# JMIR Infodemiology

# Contents

## Original Papers

## Corrigenda and Addenda

Original Paper

# Analyzing Social Media to Explore the Attitudes and Behaviors Following the Announcement of Successful COVID-19 Vaccine Trials: Infodemiology Study

Jean-Christophe Boucher[1], PhD; Kirsten Cornelson[2], PhD; Jamie L Benham[3,4], MD; Madison M Fullerton[4], MSc; Theresa Tang[4], BHSc (Hons); Cora Constantinescu[5], MD; Mehdi Mourali[6], PhD; Robert J Oxoby[7], PhD; Deborah A Marshall[3,4], PhD; Hadi Hemmati[8], PhD; Abbas Badami[8], BsEng; Jia Hu[4], MD; Raynell Lang[3], MD

[1]School of Public Policy and Department of Political Science, University of Calgary, Calgary, AB, Canada

[2]Department of Economics, University of Notre Dame, Notre Dame, IN, United States

[3]Department of Medicine, Cumming School of Medicine, University of Calgary, Calgary, AB, Canada

[4]Department of Community Health Sciences, Cumming School of Medicine, University of Calgary, Calgary, AB, Canada

[5]Department of Pediatrics, University of Calgary, Calgary, AB, Canada

[6]Haskayne School of Business, University of Calgary, Calgary, AB, Canada

[7]Department of Economics, Faculty of Arts, University of Calgary, Calgary, AB, Canada

[8]Department of Electrical and Computer Engineering, Schulich Faculty of Engineering, University of Calgary, Calgary, AB, Canada

**Corresponding Author:**
Jean-Christophe Boucher, PhD
School of Public Policy and Department of Political Science
University of Calgary
906 8 Ave SW 5th Floor
Calgary, AB
Canada
Phone: 1 403 220 8565
Email: jc.boucher@ucalgary.ca

## *Abstract*

**Background:**   The rollout of COVID-19 vaccines has brought vaccine hesitancy to the forefront in managing this pandemic. COVID-19 vaccine hesitancy is fundamentally different from that of other vaccines due to the new technologies being used, rapid development, and widespread global distribution. Attitudes on vaccines are largely driven by online information, particularly information on social media. The first step toward influencing attitudes about immunization is understanding the current patterns of communication that characterize the immunization debate on social media platforms.

**Objective:**   We aimed to evaluate societal attitudes, communication trends, and barriers to COVID-19 vaccine uptake through social media content analysis to inform communication strategies promoting vaccine acceptance.

**Methods:**   Social network analysis (SNA) and unsupervised machine learning were used to characterize COVID-19 vaccine content on Twitter globally. Tweets published in English and French were collected through the Twitter application programming interface between November 19 and 26, 2020, just following the announcement of initial COVID-19 vaccine trials. SNA was used to identify social media clusters expressing mistrustful opinions on COVID-19 vaccination. Based on the SNA results, an unsupervised machine learning approach to natural language processing using a sentence-level algorithm transfer function to detect semantic textual similarity was performed in order to identify the main themes of vaccine hesitancy.

**Results:**   The tweets (n=636,516) identified that the main themes driving the vaccine hesitancy conversation were concerns of safety, efficacy, and freedom, and mistrust in institutions (either the government or multinational corporations). A main theme was the safety and efficacy of mRNA technology and side effects. The conversation around efficacy was that vaccines were unlikely to completely rid the population of COVID-19, polymerase chain reaction testing is flawed, and there is no indication of long-term T-cell immunity for COVID-19. Nearly one-third (45,628/146,191, 31.2%) of the conversations on COVID-19 vaccine hesitancy clusters expressed concerns for freedom or mistrust of institutions (either the government or multinational corporations) and nearly a quarter (34,756/146,191, 23.8%) expressed criticism toward the government's handling of the pandemic.

**Conclusions:** Social media content analysis combined with social network analysis provides insights into the themes of the vaccination conversation on Twitter. The themes of safety, efficacy, and trust in institutions will need to be considered, as targeted outreach programs and intervention strategies are deployed on Twitter to improve the uptake of COVID-19 vaccination.

## Introduction

The COVID-19 pandemic is emerging as one of the greatest public health threats in history, with over 140 million infections and 3 million deaths worldwide attributed to the SARS-CoV-2 virus as of April 2021 [1]. As transmission of COVID-19 continues around the globe, a COVID-19 vaccine is an important and valuable tool to reduce the spread of infection. Due to the critical need, the speed of vaccine development, production, and mass rollout has been faster than ever seen before, leading to concerns about vaccine efficacy and safety [2,3]. The intricacy surrounding COVID-19 vaccine hesitancy appears to be further reaching and fundamentally different than other immunizations [2]. Vaccine production cannot meet demand, requiring rollout plans for targeted subpopulations, and there are several types of COVID-19 vaccines being used within one country [2]. This has led to increased concern, mistrust, and confusion surrounding COVID-19 vaccination.

Despite the massive undertaking of vaccine development, a vaccination program is only as successful as its uptake. Vaccine hesitancy, defined as a delay in acceptance or refusal of vaccines despite their availability [4], was mentioned by the World Health Organization as one of the top 10 threats to public health in 2019 [5]. The effective rollout of a COVID-19 vaccine strategy may be obstructed by the beliefs and attitudes of vaccine-hesitant individuals worldwide [6]. A recent survey of US adults in April 2020 found that 23% of persons would not be willing to get vaccinated against COVID-19 [7]. In Canada, a survey in March 2021 demonstrated that 76.9% of Canadians were very or somewhat willing to receive a COVID-19 vaccine [8].

Vaccine hesitancy is a "multifaceted, deeply complex construct that may be rooted in the moral composition that guides our daily decision making" [8,9]. Studies have identified that vaccination decisions are shaped by multiple complex interactions between individual, community, cultural, historical, political, and societal factors [2,10]. There have been several metrics and scales developed and used to measure vaccine confidence and hesitancy [11-13]. The Vaccine Confidence Index (VCI) survey tool was developed to measure individuals' perceptions on the safety, importance, effectiveness, and religious beliefs toward vaccines, which were identified as key drivers of public confidence in vaccination [11]. Another validated score measuring the psychological antecedents of vaccination is known as the 5C scale [13]. The 5C scale includes confidence (trust and attitudes), complacency (not perceiving diseases as high risk), constraints (structural and psychological barriers), collective responsibility (willingness to protect others), and calculation (information searching) [13].

Negative beliefs about vaccines may prevent the implementation of provaccination policies. Public health officials need to prioritize implementing strategies to help reduce these negative perceptions [14]. Traditional approaches to promoting immunization have assumed that inadequate knowledge of the associated risks and benefits drive hesitancy; however, this stance has proven to be ineffective as an intervention [15]. The source of the hesitancy is often both from lack of information and from lack of trust in institutions such as the government, physicians, and pharmaceutical companies [16,17]. Rather than just trying to enhance knowledge, a different approach to overcoming vaccine hesitancy for COVID-19 may be to focus on changing personal attitudes [15].

Prior studies have demonstrated that social media can help understand attitudes and behaviors during public health crises and promote health messaging [18-20]. Particularly, Twitter has been used for public health research and more specifically regarding COVID-19 [18,21]. A recent study aimed to characterize the main topics of Twitter conversations related to COVID-19 and identified four main themes including the origin of the virus, its sources, the impact on people/countries/economy, and lastly methods of mitigating the risk of infection [18].

Attitudes toward vaccination are, in large part, shaped by information and ideas individuals encounter through social media [22,23]. Social media is a principal informational forum for vaccination uptake with large proportions of content involving antivaccination messaging [22,24]. Reliable and accurate information on social media is often mixed with inaccurate, conspiratorial, incomplete, or biased messages [22]. A recent study evaluating vaccine hesitancy through content analysis of tweets in Canada identified major themes including safety, suspicion of economic or political motivation, knowledge deficit, opinions of authority figures, and lack of liability from pharmaceutical companies [25]. This study demonstrates the significant utility of using Twitter to better understand vaccine hesitancy at a population level [25]. However, a broader reaching and deeper understanding of current patterns of communication that characterize the immunization debate on social media platforms is needed to inform public health interventions aimed toward influencing attitudes on immunization [15,17,22,23].

This study used social network analysis (SNA) and unsupervised machine learning to characterize COVID-19 vaccine content on Twitter, which allows researchers to access the application programming interface (API). The use of machine learning in social sciences is expanding and has generated interesting methodological conversations on different approaches to study

COVID-19 sentiment, attitude, and emotion [26-28]. With this study, we aimed to use social media content analysis to provide a further understanding of societal attitudes, communication trends, and barriers to COVID-19 vaccine uptake, at a critical time during the COVID-19 pandemic, just following reporting of initial vaccination trials. We hypothesize that these data will be critical globally for developing targeted outreach programs and intervention strategies on Twitter to impact COVID-19 vaccine hesitancy.

## Methods

### Study Data

Tweets published in English and French using specific words ("COVID" AND ("vacc" OR "vax" OR "immu")) in either the content of the tweet or hashtag derivatives were collected through the Twitter API between November 19 and 26, 2020. This Boolean query was selected in order to maximize inclusivity without adding unnecessary noise to our data set. For example, with the query "vacc," we effectively targeted all derivatives such as "vaccine," "vaccines," "vaccinated," "vaccination," "vaccinations," and associated hashtags. Furthermore, it follows search queries from similar studies examining vaccine hesitancy and social media [29-31]. English and French were selected as they are two of the most common languages used on social media. The data set included several features, such as descriptive information about the user, username, content of tweets (hashtags, relationship among users such as retweet, replies, and mentions, etc), self-reported location of the user, number of followers, date of account creation, and time of tweet posting. Tweets were extracted using Twitter's public streaming API allowing researchers to collect a random sample of tweets in real time of up to 1% of all public tweets published daily.

This time period was crucial in COVID-19 vaccination conversations on social media as it came a week after pharmaceutical companies (Pfizer and Moderna) announced successful trials of their COVID-19 vaccines [32,33]. For the first time since the beginning of the COVID-19 pandemic in early 2020, social media conversations on vaccination were based, at least in part, on plausible empirical information about the efficacy and availability of a vaccine. We therefore wanted to evaluate the public's initial reaction and response to COVID-19 vaccination. Few prior studies have included this time frame during the analysis period [30,34-36].

### Social Network Analysis

In this study, we designed a data analytics workflow by first using SNA to identify social media accounts most likely expressing doubtful or mistrustful opinions on COVID-19 vaccination. This network was created using a weighted retweet directed network to represent connections between accounts. Although retweets are not a perfect indicator of like-mindedness, on aggregate, users have a proclivity to engage more with accounts that reflect some form of social or intellectual homophily [37]. Through SNA, we examined the underlying structure of community clustering within the broader social network exposed by online interactions, isolating different "communities" of like-minded users. The Louvain modularity method [38] was used to detect subclusters of online communities mentioning COVID-19 vaccination.

### Natural Language Programming Analysis

Based on our SNA results, we developed an unsupervised machine learning approach to natural language processing by using a sentence-level algorithm transfer function to detect semantic textual similarity [15]. Our goal was to examine how antivaccine conversation clusters talk about or frame a possible COVID-19 vaccine without prior assumptions about the nature of the conversation. For this analysis, we first tokenized our sentences (tweets) and cleaned the data set, removing duplicates, stop words, symbols, numbers, punctuation, URLs, whitespaces, and stemmed words to their roots. We also added a language identifier to remove tweets in languages other than French or English and reduce noise in the data set. We then used DistilBERT [39], a knowledge distillation learning model, for sentence-level embedding. DistilBERT is a compressed version of BERT (Bidirectional Encoder Representations from Transformers), which retains much of the computational accuracy of BERT without the environmental costs associated with high-dimensionality embeddings. DistilBERT positions all sentences (here tweets from the antivaccine conversation) in a multidimensionality vector space from which we can compare the semantic similarity of tweets.

We used an agglomerative hierarchical cluster model to identify a relevant number of clusters from the multidimensionality output produced by the DistilBERT computation, which provided us a measure to identify different topics of similar tweets. To infer topics from our clustering modeling, we used a bi-gram of term frequency-inverse document frequency (tf-idf), which measures the originality of a word by comparing the number of occurrences of the word in a document (term frequency) and the number of documents with the word (inverse document frequency). This measure allows us to undervalue words that appear frequently in most documents (such as "the") and provide little information, and overvalue words that appear sporadically in the corpus, but often in some documents. Topics were then inferred manually based on the cluster model and informed by the tf-idf output. Figure 1 illustrates our data analytics pipeline.

**Figure 1.** Data analytics workflow. NLP: natural language processing.



## Results

In total, 636,516 tweets were collected from 428,535 accounts, for an average of 79,564 tweets per day. Figure 2 presents the cluster map of COVID-19 immunization conversations on Twitter between November 19 and November 26, 2020. We found a polarized conversation about immunization on social media.

During this observation period, a large proportion of accounts debating COVID-19 immunization on Twitter were connected and exposed to social media conversations promoting vaccination narratives. The largest cluster (green), comprising approximately 49.4% (n=211,549) of Twitter accounts, revealed a vaccine acceptant point of view. Based on degree centrality, the cluster seemed to overlap with more progressive-leaning political leaders, such as Hillary Clinton and US Representative Alexandria Ocasio-Cortez (D-14-NY), and mainstream news

media, such as the NY Post, the Hill, ABC, and Reuters. A second provaccine cluster (orange) centered on the Indian COVID-19 immunization online debate, with political leaders, such as Rahul Gandhi and Press Trust of India (the largest news agency in India), leading the conversation.

As Figure 2 indicates, we also found two clusters opposite to these vaccine acceptant clusters exhibiting more vaccine hesitant narratives. There were 23.4% (n=146,191) of conversations on Twitter during this period of observation that can be directly attributed to vaccine hesitancy. First, in red, our study identified a large cluster comprising 88,892 Twitter handles accounting for 18.4% of all accounts in the cluster. These interactions from the Twittersphere in English originated mostly from the United States, the United Kingdom, and Canada, and gravitated around accounts of prominent antivaccine physicians and organizations, right-wing activists, show hosts, such as Rush Limbaugh in the United States and Simon Dolan or Michael J Blair in the United Kingdom, and some alternative news organizations, such as

Breitbart News. In this conversation cluster, we found an overlap between ideologically leaning advocates, especially those associated at the margins of right-wing or Conservative parties, and antivaccine online conversations. We found considerable cross-pollination between accounts originating from all across the English Twittersphere, demonstrating the internationalization of COVID-19 vaccination conversations on social media. Second, in blue, a smaller cluster representing 2.7% (n=11,509) of accounts appeared to be shaped around Francophone (from France and Québec) vaccine hesitant conversations.

**Figure 2.** Twitter retweet cluster of COVID-19 vaccination (November 19 to November 26, 2020). The vaccine acceptant cluster (n=211,549), vaccine hesitant cluster (n=88,892), Indian vaccine acceptant cluster (n=28,713), and French vaccine hesitant cluster (n=11,509) are seen. Nodes represent specific Twitter accounts, while edges represent retweet activity between accounts. Presented are the four largest online communities mentioning COVID-19 vaccination.



Examining the COVID-19 vaccine hesitancy narrative during this observation period provided key insights on what drives attitudes. As shown in Table 1, our unsupervised machine learning analysis identified 12 specific archetypical vaccine hesitancy tweets (we tested a different number of clusters between 9 to 20). Social media content highlighted concerns over the efficacy and the safety of a possible COVID-19 vaccine fitting into the 5C scale domain of confidence. Five broad categories focused on vaccine efficacy in our data set. The first topic (topic 3 in Table 1) regrouped tweets suggesting that attempts to produce a COVID-19 vaccine, especially using mRNA technology, remain tentative. These tweets highlighted our failure to develop an HIV vaccine with mRNA and that existing mRNA vaccines for COVID-19 have not been tested enough to demonstrate efficacy. The second category of vaccine hesitancy tweets concentrating on efficacy (topic 2) focused on a comment made by Sir John Irving Bell on November 20, 2020, that existing vaccines were unlikely to completely rid the population of COVID-19. The third and fourth vaccine efficacy topics (topics 7 and 8) grouped together tweets from prominent physicians arguing that existing polymerase chain reaction testing is flawed and that there are no indications of long-term T-cell immunity for COVID-19. The last topic highlighted (topic 10), which obtained subsequent media attention, suggested that as much as 25% of the population would have contracted COVID-19 by the time a vaccine would be rolled out, and consequently, the vaccine would be unnecessary to reach herd immunity.

As for the concern of vaccine safety, we found two broad groups of tweets. The first (topic 5) highlighted a classic antivaccine story where an emergency medical technician/fire rescue was required to get TDAP (tetanus, diphtheria, and pertussis) boosters and had complications. Although not specifically addressing the issue of a COVID-19 vaccine, these tweets emphasize how existing antivaccine narratives have created a baseline from which some individuals frame a COVID-19 vaccine. The second topic (topic 9) centered on tweets in French framing a COVID-19 vaccine as a poison and suggesting that mRNA technology has not been tested yet and would be harmful.

Although safety and efficacy concerns remain a major component explaining COVID-19 vaccine hesitancy, it is only half of the story. As we examined the content of tweets of users integrated in the vaccine hesitancy online cluster, we found a large proportion of those interactions that emphasized a concern for personal freedom and/or some form of mistrust of institutions. In the 5C scale of psychological antecedents of vaccination, confidence in vaccines includes trust in the system that delivers the vaccines and the motivation of policy makers who decide on the need of vaccines [13]. Tweets were framed in three directions. First, a large percentage of such tweets (n=45,628, 31.2%) expressed some criticism toward the government's handling of the COVID-19 crisis, especially decisions to curb individual freedom. COVID-19 immunization is framed as the right of individuals to decide for themselves whether or not to be vaccinated. Any indication that governmental authorities might require vaccination is perceived

as a direct attempt to limit individual freedom. Two main topics (topics 1 and 12) highlighted mistrust in government policy, with one condemning the Johnson government in the United Kingdom for proposing that individuals vaccinated for COVID-19 could receive "freedom passes" and one criticizing Denmark's decision to cull its mink population to halt the spread of a coronavirus variant. Some topics (topics 6 and 11) expressed mistrust in multinational corporations, most notably airline companies, such a Qantas, who suggested that vaccination should be made compulsory for international travel. This decision is framed as a direct restriction of personal freedom by multinational corporations, and some users made a direct reference to populist narratives suggesting that these policies would only affect commercial flights and elites would be able to avoid vaccination while using private flights.

