ACM 2007 Oct;50(10):94-100 [FREE Full text] [doi: 10.1145/1290958.1290968]
- 48. Stringhini G, Kruegel C, Vigna G. Detecting spammers on social networks. 2010 Presented at: Annual Computer Security Applications Conference; December 6-10, 2010; Austin, TX p. 1-9. [doi: 10.1145/1920261.1920263]
- 49. Kudugunta S, Ferrara E. Deep neural networks for bot detection. Information Sciences 2018 Oct;467:312-322. [doi: 10.1016/j.ins.2018.08.019]
- 50. Mazza M, Cresci S, Avvenuti M, Quattrociocchi W, Tesconi M. RTbust: Exploiting temporal patterns for botnet detection on Twitter. 2019 Presented at: ACM Conference on Web Science; June 30, 2019; Boston, MA, USA. [doi: 10.1145/3292522.3326015]



- 51. Santia GC, Mujib MI, Williams JR. Detecting Social Bots on Facebook in an Information Veracity Context. 2019 Presented at: International AAAI Conference on Web and Social Media; June 11, 2019; Munich, Germany URL: <a href="https://ojs.aaai.org/index.php/ICWSM/article/view/3244">https://ojs.aaai.org/index.php/ICWSM/article/view/3244</a>
- 52. Yang KC, Varol O, Hui PM, Menczer F. Scalable and Generalizable Social Bot Detection through Data Selection. 2020 Presented at: AAAI Conference on Artificial Intelligence; February 7-12, 2020; New York, NY p. 1096-1103. [doi: 10.1609/aaai.v34i01.5460]
- 53. Treen KM, Williams HT, O'Neill SJ. Online misinformation about climate change. WIREs Clim Change 2020 Jun 18;11(5):e665 [FREE Full text] [doi: 10.1002/wcc.665]
- 54. Shahi GK, Dirkson A, Majchrzak TA. An exploratory study of COVID-19 misinformation on Twitter. Online Soc Netw Media 2021 Mar;22:100104 [FREE Full text] [doi: 10.1016/j.osnem.2020.100104] [Medline: 33623836]
- 55. Dwoskin E, Timberg C. Misinformation dropped dramatically the week after Twitter banned Trump and some allies. The Washington Post. 2021 Jan 16. URL: <a href="https://www.washingtonpost.com/technology/2021/01/16/misinformation-trump-twitter/">https://www.washingtonpost.com/technology/2021/01/16/misinformation-trump-twitter/</a> [accessed 2021-05-09]
- 56. Vosoughi S, Roy D, Aral S. The spread of true and false news online. Science 2018 Mar 09;359(6380):1146-1151. [doi: 10.1126/science.aap9559] [Medline: 29590045]
- 57. Wang P, Angarita R, Renna I. Is this the era of misinformation yet: combining social bots and fake news to deceive the masses. 2018 Presented at: Companion Proceedings of The Web Conference 2018; April 23-27, 2018; Lyon, France p. 1557-1561. [doi: 10.1145/3184558.3191610]
- 58. Bovet A, Makse HA. Influence of fake news in Twitter during the 2016 US presidential election. Nat Commun 2019 Jan 2;10(1):1. [doi: 10.1038/s41467-018-07761-2]
- 59. Ferrara E. Disinformation and social bot operations in the run up to the 2017 French presidential election. First Monday 2017 Jul 31;22(8):1. [doi: 10.5210/fm.v22i8.8005]
- 60. Howard PN. How Political Campaigns Weaponize Social Media Bots. IEEE Spectrum. 2018. URL: <a href="https://spectrum.ieee.org/computing/software/how-political-campaigns-weaponize-social-media-bots">https://spectrum.ieee.org/computing/software/how-political-campaigns-weaponize-social-media-bots</a> [accessed 2021-05-11]
- 61. Mueller RS. Report on the investigation into Russian interference in the 2016 presidential election. US Department of Justice. URL: <a href="https://www.justice.gov/storage/report.pdf">https://www.justice.gov/storage/report.pdf</a> [accessed 2021-05-09]