**Table 1.** Inferred topic analysis of the COVID-19 vaccination hesitancy cluster (November 19 to November 26, 2020).
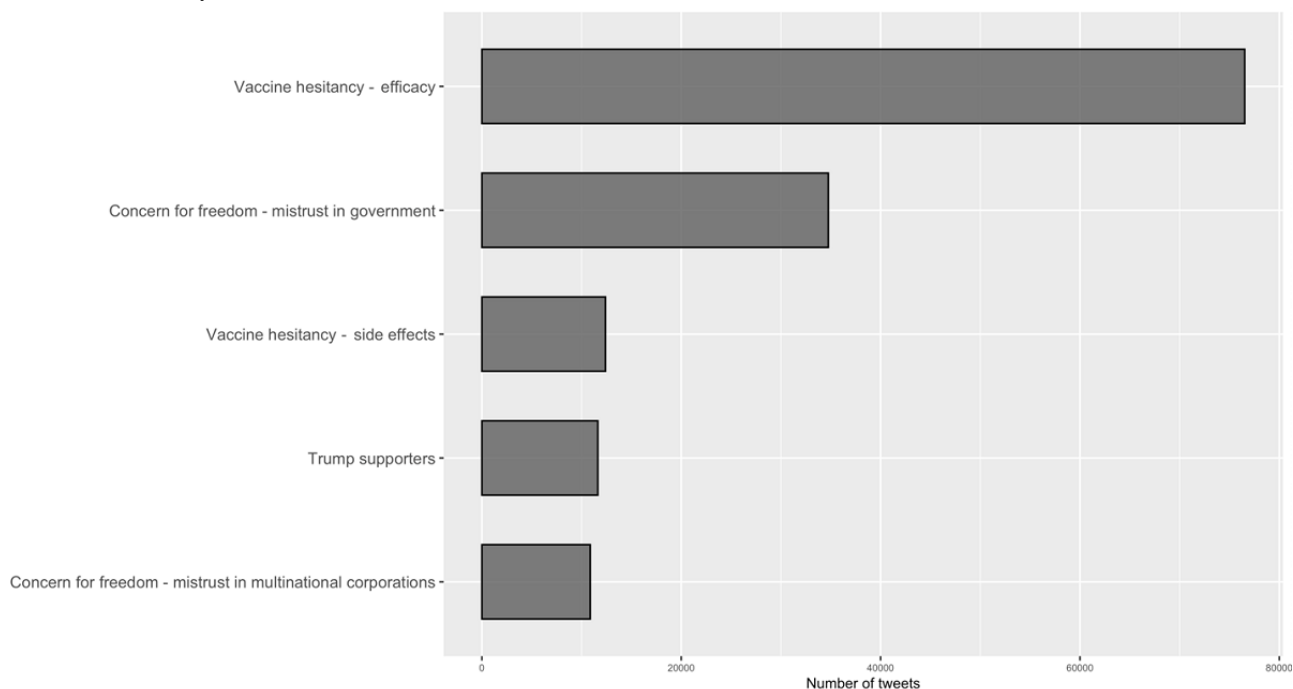
| Topic | Topic keywords bi-gram (tf-idf) | Inferred topic | Tweets (N=146,191), n (%) |
|---|---|---|---|
| 1 | Countries test, results multinational, multinational companies, the Johnson, commie, passes two, proposing freedom, commie proposing, full commie, Johnson government, gone full, government gone | Trust in the government | 33,578 (23.0%) |
| 2 | Tanked economy, tweet day, bell, disturbing tweet, day professor, most disturbing, bell talking, professor sir, irving bell, john irving, irving, sir john | COVID-19 vaccine hesitancy: Efficacy | 30,576 (20.9%) |
| 3 | Produced cell, unlicensed produced, lines aborted, shots wont, HIV/AIDS, won't walk, test enough, infection enough, antibodies past, enough effective, enough antibodies, past infection | COVID-19 vaccine hesitancy: Efficacy | 28,818 (19.7%) |
| 4 | Presi, COVID literally, one two, presi im, corner presi, literaly around, im sitting, sitting thinking, thinking incredible, incredible one | Support for Trump's management of the COVID-19 crisis | 11,631 (8.0%) |
| 5 | nick healthy, jer, emtfire, gau nick, jer gau, damage jer, dayprofessor sir, tweet dayprofessor, dayprofessor | COVID-19 vaccine hesitancy: Side effects | 9579 (6.6%) |
| 6 | Passports we, elite continue, the elite, commercial flights, fly private, breaking qantas, ceo confirms, you've vaccinated, compulsory international, international, confirms proof, proof you've | Trust in multinational corporations | 9286 (6.4%) |
| 7 | Attacked, attacked since, since apr, apr highlighting, highlighting imp, viciously, imp tcell, sarscov, despite published, published COVID, viciously attacked, ive viciously | COVID-19 vaccine hesitancy: Efficacy | 7896 (5.4%) |
| 8 | Label new, staff form, group label, jabs poison, longterm tcell, immunity cases, watch interview, interview dr, science longterm, discussing hysteria, phd discussing, dr phd. | COVID-19 vaccine hesitancy: Efficacy | 6921 (4.7%) |
| 9 | Arn caca, quand, faire, arn, bonjour, confinement vaccin, avec, caca, confinement, vaccin contre, contre le, le vaccin | COVID-19 vaccine hesitancy: Side effects | 2822 (1.9%) |
| 10 | Coronavirus im, pr quarter, newspapers, im delighted, delighted mainstream, mainstream newspapers, newspapers picking, picking pr. | COVID-19 vaccine hesitancy: Efficacy | 2320 (1.6%) |
| 11 | Fly must, want fly, begins if, must take | Trust in multinational corporations | 1586 (1.1%) |
| 12 | Protestors, force in, could force, COVID law, authorities could, Denmark, proposed forced, protesting proposed, protestors protesting, in Denmark, Denmark protestors, law authorities | Trust in the government | 1178 (0.8%) |

A final residual category (topic 4) in the vaccine hesitancy conversation singled out a tweet by radio host Rush Limbaugh mentioning that two COVID-19 vaccines were now approved and praising President Trump's management of the crisis. This tweet was widely circulated in the vaccine hesitant conversation cluster and underlines the reality that any conversation on COVID-19 immunization, and to some degree, vaccine hesitancy clusters on social media, intersect with broader clusters structured by current political polarization.

Figure 3 presents the approximate distribution of these inferred topics in COVID-19 vaccine hesitancy clusters identified during this observation period and allows us to understand what shapes an individual's attitude toward a COVID-19 vaccine. First, as this data set was collected right after the news of positive trials of COVID-19 vaccines by Pfizer and Moderna, more than half of the tweets (n=76,531, 52.4%) mentioned vaccine efficacy and raised suspicions on whether creating vaccines using mRNA technology was achievable or whether immunization would be long-lasting. With respect to vaccine safety, only 8.5% (n=12,401) of the social media debate in vaccine hesitancy clusters doubted its safety. Second, nearly one-third (n=45,628, 31.2%) of conversations on vaccine hesitancy clusters on Twitter expressed concerns for freedom or mistrust of institutions (either the government or multinational corporations). These results suggest that one key determinant of vaccine hesitancy is trust in institutions. It suggests that vaccine confidence building is a problem of shaping attitudes toward COVID-19 vaccines, which falls into the realm of health policy, as well as promoting social/political trust, which falls more into the realm of politics.

XSL•FO

RenderX

**Figure 3.** Distribution of inferred topics in COVID-19 vaccine hesitancy clusters on social media. This figure presents the aggregated results of inferred topics in vaccine hesitancy clusters from Table 1.



## Discussion

### Principal Findings

We analyzed more than 600,000 tweets over 1 week just following the announcement of successful COVID-19 vaccine trials, to characterize and understand global public perceptions and attitudes surrounding the COVID-19 vaccine and themes driving vaccine hesitancy. Our analysis revealed contrasting conversations about COVID-19 immunization on social media with both vaccine acceptant and vaccine hesitant clusters. Identified were the main themes driving the vaccine hesitant conversation at that time, including concerns of safety, efficacy, and freedom, and mistrust in institutions (either the government or multinational corporations). A main theme was the safety and efficacy of mRNA technology and side effects. The conversation around efficacy was that vaccines were unlikely to completely rid the population of COVID-19, polymerase chain reaction testing is flawed, and there is no indication of long-term T-cell immunity for COVID-19. Nearly one-third (31.2%) of the conversations on COVID-19 vaccine hesitancy clusters expressed concerns for freedom or mistrust of institutions (either the government or multinational corporations) and nearly a quarter (23.8%) expressed criticism toward the government's handling of the pandemic.

The main themes identified in this study fall under the domain of confidence using the 5C scale. Confidence is a measure of a person's level of trust in vaccine safety and efficacy, as well as in those involved in vaccine administration including policy makers and health professionals. Studies have shown that confidence is positively correlated with attitudes toward vaccination, knowledge of vaccination, and trust in health care, while it is negatively correlated with conspiracy mentality and medical harms [13]

The speed of development, production, and mass rollout of a COVID-19 vaccine has been unprecedented and has led to concerns around the safety and efficacy of vaccination. A common theme identified at that time was concern regarding mRNA technology. mRNA-based therapeutics have been used for cancer vaccines in the past; however, compared to other vaccine technology, it has not been clinically tested to the same extent [40]. Dror et al published results of a survey in which 70% of the general public responded with concerns about the safety of the COVID-19 vaccine [41]. We identified that less than 10% of the social media debate in vaccine hesitancy conversations during our observation period doubted its safety. Nevertheless, as COVID-19 vaccines are rolled out, it is highly probable that antivaccine social media conversations will transition from arguing about efficacy to questioning vaccine safety.

Mistrust in institutions has emerged as a predominant theme of vaccine hesitancy conversations on Twitter. Prior literature has reported mistrust in doctors, government sources, and pharmaceutical companies as reasons for hesitancy [16]. Governments are directly involved in many aspects of vaccine development, from funding to eventual safety approval. Individuals who believe the government is incompetent or malicious may not trust that these functions have been carried out in an appropriate way. Trust in government covaries strongly with generalized trust in other people and feelings of connectedness to others in society [42-48]. These measures of "social capital" in turn have been linked with reduced willingness to contribute to public good [43]. For example, there is a link between government trust and willingness to pay taxes [49]. Conversely, increased ethnic or political fragmentation, which creates feelings of division in society, has been shown to reduce the quality of the government and reduce physical distancing during the COVID-19 pandemic [50,51]. Because

the COVID-19 vaccine has public benefits that go beyond individual protection, individuals with low social trust may be less willing to contribute to the public good by getting the vaccine. Furthermore, recent studies using survey data have been increasingly associating trust in the government with COVID-19 behavior and vaccine hesitancy, and our social media data add to this literature [52,53].

There is extensive literature examining both individual and aggregate correlates of trust. At the individual level, trust is positively correlated with education [54] and civic engagement [42]. Aggregate measures of social trust vary with the actual performance of the government, and poor economic growth, high crime, mass protests, and political scandals appear to reduce the trust of citizens in the government [55,56]. Conversely, increasing transparency in the government appears to improve public trust in authorities [57]. Highlighted is the importance of building trust in institutions, which needs to be incorporated into models aimed at targeting vaccine hesitancy in addition to the traditional pillars of communication, information, and cognition.

Globally, persons are challenged with an overabundance of information on COVID-19 and COVID-19 vaccination, in which misinformation has been disseminated rampantly, likely fueling hesitancy [58]. Yaqub et al highlighted in a critical review of vaccine hesitancy that "hesitancy is not a rare phenomenon or confined solely to antivaccinationists; it includes people who have not yet rejected vaccination. Focusing on only vaccine uptake rates and neglecting underlying attitudes is likely to underestimate the challenge of maintaining vaccination coverage in the future" [16]. We demonstrated that social media analysis provides insights into societal attitudes, communication trends, and barriers to vaccine uptake that must be considered when developing strategies to address vaccine hesitancy.

The strength of this study lies in the methodology undertaken, which involved a bottom-up approach for the identification of cases and SNA. Previous studies have generally adopted a top-down approach to data collection, isolating known antivaccine accounts and analyzing its content and diffusion. Although critical in understanding the structure and nature of antivaccination framing on social media, such methods run the risk of selection bias. Moreover, by analyzing tweets in both French and English, we were able to broaden the scope of our vaccine conversation clusters, increasing the generalizability of this work. The noise of this analysis was minimized by linking the vaccination keywords with "COVID."

## Limitations

There are several limitations in our work. Although social media is increasingly used as a source of information and social interaction, even on matters related to health policy, it is an environment where participants self-select themselves in the population and is not a representative sample of the general population. In this sense, our results reflect a specific conversation around COVID-19 vaccination, and studies focusing on other social media platforms (eg, Facebook, WhatsApp, Reddit, and YouTube) and more traditional news media would offer a more complete overview of how antivaccine narratives are structured. Demographic segmentation of these clusters was not possible as further background information on the individuals in each cluster was not available; however, this would be a valuable component to future research. Our analysis was done on data collected right at the onset of news confirming the successful clinical trials of COVID-19 vaccines by Pfizer and Moderna. Our research offers a baseline from which we can understand the evolution of the online debate about COVID-19 vaccination. However, our assumption is that such a conversation will evolve and change throughout the pandemic, and we should expect the saliency of antivaccine arguments (safety, efficacy, and trust in institutions) to fluctuate as new information and policies are put in place. Future research is needed to monitor the COVID-19 vaccine hesitancy conversation through adopting a dynamic approach by collecting tweets over longer time periods and analyzing the patterns of change over time. Finally, as data were collected through the Twitter streaming API, the sample we analyzed may not have been fully randomized. The Twitter streaming API tends to overrepresent central users and is influenced by Twitter's sampling algorithm. Furthermore, given the size of the data set, we could not remove bots, which may have potentially skewed certain results.

## Conclusions

The recent global rollout of COVID-19 vaccination has brought vaccine hesitancy to the forefront in managing this pandemic. Hesitancy in accepting COVID-19 vaccination is fundamentally different from other vaccinations due to the new technologies being used, rapid development, and widespread global distribution. Attitudes on vaccines are largely driven by online information, particularly information on social media. We demonstrated that social media content and network analysis provides insights into societal attitudes, communication trends, and barriers to vaccine uptake. Identified themes driving the vaccine hesitant conversation included concerns of safety, efficacy, and freedom, and mistrust in institutions (either the government or multinational corporations). These themes will need to be considered as targeted outreach programs and intervention strategies are deployed globally in attempts to change personal attitudes on Twitter and improve the uptake of COVID-19 vaccination.

## Authors' Contributions

JCB, RL, KC, JLB, TT, MM, MMF, CC, DAM, RJO, HH, and JH were involved in initial concepts, literature searches, and funding applications. JCB and AB obtained all data and performed all analyses. JCB, RL, and KC wrote the initial draft of the paper. All authors made edits and contributions to the final draft of this manuscript.

## Conflicts of Interest

None declared.

## References

1.  Coronavirus Resource Center. Johns Hopkins University & Medicine. URL: https://coronavirus.jhu.edu/ [accessed 2021-03-12]
2.  Dubé E, MacDonald NE. How can a global pandemic affect vaccine hesitancy? Expert Rev Vaccines 2020 Oct;19(10):899-901. [doi: 10.1080/14760584.2020.1825944] [Medline: 32945213]
3.  Lurie N, Saville M, Hatchett R, Halton J. Developing Covid-19 Vaccines at Pandemic Speed. N Engl J Med 2020 May 21;382(21):1969-1973. [doi: 10.1056/NEJMp2005630] [Medline: 32227757]
4.  Report of the SAGE Working Group on Vaccine Hesitancy. World Health Organization. 2014. URL: http://www.who.int/immunization/sage/meetings/2014/october/SAGE_working_group_revised_report_vaccine_hesitancy.pdf [accessed 2021-08-02]
5.  Ten Threats to Global Health in 2019. World Health Organization. URL: https://www.who.int/news-room/spotlight/ten-threats-to-global-health-in-2019 [accessed 2020-10-05]
6.  Schaffer DeRoo S, Pudalov NJ, Fu LY. Planning for a COVID-19 Vaccination Program. JAMA 2020 Jun 23;323(24):2458-2459. [doi: 10.1001/jama.2020.8711] [Medline: 32421155]
7.  Trujillo K, Motta M. A majority of vaccine skeptics plan to refuse a COVID-19 vaccine, a study suggests, and that could be a big problem. The Conversation. 2020. URL: https://theconversation.com/a-majority-of-vaccine-skeptics-plan-to-refuse-a-covid-19-vaccine-a-study-suggests-and-that-could-be-a-big-problem-137559 [accessed 2020-10-06]
8.  COVID-19 vaccine willingness among Canadian population groups. Statistics Canada. 2021. URL: https://www150.statcan.gc.ca/n1/pub/45-28-0001/2021001/article/00011-eng.htm [accessed 2021-08-02]
9.  McAteer J, Yildirim I, Chahroudi A. The VACCINES Act: Deciphering Vaccine Hesitancy in the Time of COVID-19. Clin Infect Dis 2020 Jul 28;71(15):703-705 [FREE Full text] [doi: 10.1093/cid/ciaa433] [Medline: 32282038]
10. Dubé E, Gagnon D, MacDonald N, Bocquier A, Peretti-Watel P, Verger P. Underlying factors impacting vaccine hesitancy in high income countries: a review of qualitative studies. Expert Rev Vaccines 2018 Nov;17(11):989-1004. [doi: 10.1080/14760584.2018.1541406] [Medline: 30359151]
11. de Figueiredo A, Simas C, Karafillakis E, Paterson P, Larson HJ. Mapping global trends in vaccine confidence and investigating barriers to vaccine uptake: a large-scale retrospective temporal modelling study. The Lancet 2020 Sep 26;396(10255):898-908 [FREE Full text] [doi: 10.1016/S0140-6736(20)31558-0] [Medline: 32919524]
12. Shapiro GK, Tatar O, Dube E, Amsel R, Knauper B, Naz A, et al. The vaccine hesitancy scale: Psychometric properties and validation. Vaccine 2018 Jan 29;36(5):660-667. [doi: 10.1016/j.vaccine.2017.12.043] [Medline: 29289384]
13. Betsch C, Schmid P, Heinemeier D, Korn L, Holtmann C, Böhm R. Beyond confidence: Development of a measure assessing the 5C psychological antecedents of vaccination. PLoS One 2018;13(12):e0208601 [FREE Full text] [doi: 10.1371/journal.pone.0208601] [Medline: 30532274]
14. Stecula DA, Kuru O, Albarracin D, Jamieson KH. Policy Views and Negative Beliefs About Vaccines in the United States, 2019. Am J Public Health 2020 Oct;110(10):1561-1563. [doi: 10.2105/AJPH.2020.305828] [Medline: 32816542]
15. Du J, Luo C, Shegog R, Bian J, Cunningham RM, Boom JA, et al. Use of Deep Learning to Analyze Social Media Discussions About the Human Papillomavirus Vaccine. JAMA Netw Open 2020 Nov 02;3(11):e2022025 [FREE Full text] [doi: 10.1001/jamanetworkopen.2020.22025] [Medline: 33185676]
16. Yaqub O, Castle-Clarke S, Sevdalis N, Chataway J. Attitudes to vaccination: a critical review. Soc Sci Med 2014 Jul;112:1-11 [FREE Full text] [doi: 10.1016/j.socscimed.2014.04.018] [Medline: 24788111]
17. 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;5(3):337-348. [doi: 10.1038/s41562-021-01056-1] [Medline: 33547453]
18. Abd-Alrazaq A, Alhuwail D, Househ M, Hamdi M, Shah Z. Top Concerns of Tweeters During the COVID-19 Pandemic: Infoveillance Study. J Med Internet Res 2020 Apr 21;22(4):e19016 [FREE Full text] [doi: 10.2196/19016] [Medline: 32287039]
19. Steffens MS, Dunn AG, Wiley KE, Leask J. How organisations promoting vaccination respond to misinformation on social media: a qualitative investigation. BMC Public Health 2019 Oct 23;19(1):1348 [FREE Full text] [doi: 10.1186/s12889-019-7659-3] [Medline: 31640660]

20. Sinnenberg L, Buttenheim AM, Padrez K, Mancheno C, Ungar L, Merchant RM. Twitter as a Tool for Health Research: A Systematic Review. Am J Public Health 2017 Jan;107(1):e1-e8. [doi: 10.2105/AJPH.2016.303512] [Medline: 27854532]

21. Tsao S, Chen H, Tisseverasinghe T, Yang Y, Li L, Butt ZA. What social media told us in the time of COVID-19: a scoping review. The Lancet Digital Health 2021 Mar;3(3):e175-e194 [FREE Full text] [doi: 10.1016/S2589-7500(20)30315-0] [Medline: 33518503]

22. Wilson SL, Wiysonge C. Social media and vaccine hesitancy. BMJ Glob Health 2020 Oct;5(10):e004206 [FREE Full text] [doi: 10.1136/bmjgh-2020-004206] [Medline: 33097547]

23. Tavoschi L, Quattrone F, D'Andrea E, Ducange P, Vabanesi M, Marcelloni F, et al. Twitter as a sentinel tool to monitor public opinion on vaccination: an opinion mining analysis from September 2016 to August 2017 in Italy. Hum Vaccin Immunother 2020 May 03;16(5):1062-1069 [FREE Full text] [doi: 10.1080/21645515.2020.1714311] [Medline: 32118519]

24. Johnson NF, Velásquez N, Restrepo NJ, Leahy R, Gabriel N, El Oud S, et al. The online competition between pro- and anti-vaccination views. Nature 2020 Jun;582(7811):230-233. [doi: 10.1038/s41586-020-2281-1] [Medline: 32499650]

25. Griffith J, Marani H, Monkman H. COVID-19 Vaccine Hesitancy in Canada: Content Analysis of Tweets Using the Theoretical Domains Framework. J Med Internet Res 2021 Apr 13;23(4):e26874 [FREE Full text] [doi: 10.2196/26874] [Medline: 33769946]

26. Lucas C, Nielsen RA, Roberts ME, Stewart BM, Storer A, Tingley D. Computer-Assisted Text Analysis for Comparative Politics. Polit. anal 2017 Jan 04;23(2):254-277. [doi: 10.1093/pan/mpu019]

27. Greene D, Cross JP. Exploring the Political Agenda of the European Parliament Using a Dynamic Topic Modeling Approach. Polit. Anal 2017 Mar 13;25(1):77-94. [doi: 10.1017/pan.2016.7]

28. Hung M, Lauren E, Hon ES, Birmingham WC, Xu J, Su S, et al. Social Network Analysis of COVID-19 Sentiments: Application of Artificial Intelligence. J Med Internet Res 2020 Aug 18;22(8):e22590 [FREE Full text] [doi: 10.2196/22590] [Medline: 32750001]

29. Himelboim I, Xiao X, Lee DKL, Wang MY, Borah P. A Social Networks Approach to Understanding Vaccine Conversations on Twitter: Network Clusters, Sentiment, and Certainty in HPV Social Networks. Health Commun 2020 May;35(5):607-615. [doi: 10.1080/10410236.2019.1573446] [Medline: 31199698]

30. Hussain A, Tahir A, Hussain Z, Sheikh Z, Gogate M, Dashtipour K, et al. Artificial Intelligence-Enabled Analysis of Public Attitudes on Facebook and Twitter Toward COVID-19 Vaccines in the United Kingdom and the United States: Observational Study. J Med Internet Res 2021 Apr 05;23(4):e26627 [FREE Full text] [doi: 10.2196/26627] [Medline: 33724919]

31. Luo X, Zimet G, Shah S. A natural language processing framework to analyse the opinions on HPV vaccination reflected in twitter over 10 years (2008 - 2017). Hum Vaccin Immunother 2019;15(7-8):1496-1504 [FREE Full text] [doi: 10.1080/21645515.2019.1627821] [Medline: 31194609]

32. Polack FP, Thomas SJ, Kitchin N, Absalon J, Gurtman A, Lockhart S, C4591001 Clinical Trial Group. Safety and Efficacy of the BNT162b2 mRNA Covid-19 Vaccine. N Engl J Med 2020 Dec 31;383(27):2603-2615 [FREE Full text] [doi: 10.1056/NEJMoa2034577] [Medline: 33301246]

33. Jackson LA, Anderson EJ, Rouphael NG, Roberts PC, Makhene M, Coler RN, mRNA-1273 Study Group. An mRNA Vaccine against SARS-CoV-2 - Preliminary Report. N Engl J Med 2020 Nov 12;383(20):1920-1931 [FREE Full text] [doi: 10.1056/NEJMoa2022483] [Medline: 32663912]

34. Lyu JC, Han EL, Luli GK. COVID-19 Vaccine-Related Discussion on Twitter: Topic Modeling and Sentiment Analysis. J Med Internet Res 2021 Jun 29;23(6):e24435 [FREE Full text] [doi: 10.2196/24435] [Medline: 34115608]

35. Liu S, Liu J. Understanding Behavioral Intentions Toward COVID-19 Vaccines: Theory-Based Content Analysis of Tweets. J Med Internet Res 2021 May 12;23(5):e28118 [FREE Full text] [doi: 10.2196/28118] [Medline: 33939625]

36. Allington D, McAndrew S, Moxham-Hall V, Duffy B. Coronavirus conspiracy suspicions, general vaccine attitudes, trust and coronavirus information source as predictors of vaccine hesitancy among UK residents during the COVID-19 pandemic. Psychol Med 2021 Apr 12:1-12 [FREE Full text] [doi: 10.1017/S0033291721001434] [Medline: 33843509]

37. Yang Z, Guo J, Cai K, Tang J, Li J, Zhang L, et al. Understanding retweeting behaviors in social networks. In: CIKM '10: Proceedings of the 19th ACM International Conference on Information and Knowledge Management. 2010 Presented at: 19th ACM International Conference on Information Knowledge Management; October 26-30, 2010; Toronto, ON, Canada p. 1633-1636. [doi: 10.1145/1871437.1871691]

38. Blondel VD, Guillaume J, Lambiotte R, Lefebvre E. Fast unfolding of communities in large networks. J. Stat. Mech 2008 Oct 09;2008(10):P10008. [doi: 10.1088/1742-5468/2008/10/P10008]

39. Sanh V, Debut L, Chaumond J, Wolf T. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. arXiv. 2020. URL: https://arxiv.org/abs/1910.01108 [accessed 2021-08-05]

40. Abbasi J. COVID-19 and mRNA Vaccines-First Large Test for a New Approach. JAMA 2020 Sep 22;324(12):1125-1127. [doi: 10.1001/jama.2020.16866] [Medline: 32880613]

41. Dror AA, Eisenbach N, Taiber S, Morozov NG, Mizrachi M, Zigron A, et al. Vaccine hesitancy: the next challenge in the fight against COVID-19. Eur J Epidemiol 2020 Aug;35(8):775-779. [doi: 10.1007/s10654-020-00671-y] [Medline: 32785815]

42. Keele L. Social Capital and the Dynamics of Trust in Government. Am J Political Science 2007 Apr;51(2):241-254. [doi: 10.1111/j.1540-5907.2007.00248.x]

43. Putnam R. Bowling Alone: The Collapse and Revival of American Community. New York, NY: Simon & Schuster; 2000.

XSL•FO
RenderX

44.  Nannestad P. What Have We Learned About Generalized Trust, If Anything? Annu. Rev. Polit. Sci 2008 Jun;11(1):413-436. [doi: 10.1146/annurev.polisci.11.060606.135412]

45.  Rothstein B, Stolle D. The State and Social Capital: An Institutional Theory of Generalized Trust. Comp Politics 2008 Jul 01;40(4):441-459. [doi: 10.5129/001041508x12911362383354]

46.  Knack S, Keefer P. Does Social Capital Have an Economic Payoff? A Cross-Country Investigation. The Quarterly Journal of Economics 1997 Nov 01;112(4):1251-1288. [doi: 10.1162/003355300555475]

47.  Knack S. Social Capital and the Quality of Government: Evidence from the States. American Journal of Political Science 2002 Oct;46(4):772-785. [doi: 10.2307/3088433]

48.  Jennings MK, Stoker L. Social Trust and Civic Engagement across Time and Generations. Acta Polit 2004 Dec 13;39(4):342-379. [doi: 10.1057/palgrave.ap.5500077]

49.  Antinyan A, Corazzini L, Pavesi F. Does trust in the government matter for whistleblowing on tax evaders? Survey and experimental evidence. Journal of Economic Behavior & Organization 2020 Mar;171:77-95. [doi: 10.1016/j.jebo.2020.01.014]

50.  Alesina A, Zhuravskaya E. Segregation and the Quality of Government in a Cross Section of Countries. American Economic Review 2011 Aug 01;101(5):1872-1911. [doi: 10.1257/aer.101.5.1872]

51.  Algan Y, Hémet C, Laitin DD. The Social Effects of Ethnic Diversity at the Local Level: A Natural Experiment with Exogenous Residential Allocation. Journal of Political Economy 2016 Jun;124(3):696-733. [doi: 10.1086/686010]

52.  Merkley E, Loewen PJ. Anti-intellectualism and the mass public's response to the COVID-19 pandemic. Nat Hum Behav 2021 Jun 28;5(6):706-715. [doi: 10.1038/s41562-021-01112-w] [Medline: 33911228]

53.  Lazarus JV, Ratzan SC, Palayew A, Gostin LO, Larson HJ, Rabin K, et al. A global survey of potential acceptance of a COVID-19 vaccine. Nat Med 2021 Feb;27(2):225-228 [FREE Full text] [doi: 10.1038/s41591-020-1124-9] [Medline: 33082575]

54.  Brehm J, Rahn W. Individual-Level Evidence for the Causes and Consequences of Social Capital. American Journal of Political Science 1997 Jul;41(3):999-1023. [doi: 10.2307/2111684]

55.  Chanley VA, Rudolph TJ, Rahn WM. The origins and consequences of public trust in government: a time series analysis. Public Opin Q 2000;64(3):239-256. [doi: 10.1086/317987] [Medline: 11114267]

56.  Sangnier M, Zylberberg Y. Protests and trust in the state: Evidence from African countries. Journal of Public Economics 2017 Aug;152:55-67. [doi: 10.1016/j.jpubeco.2017.05.005]

57.  Buell RW, Norton MI. Surfacing the Submerged State with Operational Transparency in Government Services. SSRN Journal 2013:1-24. [doi: 10.2139/ssrn.2349801]

58.  How to report misinformation online. World Health Organization. URL: https://tinyurl.com/3mh9bpp2 [accessed 2021-08-02]

## Abbreviations

**API:** application programing interface
**BERT:** Bidirectional Encoder Representations from Transformers
**SNA:** social network analysis

<u>Original Paper</u>

# Monitoring Depression Trends on Twitter During the COVID-19 Pandemic: Observational Study

Yipeng Zhang[1], BSc; Hanjia Lyu[1*], MSc; Yubao Liu[1*], MSc; Xiyang Zhang[2], MA; Yu Wang[1], PhD; Jiebo Luo[1], PhD

[1]University of Rochester, Rochester, NY, United States

[2]University of Akron, Akron, OH, United States

[*]these authors contributed equally

**Corresponding Author:**
Jiebo Luo, PhD
University of Rochester
500 Joseph C Wilson Blvd
Rochester, NY
United States
Phone: 1 585 276 5784
Email: jluo@cs.rochester.edu

## *Abstract*

**Background:**   The COVID-19 pandemic has affected people's daily lives and has caused economic loss worldwide. Anecdotal evidence suggests that the pandemic has increased depression levels among the population. However, systematic studies of depression detection and monitoring during the pandemic are lacking.

**Objective:**   This study aims to develop a method to create a large-scale depression user data set in an automatic fashion so that the method is scalable and can be adapted to future events; verify the effectiveness of transformer-based deep learning language models in identifying depression users from their everyday language; examine psychological text features' importance when used in depression classification; and, finally, use the model for monitoring the fluctuation of depression levels of different groups as the disease propagates.

**Methods:**   To study this subject, we designed an effective regular expression-based search method and created the largest English Twitter depression data set containing 2575 distinct identified users with depression and their past tweets. To examine the effect of depression on people's Twitter language, we trained three transformer-based depression classification models on the data set, evaluated their performance with progressively increased training sizes, and compared the model's tweet chunk-level and user-level performances. Furthermore, inspired by psychological studies, we created a fusion classifier that combines deep learning model scores with psychological text features and users' demographic information, and investigated these features' relations to depression signals. Finally, we demonstrated our model's capability of monitoring both group-level and population-level depression trends by presenting two of its applications during the COVID-19 pandemic.

**Results:**   Our fusion model demonstrated an accuracy of 78.9% on a test set containing 446 people, half of which were identified as having depression. Conscientiousness, neuroticism, appearance of first person pronouns, talking about biological processes such as eat and sleep, talking about power, and exhibiting sadness were shown to be important features in depression classification. Further, when used for monitoring the depression trend, our model showed that depressive users, in general, responded to the pandemic later than the control group based on their tweets (n=500). It was also shown that three US states—New York, California, and Florida—shared a similar depression trend as the whole US population (n=9050). When compared to New York and California, people in Florida demonstrated a substantially lower level of depression.

**Conclusions:**   This study proposes an efficient method that can be used to analyze the depression level of different groups of people on Twitter. We hope this study can raise awareness among researchers and the public of COVID-19's impact on people's mental health. The noninvasive monitoring system can also be readily adapted to other big events besides COVID-19 and can be useful during future outbreaks.

**KEYWORDS**

mental health; depression; social media; Twitter; data mining; natural language processing; transformers; COVID-19

XSL•FO
**RenderX**

# Introduction

## Background

COVID-19 is an infectious disease that has been spreading rapidly worldwide since early 2020. It was first identified on December 31, 2019, and was officially declared as a pandemic by the World Health Organization on March 11, 2020 [1]. As of September 15, 2020, COVID-19 has infected 216 countries, areas, or territories with over 29 million confirmed cases and 930,000 confirmed deaths [1]. In response to the pandemic, over 190 countries have issued nationwide closures of educational facilities [2], and many governments have issued flight restrictions and stay-at-home-orders, affecting the everyday lives of people worldwide.

Mental disorders were affecting approximately 380 million people of all ages worldwide before COVID-19 [3]. Previous psychological studies have shown that mental disorders lead to many negative outcomes including suicide [4,5]. However, these studies face two challenges. First, it is known that individuals with mental disorders are sometimes unwilling or ashamed to seek help [6]. Second, it is oftentimes infeasible for psychological studies to obtain and track a large sample of diagnosed individuals and perform statistically significant numerical analysis.

Multiple studies have investigated the economic and social impacts of COVID-19 [7,8], and various studies have shown that COVID-19 has greatly impacted people's mental health worldwide. These studies found that there are higher rates of depression, anxiety, posttraumatic stress disorder (PTSD), and stress symptoms reported during COVID-19 than before [9]. Females, young age groups, students, and low education groups are especially susceptible to depression during the pandemic [9]. The pandemic negatively affected individuals' mental health because of the changes that it brought to life. For example, it has been shown that after nationwide lockdowns people experienced high levels of stress because of social isolation [10]; the fact that a large proportion of the population is not wearing masks also makes people experience high levels of anxiety and depression [11]. For individuals with mental disorders, their need is amplified; the study by Hao et al [12] suggests that, during the pandemic, psychiatric patients reported more moderate to severe anger and impulsivity as well as concerns about their physical health, as opposed to the healthy controls, and that ideal remote mental health services such as telepsychiatry consultation and home delivery of medications could not be fully established due to the sudden lockdown [12].