- 62. Twitter. Update on Twitter's review of the 2016 US election. Twitter Blog. 2018 Jan 19. URL: <a href="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</a> [accessed 2021-05-09]
- 63. Rheault L, Musulan A. Efficient detection of online communities and social bot activity during electoral campaigns. Journal of Information Technology & Politics 2021 Feb 02:1-14. [doi: 10.1080/19331681.2021.1879705]
- 64. Howard PN, Woolley S, Calo R. Algorithms, bots, and political communication in the US 2016 election: The challenge of automated political communication for election law and administration. Journal of Information Technology & Politics 2018 Apr 11;15(2):81-93 [FREE Full text] [doi: 10.1080/19331681.2018.1448735]
- 65. Assenmacher D, Clever L, Frischlich L, Quandt T, Trautmann H, Grimme C. Demystifying Social Bots: On the Intelligence of Automated Social Media Actors. Social Media + Society 2020 Sep 01;6(3):205630512093926. [doi: 10.1177/2056305120939264]
- 66. Varol O. Bot Repository. URL: <a href="https://botometer.osome.iu.edu/bot-repository/index.html">https://botometer.osome.iu.edu/bot-repository/index.html</a> [accessed 2020-11-20]
- 67. Giorgi S, Sap M. TwitterMySQL. URL: <a href="https://github.com/dlatk/TwitterMySQL">https://github.com/dlatk/TwitterMySQL</a> [accessed 2020-11-20]
- 68. Al-Rawi A, Shukla V. Bots as Active News Promoters: A Digital Analysis of COVID-19 Tweets. Information 2020 Sep 27;11(10):461. [doi: 10.3390/info11100461]
- 69. Gilani Z, Farahbakhsh R, Tyson G, Wang L, Crowcroft J. Of bots and humans (on Twitter). 2017 Presented at: IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining; July 31, 2017; Sydney, Australia. [doi: 10.1145/3110025.3110090]
- 70. Shi W, Liu D, Yang J, Zhang J, Wen S, Su J. Social Bots' Sentiment Engagement in Health Emergencies: A Topic-Based Analysis of the COVID-19 Pandemic Discussions on Twitter. Int J Environ Res Public Health 2020 Nov 23;17(22):8701. [doi: 10.3390/ijerph17228701] [Medline: 33238567]
- 71. Allyn B. Researchers: Nearly half of accounts tweeting about coronavirus are likely bots. NPR. 2020 May 20. URL: <a href="https://www.npr.org/sections/coronavirus-live-updates/2020/05/20/859814085/">https://www.npr.org/sections/coronavirus-live-updates/2020/05/20/859814085/</a>
  researchers-nearly-half-of-accounts-tweeting-about-coronavirus-are-likely-bots [accessed 2021-05-09]
- 72. Roth Y, Pickles N. Bot or not? The facts about platform manipulation on Twitter. Twitter Blog. 2020 May 18. URL: <a href="https://blog.twitter.com/en\_us/topics/company/2020/bot-or-not.html">https://blog.twitter.com/en\_us/topics/company/2020/bot-or-not.html</a> [accessed 2021-05-09]
- 73. Ferrara E. What types of COVID-19 conspiracies are populated by Twitter bots? First Monday 2020 May 19;25(6):1. [doi: 10.5210/fm.v25i6.10633]
- 74. Allem J, Ferrara E, Uppu SP, Cruz TB, Unger JB. E-Cigarette Surveillance With Social Media Data: Social Bots, Emerging Topics, and Trends. JMIR Public Health Surveill 2017 Dec 20;3(4):e98 [FREE Full text] [doi: 10.2196/publichealth.8641] [Medline: 29263018]
- 75. Allem J, Escobedo P, Dharmapuri L. Cannabis Surveillance With Twitter Data: Emerging Topics and Social Bots. Am J Public Health 2020 Mar;110(3):357-362. [doi: 10.2105/AJPH.2019.305461] [Medline: 31855475]