Given this pressing situation, we would like to quantify mental health conditions of the general population during the pandemic. Nevertheless, the data source selection is critical for overcoming the two challenges mentioned previously. In the past decade, people have been increasingly relying on social media platforms such as Facebook, Twitter, and Instagram to express their feelings. Social media can thus serve as a resourceful medium for mining information about the public's mental health conditions [13-17]. The public have long been known to search online for information about diseases and medical issues [18]. COVID-19 is no exception. Indeed, using social media, public opinions on personal face mask use [19] and COVID-19 vaccine uptake [20,21] have been investigated. Existing research has also studied the predictive power of online medical consultation, online medical appointment, and online medical search in forecasting regional outbreaks and found online medical consultation to be the most predicative [22]. Furthermore, a recent longitudinal study on the mental health of the Chinese population during the pandemic has found that dissemination of health information via radio was associated with higher levels of anxiety and depression, and suggested television and the internet as alternatives [23]. Therefore, we believe social media platforms like Twitter offer a solution to the challenges, as they enable us to perform a large-scale quantitative study on mental disorders in a noninvasive fashion.

As shown in Figure 1, we used data from the ForSight by Crimson Hexagonplot [24] to plot the word frequencies of several mental disorders on Twitter, including "depression," "PTSD," "bipolar disorder," and "autism," from January 1 to May 4, 2020. Note that we excluded false-positive tweets that contained misleading phrases such as "economic depression" or "great depression." We noticed a rapid growth of the word frequencies of these mental disorders starting from March 17, when the pandemic spread across most of the world. Past research has suggested that depression is more pervasive than other psychological disorders during the COVID-19 period [9]. Similarly, we found that the word "depression" occurs substantially more frequently on Twitter compared to the other three mental disorders. Accordingly, depression is likely to be triggered most frequently by COVID-19, and we focused on understanding COVID-19's impact on depression in this study.

**Figure 1.** Density of Twitter coverage regarding "depression," "ptsd," "bipolar disorder," and "autism." ptsd: posttraumatic stress disorder.



## Prior Work

The potential of machine learning models for identifying Twitter users who have been diagnosed with depression was pioneered by De Choudhury et al [25], who analyzed how features obtained by Linguistic Inquiry and Word Count (LIWC) were related to depression signals on social media and how that can be used for user-level classification on a data set containing 171 depression users. The data was collected by designing surveys for volunteers through crowdsourcing. Following this work, Coppersmith et al [26] used LIWC, 1-gram language model, character 5-gram model, and user's engagement on social media (user mention rate, tweet frequency, etc) to perform tweet-level classification on a data set containing 441 depression users.

The CLPsych 2015 Shared Task data set containing 447 diagnosed depression users [27] was published in 2015 and was favored by a wide range of studies [28-30]. The data was gathered by regular expression search in tweets in combination with manual annotation. Among these studies, the performance of traditional machine learning classification algorithms (decision trees, support vector machines [SVMs], naive Bayes, logistic regression) on 1-grams and 2-grams was investigated by Nadeem [30]; Jamil et al [28] used SVM on bag of words (BOW) and depression word count along with LIWC features and NRC sentiment features; Orabi et al [29] explored the performance of small deep neural network [architectures]—one-dimensional convolutional neural network (CNN) and bidirectional long short-term memory (BiLSTM) with context-aware attention—and achieved the best performance (87% accuracy) on the task.

The CLPsych 2019 Shared Task [31] focused on evaluating Reddit users' suicide risk based on their posts, for which Matero et al [32] applied a pretrained Bidirectional Encoder Representations from Transformers (BERT) [33] embedding to encode the data. Suicide risk assessment on Spanish tweets was also studied by Ramírez-Cifuentes et al [34]. We argue that our task is different since few detected depressive Twitter users express suicide intent, while all the positive suicidal users in the suicide risk data sets should be viewed as in late stages of depression [35,36]. There are also some studies that performed depression detection on Reddit users [37-39] with sample sizes of less than 1300 Reddit posts. By contrast, we used the transformer-based models in our study, which have been shown to achieve state-of-the-art results in a wide range of natural language processing problems [33,40,41].

In addition to these two challenge data sets, several studies attempted to gather their own data of various forms. Tsugawa et al [42] performed analysis of models using BOW, latent Drichlet allocation (LDA) [43], and social media engagement features on a data set containing 81 Japanese-speaking depression Twitter users collected by crowdsourcing. Zhou et al [44] used ubiquitous multimodal sensors and performed in-depth analysis on users' social media content, social network, webcam video, and user interaction on a sample of 5 depression users. Detecting depression from Spanish tweets using sentiment and emotion lexicons was used by Leis et al [45]. Zhang et al [46] performed observational analysis of the relationship between deteriorating depression and behavior changes when engaging with Google search and YouTube on 49 depressive college students. Shen et al [47] proposed a multimodal dictionary learning method that used topic, social media

engagement, profile image, and emotional features to learn a latent feature dictionary that performed well on a data set of 1402 users with depression, the largest Twitter depression data set used to the best of our knowledge. Given the skyrocketing word density of "depression" in Figure 1, we show that a substantially larger depression data set can be quickly constructed from the COVID-19–related tweets within several months.

## Goal of the Study

Although the time series plots of keyword frequencies in Figure 1 offer an intuitive reading of depression's general trend in the population, they are apparently filled with noise and lack plausible explanation to be an accurate representation. To generalize beyond keywords, we would like to train machine learning–based models to identify depression on social media. Reddit automatically gathers posts of the same topic into "subreddits"; however, as pointed out by Pirina and Çöltekin [38], labeling posts completely according to subreddit names causes categories to be topically specific and cannot be generalized to regular social media text. Moreover, depression prediction models can potentially be used on the population level [48], but none of the work mentioned in the previous section applied their models to the general Twitter population on the fly.

Therefore, the main objectives of this study are to develop a method to create a large-scale depression user data set in an automatic fashion so that the method is scalable and can be adapted to future events; to verify the effectiveness of transformer-based deep learning language models in identifying depression users from their everyday language; to further improve the depression classification model using explainable psychological text features and to examine their importance in classification; and, finally, to use the model for monitoring the fluctuation of depression levels of different groups as the disease propagates.

# Methods

## Data Collection

First, we identified users with depression from 41.3 million COVID-19–related tweets posted by about 36.6 million users from March 23 to April 18, 2020. We collected the COVID-19–related tweets using the keywords "corona," "covid19," "covid_19," "coronavirus," "#Corona," "#Covid_19," and "#coronavirus." From these tweets, we looked for signals that can tell whether the user has depression from both the text and the user profile description.

Empirically, we observed that many Twitter users with depression described themselves as "depression fighters" in their descriptions. Some of them may also post relevant tweets to declare that they have been diagnosed with depression. Inspired by Coppersmith et al [26], we used regular expressions to find these authors by examining their tweets and descriptions. Building upon their method, we further extended our regular expression search based on some patterns we noticed on manually identified depression users, in pursuit of efficacy. In tweets, we searched for phrases such as "I

have/developed/got/suffer(ed) from X depression," "my X depression," "I'm healing from X depression," and "I'm diagnosed with X depression," where X is a descriptive word such as "severe" and "major" (X can be empty as well). In descriptions, we further added phrases such as "depression fighter/sufferer/survivor" to the regular expression list; we removed users that had "practitioner" and "counselor" in their descriptions to exclude mental health practitioners. The remaining users captured by the regular expressions were considered to have depression.

In the end, 2575 distinct Twitter users were classified into the depression group. Of 200 randomly sampled users in the depression set, 86% were labeled positive by human annotators. We randomly selected another 2575 distinct users so that depression-related terms did not appear in their past 200 tweets or descriptions as our control group. Users in this group were not considered to have depression (nondepression group). Once we found the targeted Twitter users, we used the Tweepy application programming interface (API) to retrieve the public tweets posted by these users within the last 3 months since the time of posting the depression-related tweet, with a maximum of 200 tweets per user. We chose 200 tweets because, on average, it is roughly the number of tweets posted by an individual within a 3-month time span, which is the length commonly adopted by previous work [25,26]. If a user was identified from the description, we limited the time scope from January 18 to April 18, 2020.

## Data Analysis

### Personality

Previous psychological research has shown that the big five personality traits (openness, conscientiousness, extraversion, agreeableness, and neuroticism) are related to depression [49,50]. In particular, low extraversion, high neuroticism, and low conscientiousness were associated with depressive symptoms [50]. We estimated individuals' personality scores using IBM's Personality Insights service [51]. For each individual, we aggregated all their tweets into a single textual input and used the Personality Insights API to obtain the scores. The minimum number of words for using the API was 100, and we were able to retrieve 4697 (91.2%) of the 5150 users' scores. Summary statistics are shown in Multimedia Appendix 1.

### Sentiments

Besides personality, we hypothesized that individuals' sentiments and emotions could also reflect whether they were experiencing depression or not. Sentiment analysis is widely-used in deciphering people's health and well-being from text data [52]. We estimated individuals' sentiments using the Valence Aware Dictionary and Sentiment Reasoner (VADER). VADER is a lexicon and rule-based model developed by researchers from the Georgia Institute of Technology [53]. We aggregated a user's tweets into a single chunk, applied VADER, and retrieved its scores for positive and negative emotions. In Figure 2, we reported the VADER score distributions of positive emotions and negative emotions among the depression and nondepression groups. Compared with individuals with no

depression, those with depression tended to exhibit both stronger

positive and negative emotions.

**Figure 2.** Distributions of positive and negative emotion scores among the depression and nondepression groups. VADER: Valence Aware Dictionary for Sentiment Reasoning.



### Demographics

Previous psychological studies have shown differences in depression rates among people of different ages and of different genders [54-56]. Research has shown a U-shaped relationship between age and depression, with depression reaching its lowest level around the age of 45 years [54]. Women are known to be substantially more likely to have depression [57]. To estimate the age and gender of the user, we adopted the M3-inference model proposed by Wang et al [58]. The M3 model performs multimodal analysis on a user's profile image, username, and description. Following M3's structure, we labeled each user with a binary gender label (as approximation) and a one-hot age label among four age intervals (≤18 years, 19-29 years, 30-39 years, ≥40 years), which were then used in our fusion model. Of the 5150 users, we were able to retrieve 5059 (98.2%) users' demographic information.

### Linguistic Inquiry Word Count

We used LIWC—a well-validated psycholinguistic dictionary [59]—to capture people's psychological states by analyzing the contents of their tweets. LIWC is a dictionary-based linguistic analysis tool that can count the percentage of words that reflect different emotions, thinking styles, and social concerns, and captures people's psychological states. Zhang et al [60] applied LIWC to the tweets of US working adults to analyze the influence of COVID-19 on their well-being; some LIWC features in college students' YouTube and Google search logs have been shown to correlate with their Patient Health Questionnaire-9 depression scores [46]; Coppersmith et al [26] showed the relationship between the use of the first person pronoun (which is one of the LIWC features) and depression [26].

We chose 8 features that were analyzed in previous works [26,61,62] and 7 other features that we found relevant to our study. Similar to the methods of Chen et al [63], we then applied LIWC to the concatenated tweets of individuals. Figure 3 shows the linguistic profiles for the tweets of the depression and nondepression groups. Both the depression and nondepression groups exhibited slightly positive tones, with negligible differences. The tweets of the nondepression group showed more analytical thinking, more clout, and less authentic expression than those of the depression group. The tweets of the depression group scored higher in both positive and negative emotion categories than the ones of the nondepression groups, which suggests a higher degree of immersion [64]. Moreover, the tweets of the depression group also showed more anxiety and anger emotions, and included more *swear* words—the *anxiety*, *anger*, and *swear* scores of the depression group were 50%, 22%, and 45% higher than that of the nondepression group, respectively—which is consistent with the findings of Coppersmith et al [26]. Death-related words appeared more frequently in the tweets of the depression group, which echoes Stirman and Pennebaker [62]. Similar to these 2 studies, we found more first person singular usage in the tweets of the depression group.

We also found that the tweets of the depression group expressed more sadness emotion and used words related to the biological process more frequently. Although there is no clear link between biological process–related words and depression, this finding shows that people with depression may pay more attention to their biological statuses. The *power* score for the tweets of the nondepression group was higher, which reflects a higher need for the power according to the findings of McClelland [65]. By comparing the *work* scores of the depression and nondepression groups, we found that the users of the nondepression group paid more attention to work-related issues as well.

**Figure 3.** Linguistic profiles for the depression and nondepression tweets. LIWC: Linguistic Inquiry and Word Count.



## Social Media Engagement

We used the proportion of tweets with mentions, number of responses, unique user mentions, user mentions, and tweets to measure the social media engagement of each user, as did Coppersmith et al [26]. To better understand the difference of social media engagement between the depression and nondepression groups, we added 0.1 to the number of responses, unique users mentions, users mentions, and tweets, and took the logarithm. By applying the Mann-Whitney rank test, we found that, except for the number of unique user mentions, other features were statistically different ($P<.05$) between the depression and nondepression groups. The users of the depression group posted more tweets and replied more. They tended to post fewer tweets with mentions, while the number of mentions for the depression group was larger, which suggests that when users of the depression group posted tweets to interact with other users, it involved more users.

## Modeling

### Task Definition

We formulated our task as a classification task, where the model was trained to predict whether a particular tweet or a chunk of tweets comes from a user from the depression set. Note that not all tweets by people in the depression set were explicitly referring to depression per se. By definition, though, they were all posted by users with depression and were thus labeled true. To help improve the model's generalizability, during training and testing, we excluded all the tweets used to identify the users with depression by regular expressions that contained trivial patterns and keywords. We assumed there were subtle differences in the language used between the depression and nondepression groups. Our goal was to build a model capable of capturing these subtleties and classifying users correctly.

## Tweet Chunking and Preprocessing

We performed stratified random sampling on our data set. We first sampled 500 users to form our testing set. On the rest of the users, we progressively added users to the training sets and recorded the performance of the models trained on sets of 1000, 2000, and 4650 users. All the training and testing sets have a 1:1 (depression:nondepression) ratio.

Jamil et al [28] have shown that one single tweet does not contain enough signals to determine whether a user has depression. Thus, we concatenated consecutive tweets of the same user together to create tweet chunks of 250 words and labeled the chunks based on the user's label. Given an input sentence, the transformer tokenizer first splits each word from the input sentence into *word-pieces* and then vectorizes them for computation. The 250 words roughly corresponded to the maximum 512 input word-pieces allowed by transformer-based language models including BERT [33] and Robustly Optimized BiLSTM Memory Pretraining Approach (RoBERTa) [40]. This limitation is due to the self-attention mechanism in the transformer, whose time complexity scales quadratically with the input sequence length.

We preprocessed the text using the tweet preprocessing pipeline proposed by Baziotis et al [66]. We adopted this method especially due to its capability of marking Twitter-specific text habits and converting them to special tokens such as "<allcaps>" (capitalized words), "<elongated>" (repeated letters), "<repeated>" (repeated words), etc. For example, "YESSSSS, I love it so much!!!" after preprocessing will be in the form of "Yes <allcaps> <elongated>, I love it so <repeated> much! <elongated>."

After chunking and preprocessing, on average, each user had 6-7 text chunks, making the actual sizes of the 4650-user train-validation set and the 500-user testing set to be 29,315

and 3105, respectively. The preprocessed tweet chunk data sets were then passed to deep learning models for training.

### Deep Learning Models

We used deep learning models to perform chunk-level classification. We set up two baseline models, multi-channel CNN and BiLSTM with context-aware attention (attention BiLSTM), as described in Orabi et al [29], which achieved the best performance on the CLPsych 2015 data set. We used the pretrained GloVe embedding (840B tokens, 300d vectors) [67] augmented with the special tokens added during preprocessing. The embedding weights were further trained jointly with the model. Recently, transformer-based deep learning language models have achieved state-of-the-art performance in multiple language modeling tasks. We trained three representative transformer-based sequence classification models—BERT [33], RoBERTa [40], and XLNet [41]—with their own pretrained tokenizers augmented with the special tokens for tokenization. We chose to use the base models for all of them since we found no noticeable performance gains using their larger counterparts.

### Signal Fusion

We ran the models on all the tweet chunks of the same user and took the average of the confidence scores to get the user-level confidence score. There were 4163 (89.5%) out of 4650 users remaining in the training set and 446 (89.2%) out of 500 users in the testing set whose entire features were retrievable. We then passed different combinations of user-level scores (personality, VADER, demographics, engagement, LIWC, and average confidence) to machine learning classification algorithms including random forest, logistic regression, and SVM provided by the *scikit-learn* library [68]. We only used the explainable LIWC features mentioned in the data collection section for training the classifiers.

### Training Details

During training, we randomly split the train-validation set to training and validation sets with a ratio of 9:1. We used Adam optimizer with a learning rate of 7e-3 and weight decay of 1e-4 for training attention BiLSTM. We used Adam optimizer with a learning rate of 5e-4 for training CNN. We used AdamW optimizer with a learning rate of 2e-5 for training BERT and RoBERTa, and 8e-6 for training XLNet. We used the cross-entropy loss for all our models during training. We used the stochastic gradient descent optimizer with adaptive learning rate, with initial learning rate as 0.1 for training SVM and logistic regression classifier. We recorded the models' performances on the validation set after each epoch and kept the model with the highest accuracy and F1 scores while training until convergence. We manually selected the hyperparameters that gave the best accuracy and F1 scores on the deep learning models.

## Results

### Chunk-Level Classification

In Table 1, we report our classification results at the chunk level on the testing set. Our evaluation metrics included accuracy, F1 score, area under the receiver operating characteristic curve (AUC), precision, and recall. One immediate observation was that, regardless of the model type, the classification performance improved as we increased the size of our train-validation set. This shows that for building depression classification models it is imperative to have a large number of training samples. At the same time, it also confirms that the larger number of training samples in our experiments was indeed an advantage.

Another observation was the performance gain of transformer-based models over BiLSTM and CNN models. The CNN model slightly outperformed BiLSTM, which replicated the findings of Orabi et al [29]. We observed that BERT, RoBERTa, and XLnet invariably outperformed BiLSTM and CNN regardless of the size of our training set. In particular, the XLNet model recorded the best AUC and accuracy of all the models when trained with our full training set.

**Table 1.** Chunk-level performance (%) of all 5 models on the 500-user testing set using training-validation sets of different sizes.[a]

| Model and training-validation set | Accuracy | F1 | AUC[b] | Precision | Recall |
|---|---|---|---|---|---|
| **Attention BiLSTM[c]** | | | | | |
| 1000 users | 70.7 | 69.0 | 76.5 | 70.9 | 67.3 |
| 2000 users | 70.3 | 68.3 | 77.4 | 70.7 | 66.1 |
| 4650 users | 72.7 | 71.6 | 79.3 | 72.1 | 71.1 |
| **CNN[d]** | | | | | |
| 1000 users | 71.8 | 72.6 | 77.4 | 72.7 | 72.6 |
| 2000 users | 72.8 | 74.5 | 80.3 | 72.2 | 76.9 |
| 4650 users | 74.0 | 70.9 | 81.0 | 77.4 | 68.9 |
| **BERT[e]** | | | | | |
| 1000 users | 72.7 | 74.4 | 79.8 | 72.0 | 76.9 |
| 2000 users | 75.7 | 76.3 | 82.9 | 76.1 | 75.7 |
| 4650 users | 76.5 | 77.5 | 83.9 | 76.3 | 78.8 |
| **RoBERTa[f]** | | | | | |
| 1000 users | 74.4 | 75.7 | 82.0 | 74.2 | 77.3 |
| 2000 users | 75.9 | 77.9 | 83.2 | 73.8 | *82.5* [g] |
| 4650 users | 76.2 | *78.0* | 84.1 | 74.4 | 81.9 |
| **XLNet** | | | | | |
| 1000 users | 73.7 | 75.1 | 80.7 | 73.2 | 77.2 |
| 2000 users | 74.6 | 76.8 | 82.6 | 72.6 | 81.5 |
| 4650 users | *77.1* | 77.9 | *84.4* | *77.5* | 78.3 |

[a]We used 0.5 as the threshold when calculating the scores.

[b]AUC: area under the receiver operating characteristic curve.

[c]BiLSTM: bidirectional long short-term memory.

[d]CNN: convolutional neural network.

[e]BERT: Bidirectional Encoder Representations from Transformers.

[f]RoBERTa: Robustly Optimized BiLSTM Pretraining Approach.

[g]Italics indicate the best performing model in each column.

## User-Level Classification

Next, we report our experiment results at the user level. Since XLNet trained on the 4650-user data set outperformed the other models, we took it for user-level performance comparison. Our experimental results demonstrated a substantial increase on the user-level scores of XLNet shown in Table 2 compared to the chunk-level score shown in Table 1. This indicates that more textual information of a user yields more reliable results on determining whether the user has depression. Building on the user-level XLNet scores, we further included VADER, demographic, engagement, personality, and LIWC scores as signals. We first used all features and compared the performance of random forest, logistic regression, and SVM. We noticed that SVM achieved the best scores on accuracy and F1, slightly surpassing logistic regression. Thus, we used SVM for testing the performance when using part of the features collected.

The results are shown in Table 2. The results have shown that using VADER, demographics, and social media engagement features alone does not help the classification by much. Classifiers using personality features and LIWC features perform relatively better. We then used these five feature groups and obtained a better result (accuracy 71.5%; F1 score 72.0%). However, the classifier was still outperformed by XLNet, showing that the transformer-based models indeed worked better on depressive Twitter text modeling compared with other approaches. We further increased the classifier's performance by using all the features, namely, VADER, demographics, engagement, personality, and LIWC features, and the averaged XLNet confidence score; the performance of the three machine learning algorithms did not vary much, and the SVM classifier achieved the best accuracy (78.9%) and F1 (79.2%) scores.

In an attempt to investigate what specific textual features besides those extracted by XLNet have the most impact on depression classification, we calculated the permutation feature importance [69] on the trained random forest classifier using the VADER, engagement, personality, and LIWC features with 10 repeats. The importance scores of individual features are shown in Figure

4. Among the LIWC features, "i," "bio," "power," "sad," and "authentic" are shown to be important in classification. Among the five personality features, "conscientiousness" and "neuroticism" were shown to be closely related to depression cues. We did not observe a strong relation between VADER sentiment features or social media engagement features and the depression signals. As for the LIWC sentiment features, only

"sad" and "anxiety" were shown to be relatively important. It is worth noting that LIWC's "sad" and "anxiety" categories each referred to about 150 words. By contrast, more than 7500 words or features fell in to the negative category in VADER. The insignificance of VADER features can be attributed to the more focused nature of LIWC.

**Table 2.** User-level performance (%) using different features.

| Features[a] | Accuracy | F1 | AUC[b] |
|---|---|---|---|
| VADER[c] | 54.9 | 61.7 | 54.6 |
| Demographics | 58.7 | 56.0 | 61.4 |
| Engagement | 58.7 | 62.3 | 61.7 |
| Personality | 64.8 | 67.8 | 72.4 |
| LIWC[d] | 70.6 | 70.8 | 76.0 |
| V + D + E + P + L[e] | 71.5 | 72.0 | 78.3 |
| XLNet | 78.1 | 77.9 | 84.9 |
| All (random forest) | 78.4 | 78.1 | 84.9 |
| All (logistic regression) | 78.3 | 78.5 | *86.4* [f] |
| All (SVM[g]) | *78.9* | *79.2* | 86.1 |

[a]We used SVM for classifying individual features.

[b]AUC: area under the receiver operating characteristic curve.

[c]VADER: Valence Aware Dictionary and Sentiment Reasoner.

[d]LIWC: Linguistic Inquiry and Word Count.

[e]V + D + E + P + L: VADER + demographics + engagement + personality + LIWC.

[f]Italics indicate the best performing model in each column.

[g]SVM: support vector machine.

**Figure 4.** Permutation importance of different features. LIWC: Linguistic Inquiry and Word Count; VADER: Valence Aware Dictionary for Sentiment Reasoning.



## Application Results

In this section, we report two COVID-19–related applications of our XLNet based depression classifier: (1) monitoring the

evolution of depression levels among the depression group and the nondepression group, and (2) monitoring the depression level at the US country level and state level during the pandemic.

We chose to use XLNet because of its simplicity as a stand-alone model, as it performed comparably to the fusion model.

### Depression Monitoring on Depression and Nondepression Groups

We took the 500 users from the testing set (n=500), along with their tweets from January 1 to May 22, 2020. We concatenated a user's tweets consecutively from January 1 one by one until reaching 250 words and labeled this chunk's date as the date of the author posting the tweet that was in the middle of the chunk. We grouped 3 days into a bin from January 1 and assigned the chunks to the bins according to the labeled date. We ran the XLNet model on the preprocessed tweet chunks and recorded the confidence scores. We trimmed the upper and lower 10% of the data to reduce the skew in the score distribution. We then took the mean of the scores for each time bin and plotted the depression trend shown in Figure 5. We further took a moving average of 5 time bins to smooth the curves.

**Figure 5.** Aggregated depression level trends of the depression and nondepression groups from January 1 to May 22, 2020. Since users with depression have a substantially higher depression level, we used different y-axes for the 2 groups' depression levels to compare them side by side.



Two immediate observations followed. First, depression level among users in the depression group was substantially higher than that in the nondepression group. This held across the entire observation period from early January to late May 2020. Second, and more importantly, the depression levels shared a strikingly similar trend among the two groups.

Delving deeper into these curves, we marked three important time points on the plot—the first confirmed case of COVID-19 in the United States (January 21, 2020), the US National Emergency announcement (March 13), and the last stay-at-home order issued (South Carolina, April 7). In January, both groups experienced a drop in depression scores. This may be caused by the fact that people's mood usually hits its lowest in winter [70]. From the day when there was the first confirmed case in the United States to the day of the announcement of the US National Emergency, the trends of the depression and nondepression groups were different. The depression level of the depression group went down slightly, while the depression level of the nondepression group went up. Aided by psychological findings, we hypothesized that depressive users were less affected by negative events happening in the outside world because they focused on their own feelings and life events, since they were mostly affected by negative events that threatened them directly [71] and more interact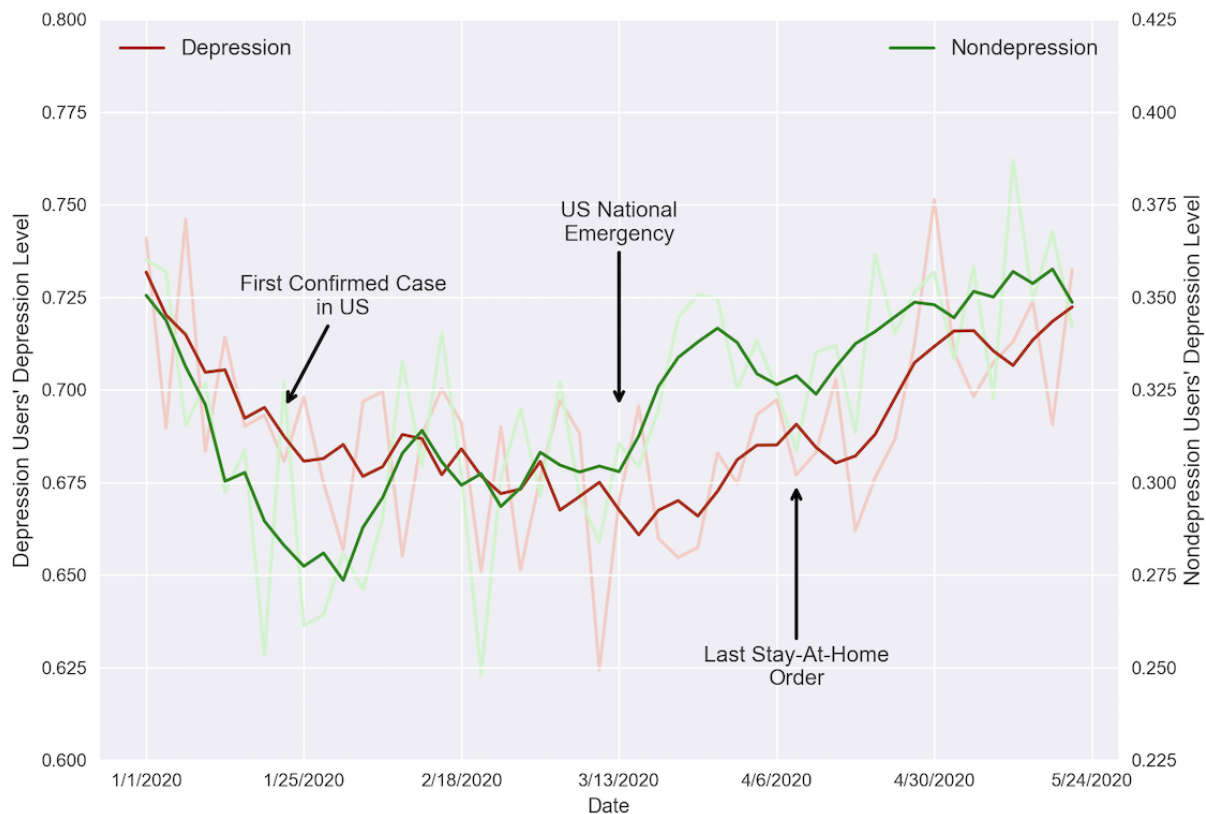ions with the outside world gave them more negative feedback [72]. Moreover, the depression levels of the depression and nondepression groups both increased after the announcement of the US National Emergency.

To better understand the trend, we applied the LDA model to retrieve the topics before and after the announcement of the US National Emergency. Each chunk of the tweets was assigned 5 weights for each of the 5 topics. We labeled the topic of the highest weight as the dominant topic of this chunk of the tweets and counted the frequency of each topic shown in Figure 6. Details about the keywords of the topics are reported in Multimedia Appendix 1. Before the announcement, the two most frequent topics of the depression and nondepression groups were the discussions about US President Donald Trump and

about school and work. The third most frequent topic of the nondepression group was about health while that of depression group was about entertainment. This supports the difference of the depression level trends of the two groups. After the announcement of the US National Emergency, the most frequent topic of the depression group was depression and anxiety during COVID-19, while this was the third most frequent topic of the nondepression group. Further, all 5 topics of each group were about COVID-19. This shows that, when people mostly talk about COVID-19, depression signals rise for both groups.

**Figure 6.** Topic distributions of depression and nondepression groups before and after the announcement of the US National Emergency.



### Aggregated Depression in COVID-19

To investigate country-level and state-level depression trends during COVID-19, we randomly sampled users who had US state locations stated in their profiles and crawled their tweets between March 3 and May 22, 2020, the period right before and after the US announced a National Emergency on March 13. Using the same logic as in the previous section, we plotted the change of depression scores of 9050 geolocated users (n=9050) sampled from the 36.6 million users mentioned, excluding those used for training, as the country-level trend. For state-level comparison, we plotted the aggregated scores of three representative states—economical center New York on the East Coast that was highly affected by the virus, tech center California on the West Coast that was also struck hard by the virus, and the less affected tourism center Florida in the southeast. Each selected state had at least 550 users in the data set to validate our findings. Their depression levels are shown in Figure 7.

The first observation of the plot is that depression scores of all three states and the United States behaved similarly during the pandemic; they experienced a decrease right before the National Emergency; a steady increase after that; a slight decrease past April 23, 2020; and another sharp increase after May 10. We also noticed that the overall depression score of Florida was substantially lower than the US average and the other two states. Since Florida had a lower score both before and after the virus outbreak, we hypothesized that it has a lower depression level overall compared to the average US level irrespective of the pandemic.

We calculated the topics at the state level after the announcement of the US National Emergency. As shown in Figure 8, the most frequent topic was the government's policy on COVID-19. California and Florida were the states that paid relatively more attention to this topic compared to the US average and New York State. Florida also talked more about the life change during COVID-19. Another finding was that people in New York talked more about the hospital news, likely because the state contained the majority of cases in the country by May 22, 2020 [73].

**Figure 7.** Aggregated depression level trends of the United States, New York, Califoria, and Florida after the announcement of the US National Emergency.
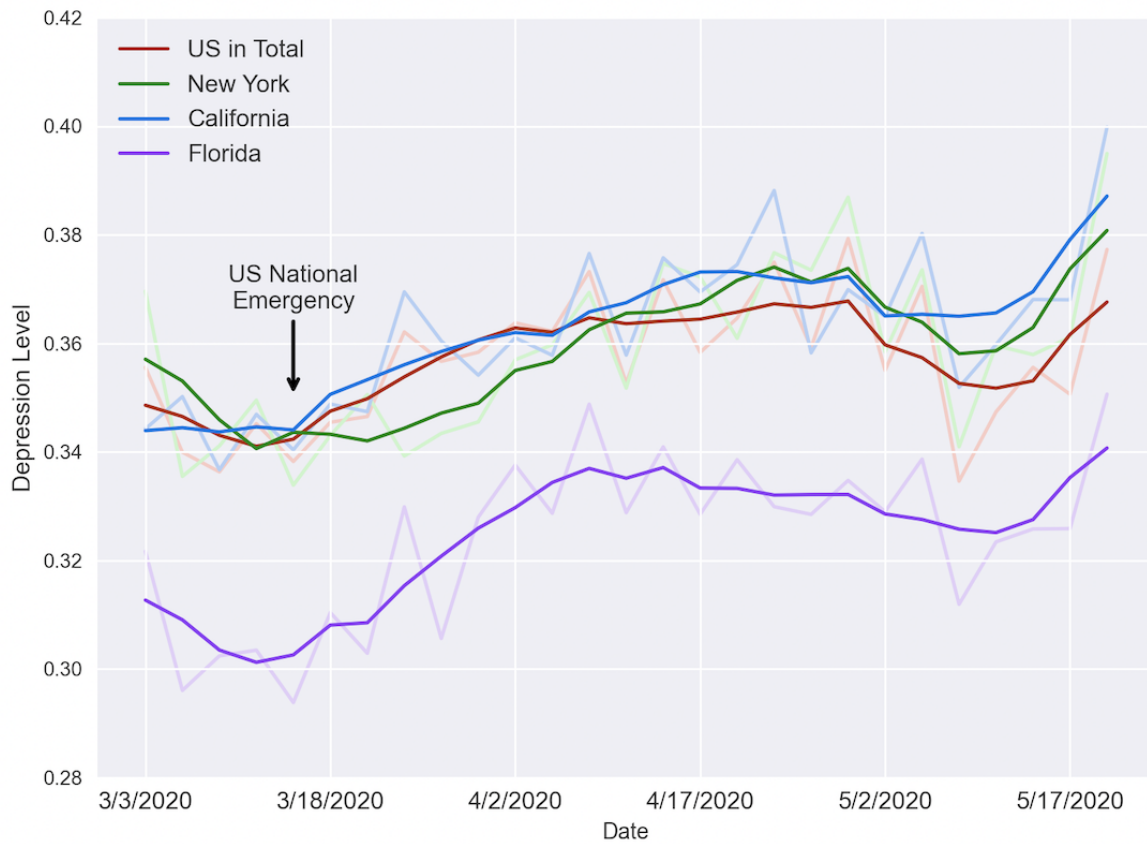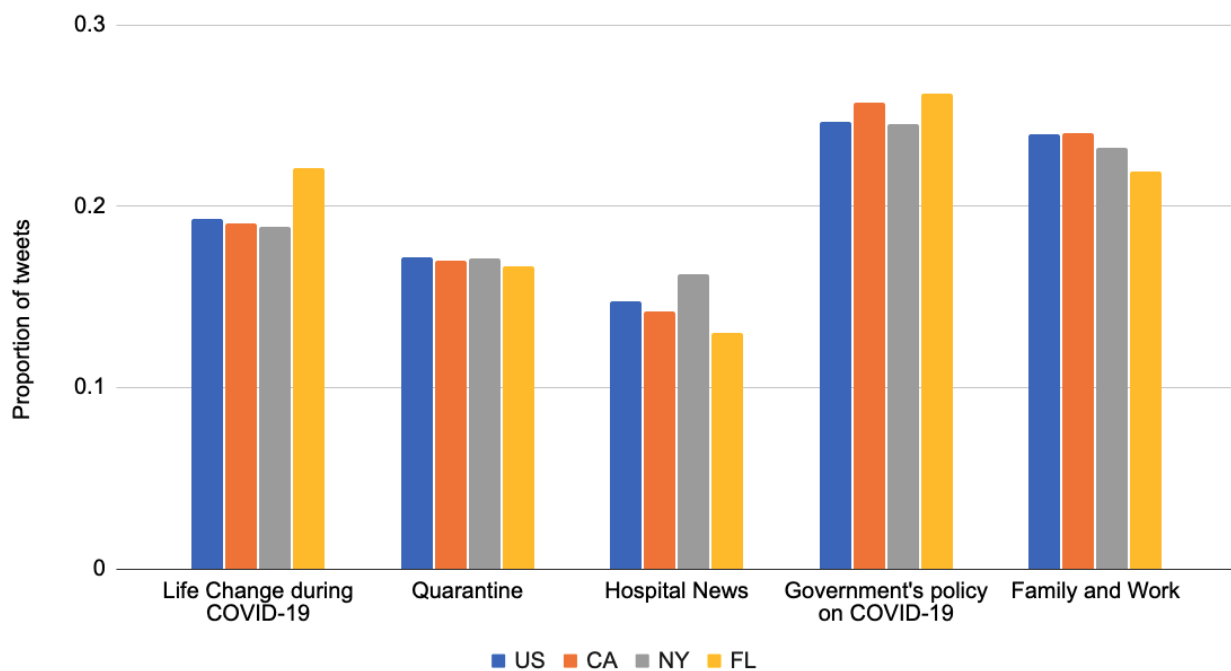


**Figure 8.** Distributions of the top 5 topics (state level) after the announcement of the US National Emergency.

## *Discussion*

### Principal Results

In this study, we developed a practical pipeline that included first gathering and cleaning a large-scale Twitter depression classification data set quickly in response to an outbreak, then training an accurate depression signal detection model on this data set, and finally applying the model to monitoring public depression trends. We analyzed the depression level trends during the COVID-19 pandemic, which shed light on the psychological impacts of the pandemic. Our main results were fourfold and corresponded to the four objectives listed in the *Goal of the Study* section.