- 76. Mehta B, Salmon J, Ibrahim S. Potential Shortages of Hydroxychloroquine for Patients with Lupus During the Coronavirus Disease 2019 Pandemic. JAMA Health Forum 2020 Apr 10;1(4):e200438. [doi: 10.1001/jamahealthforum.2020.0438]
- 77. Lupus Research Alliance. COVID-19 Caused Hydroxychloroquine Issues for Third of Lupus Patients, New LRA Survey Finds. 2020 May 28. URL: <a href="https://www.lupusresearch.org/covid-19-caused-hydroxychloroquine-issues-for-third-of-lupus-patients-new-lra-survey-finds/">https://www.lupusresearch.org/covid-19-caused-hydroxychloroquine-issues-for-third-of-lupus-patients-new-lra-survey-finds/</a> [accessed 2021-05-09]
- 78. Skipper CP, Pastick KA, Engen NW, Bangdiwala AS, Abassi M, Lofgren SM, et al. Hydroxychloroquine in Nonhospitalized Adults With Early COVID-19: A Randomized Trial. Ann Intern Med 2020 Oct 20;173(8):623-631 [FREE Full text] [doi: 10.7326/M20-4207] [Medline: 32673060]
- 79. Mitjà O, Corbacho-Monné M, Ubals M, Tebe C, Peñafiel J, Tobias A, BCN PEP-CoV-2 RESEARCH GROUP. Hydroxychloroquine for Early Treatment of Adults with Mild Covid-19: A Randomized-Controlled Trial. Clin Infect Dis 2020 Jul 16:1 [FREE Full text] [doi: 10.1093/cid/ciaa1009] [Medline: 32674126]
- 80. National Institutes of Health. NIH halts clinical trial of hydroxychloroquine. National Institutes of Health News Releases. 2020 Jun 20. URL: <a href="https://www.nih.gov/news-events/news-releases/nih-halts-clinical-trial-hydroxychloroquine">https://www.nih.gov/news-events/news-releases/nih-halts-clinical-trial-hydroxychloroquine</a> [accessed 2021-05-09]
- 81. Allington D, Duffy B, Wessely S, Dhavan N, Rubin J. Health-protective behaviour, social media usage and conspiracy belief during the COVID-19 public health emergency. Psychol Med 2020 Jun 09:1-7 [FREE Full text] [doi: 10.1017/S003329172000224X] [Medline: 32513320]
- 82. Loomba S, de Figueiredo A, Piatek SJ, de Graaf K, Larson HJ. Measuring the impact of COVID-19 vaccine misinformation on vaccination intent in the UK and USA. Nat Hum Behav 2021 Mar 05;5(3):337-348. [doi: 10.1038/s41562-021-01056-1] [Medline: 33547453]
- 83. Roth Y, Pickles N. Updating our approach to misleading information. Twitter Blog. 2020 May 11. URL: <a href="https://blog.twitter.com/en\_us/topics/product/2020/updating-our-approach-to-misleading-information.html">https://blog.twitter.com/en\_us/topics/product/2020/updating-our-approach-to-misleading-information.html</a> [accessed 2021-05-09]
- 84. Statt N. Facebook will now show a warning before you share articles about COVID-19. The Verge. 2020 Aug 12. URL: <a href="https://www.theverge.com/2020/8/12/21365305/facebook-covid-19-warning-notification-post-misinformation">https://www.theverge.com/2020/8/12/21365305/facebook-covid-19-warning-notification-post-misinformation</a> [accessed 2021-05-09]
- 85. Timberg C, Dwoskin E. Twitter is sweeping out fake accounts like never before, putting user growth at risk. The Washington Post. 2018 Jul 09. URL: <a href="https://www.washingtonpost.com/technology/2018/07/06/">https://www.washingtonpost.com/technology/2018/07/06/</a>
  <a href="twitter-is-sweeping-out-fake-accounts-like-never-before-putting-user-growth-risk/">https://www.washingtonpost.com/technology/2018/07/06/</a>
  <a href="twitter-is-sweeping-out-fake-growth-risk-user-growth-risk-user-before-putting-user-growth-risk-user-growth-risk-us
- 86. Facebook. Community Standards Enforcement Report. Facebook Transparency. URL: <a href="https://transparency.facebook.com/community-standards-enforcement">https://transparency.facebook.com/community-standards-enforcement</a> [accessed 2021-05-09]
- 87. Kamal G. California's BOT Disclosure Law, SB 1001, now in effect. The National Law Review. 2019 Jul 15. URL: <a href="https://www.natlawreview.com/article/california-s-bot-disclosure-law-sb-1001-now-effect">https://www.natlawreview.com/article/california-s-bot-disclosure-law-sb-1001-now-effect</a> [accessed 2021-05-09]
- 88. Frazin R. Feinstein introduces bill to prohibit campaigns from using social media bots. The Hill. 2019 Jul 16. URL: <a href="https://thehill.com/policy/cybersecurity/453336-dem-senator-introduces-bill-to-prohibit-campaigns-from-using-bots">https://thehill.com/policy/cybersecurity/453336-dem-senator-introduces-bill-to-prohibit-campaigns-from-using-bots</a> [accessed 2021-05-09]
- 89. Kang C, Fandos N, Issac M. Tech executives are contrite about election meddling, but make few promises on Capitol Hill. The New York Times. 2017 Oct 31. URL: <a href="https://www.nytimes.com/2017/10/31/us/politics/facebook-twitter-google-hearings-congress.html">https://www.nytimes.com/2017/10/31/us/politics/facebook-twitter-google-hearings-congress.html</a> [accessed 2021-05-09]
- 90. Gorwa R, Guilbeault D. Unpacking the Social Media Bot: A Typology to Guide Research and Policy. Policy & Internet 2018 Aug 10;12(2):225-248. [doi: 10.1002/poi3.184]
- 91. Samples J. Why the government should not regulate content moderation of social media. Cato Institute. 2019 Apr 09. URL: <a href="https://www.cato.org/publications/policy-analysis/why-government-should-not-regulate-content-moderation-social-media">https://www.cato.org/publications/policy-analysis/why-government-should-not-regulate-content-moderation-social-media</a> [accessed 2021-05-09]
- 92. Crews Jr CW. The Case against Social Media Content Regulation: Reaffirming Congress' Duty to Protect Online Bias, 'Harmful Content,' and Dissident Speech from the Administrative State. SSRN 2020 Jun 28; Competitive Enterprise Institute, Issue Analysis:1-34 [FREE Full text]
- 93. Gillespie T. Custodians of the Internet: Platforms, Content Moderation, and the Hidden Decisions That Shape Social Media. New Haven, CT: Yale University Press; 2018.
- 94. Ibrahim Y. The breastfeeding controversy and Facebook: Politicization of image, privacy and protest. International Journal of E-Politics (IJEP) 2010;1(2):16-28. [doi: 10.4018/jep.2010040102]
- 95. Dixon L, Li J, Sorensen J, Thain N, Vasserman L. Measuring and mitigating unintended bias in text classification. 2018 Presented at: AAAI/ACM Conference on AI, Ethics, and Society; 2018; New Orleans, LA p. 67-73. [doi: 10.1145/3278721.3278729]
- Sap M, Card D, Gabriel S, Choi Y, Smith N. The risk of racial bias in hate speech detection. 2019 Presented at: Annual Meeting of the Association for Computational Linguistics; July 2019; Florence, Italy p. 1668-1678. [doi: 10.18653/v1/p19-1163]
- 97. Zhao J, Xiong L, Jayashree P. Dual-agent GANs for photorealistic and identity preserving profile face synthesis. 2017 Presented at: Advances in Neural Information Processing Systems; December 4, 2017; Long Beach, CA, USA p. 66-76



URL: http://papers.neurips.cc/paper/6612-dual-agent-gans-for-photorealistic-and-identity-preserving-profile-face-synthesis.pdf

98. Brown TB, Mann B, Ryder N, Subbiah M, Kaplan J, Dhariwal P. Language models are few-shot learners. 2020 Presented at: Advances in Neural Information Processing Systems; December 6, 2020; Virtual p. 1877-1901 URL: <a href="https://arxiv.org/abs/2005.14165">https://arxiv.org/abs/2005.14165</a>

#### **Abbreviations**

**API:** application programming interface

Edited by R Kukafka, C Basch; submitted 04.01.21; peer-reviewed by A Agarwal, D Yeung, N Martinez-Martin; comments to author 16.02.21; revised version received 04.03.21; accepted 16.04.21; published 20.05.21

<u>Please cite as:</u>

Himelein-Wachowiak M, Giorgi S, Devoto A, Rahman M, Ungar L, Schwartz HA, Epstein DH, Leggio L, Curtis B Bots and Misinformation Spread on Social Media: Implications for COVID-19

J Med Internet Res 2021;23(5):e26933 URL: https://www.jmir.org/2021/5/e26933

doi: <u>10.2196/26933</u> PMID: <u>33882014</u>

©McKenzie Himelein-Wachowiak, Salvatore Giorgi, Amanda Devoto, Muhammad Rahman, Lyle Ungar, H Andrew Schwartz, David H Epstein, Lorenzo Leggio, Brenda Curtis. Originally published in the Journal of Medical Internet Research (https://www.jmir.org), 20.05.2021. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in the Journal of Medical Internet Research, is properly cited. The complete bibliographic information, a link to the original publication on https://www.jmir.org/, as well as this copyright and license information must be included.