First, using a stringent yet effective regular expression-based search method, we constructed by far the largest data set with 5150 Twitter users, including half identified as depression users and half as control users, along with their tweets within the past 3 months and their Twitter activity data.

Second, we developed a chunking and regrouping method to construct 32,420 tweet chunks, with 250 words each in the data set. We progressively added data to our training set and showed experimentally that the performance of deep learning models improves as the size of the training set grows, which validates the importance of our data set size. We compared the models' performances at the chunk level with the user level and observed further performance gain, which added credibility to our chunking method.

Third, we built a more accurate classification model (with 78.9% accuracy on n=449) upon the deep learning models along with linguistic analysis of dimensions including personality, LIWC, sentiment features, and demographic information. A permutation importance test showed that conscientiousness, neuroticism, appearance of first person pronouns, talking about biological processes such as eating and sleeping, talking about power, and exhibiting sadness are closely related to depression cues.

Finally, we showed the feasibility of the two proposed methods for monitoring the change of public depression levels as the disease propagates by aggregating individuals' past tweets within a time frame. Our method can target different groups of people, and we showed the depression trends of identified depression and nondepression groups (n=500), and of groups at different geolocations (n=9050). The temporal trends showed that the nondepression group's depression level rose earlier than that of the depression group, which we explained by psychological theories and LDA topics extracted from key time points. We also found that New York, California, Florida, and the United States in total all shared a similar depression trend, with Florida having a substantially lower depression level, which was also verified by LDA topic analysis.

### Practical Implications

Our study has practical implications. For example, upon detecting a rise in depression levels in a certain area, internet-based intervention services can be recommended by the social media platforms to the users. An intervention for depression commonly recommended is cognitive behavioral therapy (CBT), which is a type of therapy that targets one's irrational thinking patterns and unadaptable behavioral patterns [74]. During the COVID-19 period, digital-based CBT can be adopted. It has shown to be effective in reducing symptoms of mental disorders [75,76]. At the same time, it is also cost-effective and practical during the pandemic [75]. In addition to digital-based CBT, social media–based suicide prevention messages have also shown to be effective [77] and can be sent to individuals at risk.

### Limitations

Although our data collection method is fast and fully automatic, we acknowledge that the same limitations exist as noted in detail by Coppersmith et al [26]. Specifically, the users with depression captured by us can only represent a subpopulation (those who use Twitter and are willing to disclose their conditions) of the general depression population, and we cannot guarantee that the control group was not contaminated.

### Comparison With Prior Work

The data set used in this study containing 2575 depression users was much larger than those used previously, which contained 1402 depression users at most. De Choudhury et al [48] demonstrated that depression prediction models can potentially be used at the population level. However, to the best of our knowledge, all Twitter user depression identification studies reviewed in the introduction section focus on either tweet-level or user-level classification rather than applying the model to analyzing the mental health trends of a large population. To our knowledge, we were also the first to apply the transformer-based models (BERT, RoBERTa, XLNet) to identifying depression users on Twitter using a large-scale data set and to monitor the public depression trend.

### Conclusions

COVID-19 has infected over 100 million people worldwide [1], virtually bringing the whole world to a halt. During this period, social media witnessed a spike in depression terms. Against this backdrop, we have developed transformer-based models trained with by far the largest data set on depression. We have analyzed our models' performance in comparison to existing models and verified that the large training set we compiled was beneficial to improving the models' performance. We further showed that our models can be readily applied to the monitoring of stress and depression trends of targeted groups over geographical entities such as states. We noticed substantial increases in depression signals as people talked more about COVID-19. We hope researchers and mental health practitioners find our models useful and that this study raises awareness of the mental health impacts of the pandemic.

## Authors' Contributions

YZ and JL conceived and designed the study. YZ performed regular expression search and preprocessing, examined feature importance, and wrote the majority of the manuscript. HL performed data collection and applied the LDA models. HL and YZ analyzed the data and wrote part of the manuscript. YZ and YL trained the models and performed depression monitoring. XZ analyzed the findings using psychological theories. All authors helped design the study and edit the manuscript.

## Conflicts of Interest

None declared.

Multimedia Appendix 1
Supplemental data statistics and tables.
[DOCX File , 24 KB - infodemiology_v1i1e26769_app1.docx ]

## References

1. Coronavirus disease (COVID-19) pandemic. World Health Organization. URL: https://www.who.int/emergencies/diseases/novel-coronavirus-2019 [accessed 2020-12-21]
2. Education: from disruption to recovery. UNESCO. 2020 Mar 04. URL: https://en.unesco.org/covid19/educationresponse [accessed 2020-12-21]
3. Mental disorders. World Health Organization. URL: https://www.who.int/news-room/fact-sheets/detail/mental-disorders [accessed 2020-12-21]
4. Inskip H, Harris C, Barraclough B. Lifetime risk of suicide for affective disorder, alcoholism and schizophrenia. Br J Psychiatry 1998 Jan;172:35-37. [doi: 10.1192/bjp.172.1.35] [Medline: 9534829]
5. Too LS, Spittal MJ, Bugeja L, Reifels L, Butterworth P, Pirkis J. The association between mental disorders and suicide: a systematic review and meta-analysis of record linkage studies. J Affect Disord 2019 Dec 01;259:302-313 [FREE Full text] [doi: 10.1016/j.jad.2019.08.054] [Medline: 31450139]
6. Yoshikawa E, Taniguchi T, Nakamura-Taira N, Ishiguro S, Matsumura H. Factors associated with unwillingness to seek professional help for depression: a web-based survey. BMC Res Notes 2017 Dec 04;10(1):673 [FREE Full text] [doi: 10.1186/s13104-017-3010-1] [Medline: 29202791]
7. Baker SR, Bloom N, Davis SJ, Terry SJ. COVID-induced economic uncertainty. Natl Bureau Econ Res 2020:w26983. [doi: 10.3386/w26983]
8. Nicola M, Alsafi Z, Sohrabi C, Kerwan A, Al-Jabir A, Iosifidis C, et al. The socio-economic implications of the coronavirus pandemic (COVID-19): a review. Int J Surg 2020 Jul;78:185-193 [FREE Full text] [doi: 10.1016/j.ijsu.2020.04.018] [Medline: 32305533]
9. Xiong J, Lipsitz O, Nasri F, Lui LM, Gill H, Phan L, et al. Impact of COVID-19 pandemic on mental health in the general population: a systematic review. J Affect Disord 2020 Dec 01;277:55-64 [FREE Full text] [doi: 10.1016/j.jad.2020.08.001] [Medline: 32799105]
10. Le XTT, Dang AK, Toweh J, Nguyen QN, Le HT, Do TTT, et al. Evaluating the psychological impacts related to COVID-19 of Vietnamese people under the first nationwide partial lockdown in Vietnam. Front Psychiatry 2020;11:824. [doi: 10.3389/fpsyt.2020.00824] [Medline: 32982807]
11. Wang C, Chudzicka-Czupała A, Grabowski D, Pan R, Adamus K, Wan X, et al. The association between physical and mental health and face mask use during the COVID-19 pandemic: a comparison of two countries with different views and practices. Front Psychiatry 2020;11:569981. [doi: 10.3389/fpsyt.2020.569981] [Medline: 33033485]
12. Hao F, Tan W, Jiang L, Zhang L, Zhao X, Zou Y, et al. Do psychiatric patients experience more psychiatric symptoms during COVID-19 pandemic and lockdown? A case-control study with service and research implications for immunopsychiatry. Brain Behav Immun 2020 Jul;87:100-106 [FREE Full text] [doi: 10.1016/j.bbi.2020.04.069] [Medline: 32353518]
13. Mitchell M, Hollingshead K, Coppersmith G. Quantifying the language of schizophrenia in social media. 2015 Presented at: 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality; June 5, 2015; Denver, CO p. 11-20. [doi: 10.3115/v1/w15-1202]
14. Preoţiuc-Pietro D, Eichstaedt J, Park G, Sap M, Smith L, Tobolsky V, et al. The role of personality, age, and gender in tweeting about mental illness. 2015 Presented at: 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality; June 5, 2015; Denver, CO p. 21-30. [doi: 10.3115/v1/w15-1203]
15. Conway M, O'Connor D. Social media, big data, and mental health: current advances and ethical implications. Curr Opin Psychol 2016 Jun;9:77-82 [FREE Full text] [doi: 10.1016/j.copsyc.2016.01.004] [Medline: 27042689]
16. Ernala SK, Labetoulle T, Bane F, Birnbaum ML, Rizvi AF, Kane JM, et al. Characterizing audience engagement and assessing its impact on social media disclosures of mental illnesses. 2018 Presented at: Twelfth International AAAI Conference on Web and Social Media; June 25-28, 2018; Palo Alto, CA URL: https://ojs.aaai.org/index.php/ICWSM/article/view/15027

XSL•FO
RenderX

17. Mazuz K, Yom-Tov E. Analyzing trends of loneliness through large-scale analysis of social media postings: observational study. JMIR Ment Health 2020 May 20;7(4):e17188 [FREE Full text] [doi: 10.2196/17188] [Medline: 32310141]

18. Ginsberg J, Mohebbi MH, Patel RS, Brammer L, Smolinski MS, Brilliant L. Detecting influenza epidemics using search engine query data. Nature 2009 Mar 19;457(7232):1012-1014. [doi: 10.1038/nature07634] [Medline: 19020500]

19. Yeung N, Lai J, Luo J. Face off: polarized public opinions on personal face mask usage during the COVID-19 pandemic. 2020 Presented at: 2020 IEEE International Conference on Big Data (Big Data); December 10-13, 2020; Atlanta, GA. [doi: 10.1109/bigdata50022.2020.9378114]

20. Wu W, Lyu H, Luo J. Characterizing discourse about COVID-19 vaccines: a Reddit version of the pandemic story. arXiv. Preprint posted online on January 15, 2021 [FREE Full text]

21. Lyu H, Wang J, Wu W, Duong V, Zhang X, Dye TD, et al. Social media study of public opinions on potential COVID-19 vaccines: informing dissent, disparities, and dissemination. arXiv. Preprint posted online on December 3, 2020 [FREE Full text]

22. Huang W, Cao B, Yang G, Luo N, Chao N. Turn to the internet first? Using online medical behavioral data to forecast COVID-19 epidemic trend. Inf Process Manag 2021 May;58(3):102486 [FREE Full text] [doi: 10.1016/j.ipm.2020.102486] [Medline: 33519039]

23. Wang C, Pan R, Wan X, Tan Y, Xu L, McIntyre RS, et al. A longitudinal study on the mental health of general population during the COVID-19 epidemic in China. Brain Behav Immun 2020 Jul;87:40-48 [FREE Full text] [doi: 10.1016/j.bbi.2020.04.028] [Medline: 32298802]

24. Brandwatch. URL: https://www.brandwatch.com/ [accessed 2020-12-21]

25. De Choudhury M, Gamon M, Counts S, Horvitz E. Predicting depression via social media. 2013 Presented at: Seventh International AAAI Conference on Weblogs and Social Media; July 8-11, 2013; Boston, MA URL: https://www.aaai.org/ocs/index.php/ICWSM/ICWSM13/paper/view/6124

26. Coppersmith G, Dredze M, Harman C. Quantifying mental health signals in Twitter. 2014 Presented at: Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality; June 2014; Baltimore, MA p. 51-60. [doi: 10.3115/v1/w14-3207]

27. Coppersmith G, Dredze M, Harman C, Hollingshead K, Mitchell M. CLPsych 2015 Shared Task: depression and PTSD on Twitter. 2015 Presented at: 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality; June 5, 2015; Denver, CO p. 31-39. [doi: 10.3115/v1/w15-1204]

28. Jamil Z, Inkpen D, Buddhitha P, White K. Monitoring tweets for depression to detect at-risk users. 2017 Presented at: Fourth Workshop on Computational Linguistics and Clinical Psychology — From Linguistic Signal to Clinical Reality; August 2017; Vancouver, BC p. 32-40. [doi: 10.18653/v1/w17-3104]

29. Orabi AH, Buddhitha P, Orabi MH, Inkpen D. Deep learning for depression detection of Twitter users. 2018 Presented at: Fifth Workshop on Computational Linguistics and Clinical Psychology: From Keyboard to Clinic; June 2018; New Orleans, LA p. 88-97. [doi: 10.18653/v1/w18-0609]

30. Nadeem M. Identifying depression on Twitter. arXiv. 2016 Jul 25. URL: http://arxiv.org/abs/1607.07384 [accessed 2020-12-21]

31. Zirikly A, Resnik P, Uzuner Ö, Hollingshead K. CLPsych 2019 Shared Task: predicting the degree of suicide risk in Reddit posts. 2019 Presented at: Sixth Workshop on Computational Linguistics and Clinical Psychology; June 2019; Minneapolis, MN p. 24-33. [doi: 10.18653/v1/w19-3003]

32. Matero M, Idnani A, Son Y, Giorgi S, Vu H, Zamani M, et al. Suicide risk assessment with multi-level dual-context language and BERT. 2019 Presented at: Sixth Workshop on Computational Linguistics and Clinical Psychology; June 2019; Minneapolis, MN p. 39-44. [doi: 10.18653/v1/w19-3005]

33. Devlin J, Chang MW, Lee K, Toutanova K. BERT: pre-training of deep bidirectional transformers for language understanding. 2019 Presented at: 2019 Conference of the North American Chapter of the Association for Computational Linguistics; June 2019; Minneapolis, MN p. 4171-4186. [doi: 10.3115/1614108]

34. Ramírez-Cifuentes D, Freire A, Baeza-Yates R, Puntí J, Medina-Bravo P, Velazquez DA, et al. Detection of suicidal ideation on social media: multimodal, relational, and behavioral analysis. J Med Internet Res 2020 Jul 07;22(7):e17758 [FREE Full text] [doi: 10.2196/17758] [Medline: 32673256]

35. Brådvik L. Suicide risk and mental disorders. Int J Environ Res Public Health 2018 Sep 17;15(9):2028 [FREE Full text] [doi: 10.3390/ijerph15092028] [Medline: 30227658]

36. Handley T, Rich J, Davies K, Lewin T, Kelly B. The challenges of predicting suicidal thoughts and behaviours in a sample of rural Australians with depression. Int J Environ Res Public Health 2018 May 07;15(5):928 [FREE Full text] [doi: 10.3390/ijerph15050928] [Medline: 29735902]

37. Ramirez-Esparza N, Chung C, Kacewicz E, Pennebaker JW. The psychology of word use in depression forums in English and in Spanish. In: Proceedings of the 2008 International Conference on Weblogs and Social Media. 2008 Presented at: ICWSM '08; March 30-April 2, 2008; Seattle, WA.

38. Pirina I, Çöltekin Ç. Identifying depression on Reddit: the effect of training data. In: Proceedings of the 2018 EMNLP Workshop SMM4H. 2018 Presented at: 3rd Social Media Mining for Health Applications Workshop & Shared Task; October 2018; Brussels, Belgium p. 9-12. [doi: 10.18653/v1/w18-5903]

39. Tadesse MM, Lin H, Xu B, Yang L. Detection of depression-related posts in Reddit social media forum. IEEE Access 2019;7:44883-44893. [doi: 10.1109/access.2019.2909180]

40. Liu Y, Ott M, Goyal N, Du J, Joshi M, Chen D, et al. RoBERTa: a robustly optimized BERT pretraining approach. OpenReview. 2019 Sep 25. URL: https://openreview.net/forum?id=SyxS0T4tvS [accessed 2020-12-21]

41. Yang Z, Dai Z, Yang Y, Carbonell J, Salakhutdinov R, Le QV. XLNet: generalized autoregressive pretraining for language understanding. Adv Neural Inf Process Syst 2019;32:5753-5763.

42. Tsugawa S, Kikuchi Y, Kishino F, Nakajima K, Itoh Y, Ohsaki H. Recognizing depression from Twitter activity. In: Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems. 2015 Presented at: CHI '15; Seoul, Republic of Korea; April 2015 p. 3187-3196. [doi: 10.1145/2702123.2702280]

43. Blei DM, Ng AY, Jordan MI. Latent Dirichlet allocation. J Machine Learning Res 2003;3:993-1022.

44. Zhou D, Luo J, Silenzio V, Zhou Y, Hu J, Currier G, et al. Tackling mental health by integrating unobtrusive multimodal sensing. In: Proceedings of the AAAI Conference on Artificial Intelligence. 2015 Presented at: AAAI-15; January 25-30, 2015; Austin, TX.

45. Leis A, Ronzano F, Mayer MA, Furlong LI, Sanz F. Detecting signs of depression in tweets in Spanish: behavioral and linguistic analysis. J Med Internet Res 2019 Jun 27;21(6):e14199 [FREE Full text] [doi: 10.2196/14199] [Medline: 31250832]

46. Zhang B, Zaman A, Silenzio V, Kautz H, Hoque E. The relationships of deteriorating depression and anxiety with longitudinal behavioral changes in Google and YouTube use during COVID-19: observational study. JMIR Ment Health 2020 Dec 23;7(11):e24012 [FREE Full text] [doi: 10.2196/24012] [Medline: 33180743]

47. Shen G, Jia J, Nie L, Feng F, Zhang C, Hu T, et al. Depression detection via harvesting social media: a multimodal dictionary learning solution. In: Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence. 2017 Presented at: IJCAI-17; August 2017; Melbourne, Australia. [doi: 10.24963/ijcai.2017/536]

48. De Choudhury M, Counts S, Horvitz E. Social media as a measurement tool of depression in populations. In: Proceedings of the 5th Annual ACM Web Science Conference. 2013 Presented at: WebSci '13; May 2013; Paris, France p. 47-56. [doi: 10.1145/2464464.2464480]

49. Winthorst WH, Roest AM, Bos EH, Meesters Y, Penninx BW, Nolen WA, et al. Seasonal affective disorder and non-seasonal affective disorders: results from the NESDA study. BJPsych Open 2017 Jul;3(4):196-203 [FREE Full text] [doi: 10.1192/bjpo.bp.116.004960] [Medline: 28904813]

50. Hakulinen C, Elovainio M, Pulkki-Råback L, Virtanen M, Kivimäki M, Jokela M. Personality and depressive symptoms: individual participant meta-analysis of 10 cohort studies. Depress Anxiety 2015 Jul;32(7):461-470 [FREE Full text] [doi: 10.1002/da.22376] [Medline: 26014798]

51. Watson Personality Insights. IBM. URL: https://www.ibm.com/cloud/watson-personality-insights [accessed 2020-12-21]

52. Zunic A, Corcoran P, Spasic I. Sentiment analysis in health and well-being: systematic review. JMIR Med Inform 2020 Jan 28;8(1):e16023 [FREE Full text] [doi: 10.2196/16023] [Medline: 32012057]

53. Hutto C, Gilbert E. VADER: a parsimonious rule-based model for sentiment analysis of social media text. In: Proceedings of the International AAAI Conference on Web and Social Media. 2014 Presented at: ICWSM-14; June 1-4, 2014; Ann Arbor, MI URL: https://www.aaai.org/ocs/index.php/ICWSM/ICWSM14/paper/view/8109

54. Mirowsky J, Ross CE. Age and depression. J Health Soc Behav 1992 Sep;33(3):187. [doi: 10.2307/2137349]

55. Wainwright N, Surtees P. Childhood adversity, gender and depression over the life-course. J Affect Disord 2002 Oct;72(1):33-44. [doi: 10.1016/s0165-0327(01)00420-7] [Medline: 12204315]

56. Wang C, Pan R, Wan X, Tan Y, Xu L, Ho CS, et al. Immediate psychological responses and associated factors during the initial stage of the 2019 coronavirus disease (COVID-19) epidemic among the general population in China. Int J Environ Res Public Health 2020 Mar 06;17(5):1729 [FREE Full text] [doi: 10.3390/ijerph17051729] [Medline: 32155789]

57. Albert P. Why is depression more prevalent in women? J Psychiatry Neurosci 2015 Jul;40(4):219-221 [FREE Full text] [doi: 10.1503/jpn.150205] [Medline: 26107348]

58. Wang Z, Hale S, Adelani DI, Grabowicz P, Hartman T, Flöck F, et al. Demographic inference and representative population estimates from multilingual social media data. 2019 Presented at: WWW '19: The World Wide Web Conference; May 2019; San Francisco, CA p. 2056-2067. [doi: 10.1145/3308558.3313684]

59. Tausczik YR, Pennebaker JW. The psychological meaning of words: LIWC and computerized text analysis methods. J Lang Soc Psychol 2009 Dec 08;29(1):24-54. [doi: 10.1177/0261927X09351676]

60. Zhang X, Wang Y, Lyu H, Zhang Y, Liu Y, Luo J. The influence of COVID-19 on well-being. PsyArXiv. Preprint posted online on May 7, 2020. [doi: 10.31234/osf.io/znj7h]

61. Rude S, Gortner E, Pennebaker J. Language use of depressed and depression-vulnerable college students. Cogn Emotion 2004 Dec;18(8):1121-1133. [doi: 10.1080/02699930441000030]

62. Stirman SW, Pennebaker JW. Word use in the poetry of suicidal and nonsuicidal poets. Psychosom Med 2001;63(4):517-522. [doi: 10.1097/00006842-200107000-00001] [Medline: 11485104]

63. Chen L, Lyu H, Yang T, Wang Y, Luo J. In the eyes of the beholder: analyzing social media use of neutral and controversial terms for COVID-19. arXiv. Preprint posted online on April 21, 2020 [FREE Full text]

64.   Holmes D, Alpers GW, Ismailji T, Classen C, Wales T, Cheasty V, et al. Cognitive and emotional processing in narratives of women abused by intimate partners. Violence Against Women 2007 Dec;13(11):1192-1205. [doi: 10.1177/1077801207307801] [Medline: 17951592]

65.   McClelland DC. Inhibited power motivation and high blood pressure in men. J Abnorm Psychol 1979 May;88(2):182-190. [doi: 10.1037//0021-843x.88.2.182] [Medline: 447901]

66.   Baziotis C, Pelekis N, Doulkeridis C. DataStories at SemEval-2017 Task 4: deep LSTM with attention for message-level and topic-based sentiment analysis. In: Proceedings of the 11th International Workshop on Semantic Evaluations. 2017 Presented at: SemEval-2017; August 3-4, 2017; Vancouver, BC p. 747-754. [doi: 10.18653/v1/s17-2126]

67.   Pennington J, Socher R, Manning C. GloVe: Global Vectors for Word Representation. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing. 2014 Presented at: EMNLP '14; October 2014; Doha, Qatar p. 1532-1543. [doi: 10.3115/v1/d14-1162]

68.   scikit-learn. URL: https://scikit-learn.org/stable/ [accessed 2020-12-21]

69.   Altmann A, Toloşi L, Sander O, Lengauer T. Permutation importance: a corrected feature importance measure. Bioinformatics 2010 May 15;26(10):1340-1347. [doi: 10.1093/bioinformatics/btq134] [Medline: 20385727]

70.   Thompson C, Stinson D, Fernandez M, Fine J, Isaacs G. A comparison of normal, bipolar and seasonal affective disorder subjects using the Seasonal Pattern Assessment Questionnaire. J Affect Disord 1988;14(3):257-264. [doi: 10.1016/0165-0327(88)90043-2] [Medline: 2968387]

71.   Li Y, Zhang D, Liang Y, Hu T. Meta-analysis of the relationship between life events and depression in adolescents. J Pediatr Care 2016;2(1):1. [doi: 10.21767/2471-805x.100008]

72.   Winer ES, Salem T. Reward devaluation: dot-probe meta-analytic evidence of avoidance of positive information in depressed persons. Psychol Bull 2016 Jan;142(1):18-78 [FREE Full text] [doi: 10.1037/bul0000022] [Medline: 26619211]

73.   COVID data tracker. Centers for Disease Control and Prevention. 2020 Mar 28. URL: https://covid.cdc.gov/covid-data-tracker [accessed 2020-12-21]

74.   Ho CS, Chee CY, Ho RC. Mental health strategies to combat the psychological impact of COVID-19 beyond paranoia and panic. Ann Acad Med Singap 2020 Mar 16;49(3):155-160 [FREE Full text] [Medline: 32200399]

75.   Zhang MW, Ho RC. Moodle: the cost effective solution for internet cognitive behavioral therapy (I-CBT) interventions. Technol Health Care 2017;25(1):163-165. [doi: 10.3233/THC-161261] [Medline: 27689560]

76.   Soh HL, Ho RC, Ho CS, Tam WW. Efficacy of digital cognitive behavioural therapy for insomnia: a meta-analysis of randomised controlled trials. Sleep Med 2020 Nov;75:315-325. [doi: 10.1016/j.sleep.2020.08.020] [Medline: 32950013]

77.   Robinson J, Bailey E, Hetrick S, Paix S, O'Donnell M, Cox G, et al. Developing social media-based suicide prevention messages in partnership with young people: exploratory study. JMIR Ment Health 2017 Oct 04;4(4):e40 [FREE Full text] [doi: 10.2196/mental.7847] [Medline: 28978499]

## Abbreviations

**API:** application programming interface
**AUC:** area under the receiver operating characteristic curve
**BERT:** Bidirectional Encoder Representations from Transformers
**BiLSTM:** bidirectional long short-term memory
**BOW:** bag of words
**CBT:** cognitive behavioral therapy
**CNN:** convolutional neural network
**LDA:** latent Drichlet allocation
**LIWC:** Linguistic Inquiry and Word Count
**PTSD:** posttraumatic stress disorder
**RoBERTa:** Robustly Optimized Bidirectional Long Short-Term Memory Pretraining Approach
**SVM:** support vector machine
**VADER:** Valence Aware Dictionary and Sentiment Reasoner

XSL•FO
RenderX

XSL•FO

**RenderX**

<u>Original Paper</u>

# Public Attitudes and Factors of COVID-19 Testing Hesitancy in the United Kingdom and China: Comparative Infodemiology Study

Leesa Lin[1,2*], PhD; Yi Song[3*], BSc; Qian Wang[3*], MPH; Jialu Pu[3], BSc; Fiona Yueqian Sun[1], MSc; Yixuan Zhang[3], BSc; Xinyu Zhou[3], BSc; Heidi J Larson[1], PhD; Zhiyuan Hou[3,4,5], PhD

[1]Department of Infectious Disease Epidemiology, London School of Hygiene & Tropical Medicine, London, United Kingdom

[2]Laboratory of Data Discovery for Health, Hong Kong Science Park, Hong Kong SAR, China

[3]School of Public Health, Fudan University, Shanghai, China

[4]NHC Key Laboratory of Health Technology Assessment, Fudan University, Shanghai, China

[5]Global Health Institute, Fudan University, Shanghai, China

[*]these authors contributed equally

**Corresponding Author:**
Zhiyuan Hou, PhD
School of Public Health
Fudan University
130 Dong'an Road
Shanghai, 200032
China
Phone: 86 2133563935
Email: zyhou@fudan.edu.cn

## *Abstract*

**Background:** Massive community-wide testing has become the cornerstone of management strategies for the COVID-19 pandemic.

**Objective:** This study was a comparative analysis between the United Kingdom and China, which aimed to assess public attitudes and uptake regarding COVID-19 testing, with a focus on factors of COVID-19 testing hesitancy, including effectiveness, access, risk perception, and communication.

**Methods:** We collected and manually coded 3856 UK tweets and 9299 Chinese Sina Weibo posts mentioning COVID-19 testing from June 1 to July 15, 2020. Adapted from the World Health Organization's 3C Model of Vaccine Hesitancy, we employed social listening analysis examining key factors of COVID-19 testing hesitancy (confidence, complacency, convenience, and communication). Descriptive analysis, time trends, geographical mapping, and chi-squared tests were performed to assess the temporal, spatial, and sociodemographic characteristics that determine the difference in attitudes or uptake of COVID-19 tests.

**Results:** The UK tweets demonstrated a higher percentage of support toward COVID-19 testing than the posts from China. There were much wider reports of public uptake of COVID-19 tests in mainland China than in the United Kingdom; however, uncomfortable experiences and logistical barriers to testing were more expressed in China. The driving forces for undergoing COVID-19 testing were personal health needs, community-wide testing, and mandatory testing policies for travel, with major differences in the ranking order between the two countries. Rumors and information inquiries about COVID-19 testing were also identified.

**Conclusions:** Public attitudes and acceptance toward COVID-19 testing constantly evolve with local epidemic situations. Policies and information campaigns that emphasize the importance of timely testing and rapid communication responses to inquiries and rumors, and provide a supportive environment for accessing tests are key to tackling COVID-19 testing hesitancy and increasing uptake.

XSL•FO
RenderX

## Introduction

As the number of COVID-19 cases accelerated globally in early 2020, many public health experts advocated for widespread rapid testing that could complement other containment strategies, such as hand washing, contact tracing, and quarantine, and that should be viewed as important as face covering, social distancing, and vaccines [1,2]. There are two widely accepted types of tests as follows: (1) a nucleic acid test, which is a polymerase chain reaction test that detects RNA (or genetic material) specific to the virus, and (2) an antigen test, which is a rapid turnaround virus test from a lateral flow device that can process COVID-19 samples on site without the need for laboratory equipment. Community-wide COVID-19 testing helped public health investigators understand the prevalence, contagiousness, and mortality of the disease [3], and has made it possible for communities to exit lockdowns and rapidly control potential resurgences while awaiting a safe vaccine. China, Singapore, Germany, and South Korea have been among the most early and aggressive countries in utilizing widespread frequent rapid tests (offered freely to residents) as a central pillar of their multipronged epidemic control strategies. In China, mass testing has been employed as a standard procedure in places where new outbreaks of COVID-19 surged [4]. During the resurgence in Beijing in June 2020, 3.56 million individuals at risk were tracked and tested [5], and the outbreak was quickly brought under control. In contrast, countries, such as the United Kingdom and Japan, had delays in rolling out mass testing [6], and asymptomatic individuals and high-risk populations (ie, health care workers) were not able to access testing in the pandemic's early stage due to limited capacity [7]. Two and a half months after nationwide lockdown, on May 28, 2020, the United Kingdom eventually launched the National Health System (NHS) Test and Trace program, a "world-beating system" that the Prime Minister had pledged to deliver as a central part of the government's COVID-19 recovery strategy. By July, people with COVID-19 symptoms could receive a test from the NHS without charge, and those engaged in high-risk jobs were promised regular testing by the UK government [8,9]. However, since it was introduced, the program has been repeatedly criticized for not meeting expectations. As the pandemic response progresses, the challenge of conducting COVID-19 mass testing will transition from inadequate testing capacity to inadequate uptake [10] due to pandemic fatigue, test anxiety, stigmatization, rumors, misinformation, fear of isolation and quarantine, and other disincentives.

Infodemiology, first introduced by Dr Gunther Eysenbach in early 2000 as the epidemiology of (mis)information [11], is an emerging field of research on the distribution and determinants of user-contributed health information and misinformation across the internet or in a population, with the ultimate aim of improving public health and public policy [12]. The COVID-19 pandemic created a paradigm shift in communication and infodemiology, as widespread negative health and socioeconomical impacts were observed to be caused by two concurrent pandemics (the novel coronavirus and misinformation). Social listening in the context of public health has been found to be an effective tool that offers real-time big data on public sentiment and opinions for informing and assessing governments' risk communication strategies and public reactions, especially during acute epidemic outbreaks, such as the 2009/2010 H1N1 [13], 2013/2014 Middle East respiratory syndrome (MERS) [14], and 2014 Ebola [13] epidemics. Unlike traditional research methods (eg, surveys or in-depth interviews), where opinion gathering is limited to interactions between researchers and participants, social listening allows for a rapid and thorough scanning of a multilevel dynamic information environment for digital opinions derived from public contributions, interactions, and interinfluences without researcher involvement. Social listening investigates public understanding and experiences of an event (ie, risks and countermeasures), which, as depicted in Stuart Hall's audience reception theory [15], are shaped by their individual sociocultural backgrounds and life experiences.

At-risk individuals refusing or avoiding testing could undermine a community's epidemic control and reopening strategies. Public health experts and decision makers must monitor public sentiment and acceptance toward testing and understand the root causes of testing hesitancy. To date, research has mainly focused on COVID-19 vaccines [16] and other nonpharmaceutical measures, such as lockdowns, social distancing, and mask wearing [17-20], leaving COVID-19 testing hesitancy and avoidance underinvestigated. The United Kingdom and China have highly active microblog users and have experienced initial COVID-19 outbreaks, lockdowns, and resurgence, yet mass testing was introduced in these two countries at different stages of response. As such, this infodemiology study aimed to assess public attitudes and uptakes of COVID-19 testing in the United Kingdom and China, with a focus on the factors of testing hesitancy, including effectiveness, access, risk perception, and communication.

## Methods

### Data Collection

We collected microblog posts from popular social media platforms in the United Kingdom and China. We assessed Twitter tweets (the United Kingdom) and Sina Weibo posts (mainland China) mentioning COVID-19 testing from June 1 to July 15, 2020, after the launch of the NHS Test and Trace system in the United Kingdom and the mass testing campaign in Beijing, China, during COVID-19's resurgence and before another resurgence in Xinjiang Autonomous Region, China. We used the Meltwater platform [21] to collect Twitter tweets and Weibo posts. The keywords used for collecting tweets or Weibo posts were "covid test," "covid19 test," "covid-19 test," "coronavirus test," "test for covid," and "test for coronavirus" ("核酸检测" in Chinese). Overall, 59,919 tweets from the United Kingdom (including 11,249 tweets from London) and 313,092 Weibo posts from China (including 82,743 Weibo posts from Beijing) were collected with the location and time they were sent. Weibo posts were downloaded daily so as to minimize possible bias resulting from posts being removed by authorities. We also downloaded the account profile of each Weibo post, from which we extracted gender, age, and education for analysis. Only human-contributed opinions/conversations on Twitter and

Weibo were included for analysis. Tweets or Weibo posts from news and organizational accounts, and tweets/posts generated by bots were identified by keyword matching, and then examined and removed by researchers. Duplicate tweets/posts, tweets/posts with identical text but from different accounts, retweets, and quotes without comments were removed. After removing posts that did not meet the inclusion criteria, we randomly sampled 10% of the tweets and posts by day for coding. In total, 3856 tweets from the United Kingdom, including 794 tweets from London, and 9299 Weibo posts from mainland China, including 3155 posts from Beijing, were included for formal analysis. Multimedia Appendix 1 shows the workflow of the inclusion and exclusion processes of the tweets and Weibo posts.

Analyses of accounts of social media users (Multimedia Appendix 2) suggested that our data of social media posts were well-representative of the entire social media user base. We found that 92.8% (3581/3856) of tweets in the United Kingdom and 97.5% (9067/9299) of Weibo posts in China were single posts sent by unique users.

To assess data representativeness, we compared Weibo users' demographic data in our study with the "Weibo 2020 User Development Report" [22]. This report showed that active female users accounted for 54.6% of the user base and that the user base is skewed toward young users; 78% of Weibo users are under 30 years old. Not all Weibo users' profiles were available, and in our data set with users' profiles, 67.6% (5940/8784) of users were female and 70.6% (2705/3830) were under 30 years old. Our results showed that our demographic profiles were comparable to the overall user base profile reported by Weibo.

## Data Analysis

This study employed content analysis of social media data in relation to COVID-19 testing [23-25], complemented by contextual epidemic data. COVID-19 epidemic data from the United Kingdom and mainland China were derived for trend analysis [26,27]. We plotted the trends of daily new COVID-19 case numbers in the United Kingdom and China to describe the epidemic context where mass testing programs were introduced.

We identified and classified social media posts that expressed personal opinions/discussions on COVID-19 testing. Public attitudes toward COVID-19 testing were manually screened and coded based on the three different positions one might take as follows: dominant (understanding and accepting the objectives of the test), negotiated (reacting with a mixture of acceptance and rejection), and oppositional (opposite to the dominant position and completely rejecting the test) [15]. To investigate the determinants of public attitudes toward COVID-19 testing, we employed a deductive approach with a coding framework that was adapted from the World Health

Organization's "3Cs" model of hesitancy toward vaccination [28], and this "3Cs" model was also applied for other public health behaviors, such as COVID-19 testing. The coding framework, presented in Multimedia Appendix 3, covers major factors of COVID-19 testing hesitancy, including *confidence* (degree of trust in the effectiveness and safety of the test), *complacency* (perceptions of personal risk associated with the disease and test), and *convenience* (influencers of the decision to get the test, eg, availability, affordability, and geographical accessibility), as well as *communication* (information inquiries and rumors about COVID-19 testing) [29].

In execution, we first developed a codebook with code definitions. Then, two researchers (YS and QW) coded a subsample of 500 posts independently, and when appropriate, refined the codebook. When necessary, SentiWordNet [30] was referenced. Using the final codebook, another subsample of 200 posts was independently coded to check the intercoder reliability. Cohen κ [31] was used to measure intercoder reliability, which reached κ=0.825 after the final revision. Lastly, during the formal coding phase, four coders were trained and divided into two pairs of coders (YS and JP, and QW and YZ). Each pair independently coded a subset of tweets/posts, with a third coder (QW or YS) checking and resolving any disagreements.

A descriptive analysis was performed to show the percentage of topics for both Twitter and Weibo data. The time trends were plotted for percentages of tweets or Weibo posts with various attitudes toward general COVID-19 tests by week. Geographical distributions of post numbers and percentages of oppositional attitudes across the United Kingdom and mainland China were plotted by regions or provinces. The chi-square test was used to determine differences in attitudes or behaviors toward COVID-19 by gender, age, and education.

## *Results*

### Epidemic Context: Daily New COVID-19 Case Numbers in the United Kingdom and China

From June 1 to July 15, 2020, daily new COVID-19 confirmed cases demonstrated a decreasing trend in the United Kingdom, stabilizing at around 50 in London (Figure 1). With the decrease in new cases, the COVID-19 Alert Level in the United Kingdom was downgraded from level 4 to level 3 on June 19, representing that the COVID-19 epidemic was in general circulation with a demonstrable reduction in the number of cases and deaths [32]. On June 29, Leicester became the first city in the United Kingdom to undergo a local lockdown after a resurgence of cases. Concurrently, daily new cases in China fluctuated under 60, with the majority being in Beijing. There were no more local confirmed cases in Beijing after July 6.

**Figure 1.** Numbers of daily new COVID-19 cases in the United Kingdom (UK), London, mainland China, and Beijing from June 1 to July 15, 2020 [26,27].



## Social Listening: Public Attitudes Toward COVID-19 Testing

### Overall Analysis

In the United Kingdom, 64.6% (2390/3700) of tweets across the country and 69.2% (520/751) from London showed dominant views on individual COVID-19 tests in general. Moreover, 30.7% (1136/3700) of tweets from the United Kingdom and 22.6% (170/751) from London showed negotiated views on individual tests, while 4.7% (174/3700) across the country and 8.1% (61/751) from London opposed it (Figure 2). In China, about 30% of posts (country wide: 2649/8879, 29.8%; Beijing: 848/2991, 28.3%) showed dominant views on COVID-19 tests in general. Moreover, over 60% of posts (country wide: 5454/8879, 61.4%; Beijing: 1839/2991, 61.5%) showed negotiated views, and 10% or less (country wide: 776/8879, 8.7%; Beijing: 304/2991, 10.2%) opposed it. For example,

tweets/posts with dominant views on individual COVID-19 tests included "*need larger testing capacity and faster results*," negotiated tweets/posts included "*does one need to have a covid-19 test before travelling to the US*," and oppositional tweets/posts included "*30% of negative coronavirus tests are wrong.*" Individuals in the United Kingdom (23/2075, 1.1%) and London (6/530, 1.1%) showed less opposition to government-led community-wide mass COVID-19 testing than did those in mainland China (76/1594, 4.8%) and Beijing (58/661, 8.8%). For example, tweets/posts with dominant views on community-wide mass COVID-19 testing included "*support NHS and care workers routine weekly COVID-19 tests*" and oppositional tweets/posts included "*why would people with no symptoms take a test that tells them they're sick.*" A total of 2487 tweets from the United Kingdom expressed discontent with governmental COVID-19 testing practices, including taxing tests, voting against routine testing for front-line workers, publishing wrong or untimely data, and other complaints. In

China, 304 Weibo posts questioned the necessity of having to     obtain test results before travelling, visiting doctors, etc.

**Figure 2.** Percentage of tweets or Weibo posts with attitudes toward individual COVID-19 tests and community-wide tests from June 1 to July 15, 2020. UK: United Kingdom.



### Time Trend Analysis and Geographical Mapping

Time trend analysis (Figure 3) showed that, in the United Kingdom, posts with dominant attitudes toward individual COVID-19 tests first increased from 54.5% (533/978) to 86.5% (1066/1233) and then dropped to 38.6% (180/466) after the enactment of The Health Protection (Coronavirus) Regulations in July 2020 [33], while tweets with oppositional attitudes increased from 1.9% (19/978) to 16.3% (76/466). In China, Weibo posts with dominant attitudes reduced from 54.1% (380/703) to 7.3% (30/412) during this period, while those with negotiated attitudes increased from 33.9% (238/703) to 85.7%

(353/412) and posts with oppositional attitudes slightly dropped from 12.1% (85/703) to 7.0% (29/412). Regional analyses (Figure 4) showed that the percentage of tweets/posts in opposition to testing generally corresponded with low cases in their respective regions, with the exceptions of London (64/852, 7.5%) and the East Midlands (8/157, 5.1%). Oppositional tweets mostly worried about false negative testing results and that someone could get infected after testing negative, leading to more cases. Weibo posts from Beijing showed a slightly higher level of oppositional attitudes than London (304/3155, 9.6%), mostly questioning the cost-effectiveness of implementing mass testing when daily new cases in China fluctuated under 60.

**Figure 3.** Percentage of tweets or Weibo posts with attitudes toward COVID-19 testing by week and daily new cases from June 1 to July 15, 2020. UK: United Kingdom.

**Figure 4.** Percentage of tweets/Weibo posts with oppositional attitude toward COVID-19 testing and number of cases by geographical distribution in the United Kingdom (UK) and mainland China from June 1 to July 15, 2020 (regions with less than 50 tweets/posts are not shown).



## Self-Reported Uptake of COVID-19 Tests

Overall, 4.6% (178/3856) of tweets across the United Kingdom and 4.9% (39/794) from London reported intending to undergo or having undergone COVID-19 tests (Table 1). In the United Kingdom, driving forces for undergoing testing included personal health needs due to possible exposure, symptoms, or worry (37/86, 43%), mandatory testing policies for travel (30/86, 35%), and mass community-wide testing led by the government (19/86, 22%). Comparatively, more Weibo users reported having undergone COVID-19 testing in China (3318/9299, 35.7%) and Beijing (1462/3155, 46.3%). A total of 1600 Weibo posts (1600/9299, 17.2%) from China reported driving forces for undergoing testing, including community-wide testing led by the government (784/1600, 49.0%), mandatory testing policies for travel (659/1600, 41.2%), and personal health needs

(163/1600, 10.2%). Government-led community-wide testing was reported to be the main driving force for undergoing testing in Beijing (543/824, 65.9%).

**Table 1.** Uptake of COVID-19 tests in the United Kingdom and mainland China.

| Uptake and driving forces | United Kingdom (N=3856), n (%) | London (N=794), n (%) | China (N=9299), n (%) | Beijing (N=3155), n (%) |
|---|---|---|---|---|
| **Self-reported uptake of COVID-19 tests** | 178 (4.6) | 39 (4.9) | 3318 (35.7) | 1462 (46.3) |
| Plan to take a test | 35 (0.9) | 7 (0.9) | 811 (8.7) | 417 (13.2) |
| Have taken a test | 143 (3.7) | 32 (4.0) | 2507 (27.0) | 1045 (33.1) |
| **Driving force for taking a COVID-19 test** | 86 (2.2) | 25 (3.2) | 1600 (17.2) | 824 (26.1) |
| Personal health needs | 37 (1.0) | 14 (1.8) | 163 (1.8) | 64 (2.0) |
| Mandatory testing policies for travel | 30 (0.8) | 9 (1.1) | 659 (7.1) | 223 (7.1) |
| Community-wide testing by governments | 19 (0.5) | 2 (0.3) | 784 (8.4) | 543 (17.2) |
| Others taking COVID-19 tests | 86 (2.2) | 16 (2.0) | 277 (3.0) | 102 (3.2) |

## Major Factors of COVID-19 Testing Hesitancy

### Convenience: Access to and Experience With COVID-19 Tests

In the United Kingdom, 1.1% (43/3856) of tweets shared their experiences of undergoing COVID-19 tests, of which, 69.8% (30/43) reported being uncomfortable, 16.3% (7/43) reported being nervous, and 14.0% (6/43) reported no discomfort (Table 2). Comparatively, more Weibo posts shared the overall experience of taking a COVID-19 test from China (753/9299, 8.1%) and Beijing (221/3155, 7.0%). Furthermore, 9.2% (356/3856) of tweets in the United Kingdom and 16.6% (132/794) of tweets in London discussed the logistical process of obtaining a test, including access to an appointment, wait time to undergo testing, including queues and heatstroke while waiting, wait time for test results, and others. Discussions about the logistical process of obtaining a test in China (2005/9299, 21.6%) and Beijing (754/3155, 23.9%) mostly about the wait time to undergo testing received significant attention (China: 1579/2005, 78.8%; Beijing: 557/754, 73.9%).

The price of COVID-19 tests was mostly mentioned in London (53/794, 6.7%) across the United Kingdom (96/3856, 2.5%). These discussions included calls for free testing and dissatisfaction with high expenses, complaints about self-paying without reimbursement by medical insurance, and satisfaction with free-of-charge testing when available. In China, 6.4% (598/9299) of posts across the country and 6.2% (196/3155) from Beijing discussed the price of COVID-19 tests. The discussions included satisfaction with free-of-charge community-wide testing by governments and different attitudes toward self-paying prices, being either reasonable or burdensome. Regarding priority groups for COVID-19 testing, tweeters in the United Kingdom mentioned support for weekly testing for NHS staff and care staff and calling for priority testing to be extended. In China, posts concerning priority groups for COVID-19 testing referred to support for people with possible risk exposure, taxi drivers, and couriers to receive priority testing.

**Table 2.** Convenience of COVID-19 tests in the United Kingdom and mainland China.

| Convenience of COVID-19 tests | United Kingdom (N=3856), n (%) | London (N=794), n (%) | China (N=9299), n (%) | Beijing (N=3155), n (%) |
|---|---|---|---|---|
| **Experience of taking a COVID-19 test** | 43 (1.1) | 8 (1.0) | 753 (8.1) | 221 (7.0) |
| Feel uncomfortable | 30 (0.8) | 5 (0.6) | 476 (5.1) | 118 (3.7) |
| Feel nervous | 7 (0.2) | 2 (0.3) | 234 (2.5) | 70 (2.2) |
| Do not feel uncomfortable | 6 (0.2) | 1 (0.1) | 111 (1.2) | 43 (1.4) |
| **Logistical process of obtaining a COVID-19 test** | 356 (9.2) | 132 (16.6) | 2005 (21.6) | 754 (23.9) |
| Access to an appointment | 25 (0.6) | 8 (1.0) | 93 (1.0) | 52 (1.6) |
| Wait time to take a test | 53 (1.4) | 19 (2.4) | 1579 (17.0) | 557 (17.7) |
| Wait time for the test result | 82 (2.1) | 10 (1.3) | 174 (1.9) | 114 (3.6) |
| Others | 211 (5.5) | 99 (12.5) | 186 (2.0) | 48 (1.5) |
| Tribute to medical staff | 69 (1.8) | 15 (1.9) | 1914 (20.6) | 480 (15.2) |
| Price of COVID-19 testing | 96 (2.5) | 53 (6.7) | 598 (6.4) | 196 (6.2) |
| Priority groups for COVID-19 testing | 665 (17.2) | 9 (1.1) | 271 (2.9) | 101 (3.2) |

## Confidence and Complacency Toward COVID-19 Tests

In both the United Kingdom and China, social media users expressed a high perceived risk of the COVID-19 pandemic (United Kingdom: 139/152, 91.4%; China: 789/813, 97.0%) (Table 3). Moreover, 4.3% (164/3856) of tweets across the United Kingdom and 6.2% (49/794) from London concerned the effectiveness of COVID-19 tests, and of them, 20.7% (34/164) across the United Kingdom and 32.7% (16/49) from

London expressed confidence in its effectiveness, while 79.3% (130/164) across the United Kingdom and 67.3% (33/49) from London expressed doubts. In China, 4.5% (417/9299) of posts across the country and 5.2% (164/3155) from Beijing concerned test effectiveness, and of them, 65.9% (275/417) across China and 62.2% (102/164) from Beijing expressed confidence in its effectiveness, while 34.1% (142/417) across China and 37.8% (62/164) from Beijing expressed doubts.

**Table 3.** Confidence and complacency toward COVID-19 tests in the United Kingdom and mainland China.

| Confidence and complacency | United Kingdom (N=3856), n (%) | London (N=794), n (%) | China (N=9299), n (%) | Beijing (N=3155), n (%) |
|---|---|---|---|---|
| **Confidence: trust in COVID-19 tests** | | | | |
| **Concern on the effectiveness of COVID-19 tests** | 164 (4.3) | 49 (6.2) | 417 (4.5) | 164 (5.2) |
| Trust tests to be effective | 34 (0.9) | 16 (2.0) | 275 (3.0) | 102 (3.2) |
| Doubt the effectiveness of tests | 130 (3.4) | 33 (4.2) | 142 (1.5) | 62 (2.0) |
| Expiration date of COVID-19 tests | 1 (0.03) | 0 (0) | 87 (0.9) | 42 (1.3) |
| Incidental risks due to COVID-19 tests | 14 (0.4) | 3 (0.4) | 206 (2.2) | 113 (3.6) |
| **Complacency: perception of COVID-19 risk** | | | | |
| **Perception of COVID-19 risk** | 152 (3.9) | 45 (5.7) | 813 (8.7) | 299 (9.5) |
| High risk | 139 (3.6) | 44 (5.5) | 789 (8.5) | 286 (9.1) |
| Low risk | 13 (0.3) | 1 (0.1) | 24 (0.3) | 13 (0.4) |

## Communication: Information Inquiries and Rumors Related to COVID-19 Tests

Information inquiries and rumors about COVID-19 tests (Table 4) could be found on both UK and Chinese social media platforms (United Kingdom: 75/3856, 1.9% and 55/3856, 1.4%; China: 517/9299, 5.6% and 192/9299, 2.1%, respectively). The main information inquiries (22/75, 29.3%) mentioned in the United Kingdom about COVID-19 testing were "*how many people have received a test*" and "*delay in sharing testing data with English councils*," while in China, 48.2% (249/517) of information inquiries were "*whether tests are needed before travelling somewhere*" and "*how much it costs to take a test*."

Concerns included the duration (expiration) of test results and incidental risks associated with COVID-19 testing, such as cross-infection and threat of asymptomatic infection from crowd gathering. Posts mentioning unproven expositions about or interpretations of COVID-19 testing–related news, events, or problems were labelled as rumors (including fake news and misinformation). For example, rumors in the United Kingdom included "*COVID-19 test results were falsified*" and fake news included "*4 Tory MPs voted against weekly COVID-19 tests for NHS and care staff*." Comparatively, rumors in China included "*medical staff earned money by COVID-19 tests*" and fake news included "*positive testing results for [person name] in [place]*."

**Table 4.** Communication around COVID-19 tests in the United Kingdom and mainland China.

| Communication around COVID-19 tests | United Kingdom (N=3856), n (%) | London (N=794), n (%) | China (N=9299), n (%) | Beijing (N=3155), n (%) |
|---|---|---|---|---|
| Information inquiries about COVID-19 tests | 75 (1.9) | 21 (2.6) | 517 (5.6) | 137 (4.3) |
| Rumors about COVID-19 tests | 55 (1.4) | 12 (1.5) | 192 (2.1) | 27 (0.9) |

## Attitude and Uptake of COVID-19 Tests by the Characteristics of Social Media Posts

Tables 5 and 6 show the univariate analysis of the attitude and uptake of COVID-19 tests across sociodemographic characteristics using Weibo data from China. Male Weibo users and those over 30 years old were more likely to have a positive attitude toward individual COVID-19 testing (*P*<.001), but male

users were less likely to have a positive attitude toward mass community-wide COVID-19 testing led by the government (*P*<.001) (Table 5). Additionally, females, users under 30 years old, and those with a bachelor's degree or higher were more likely to take a COVID-19 test (*P*≤.001), and there was no significant difference in the reasons for undergoing COVID-19 testing (Table 6).

XSL•FO
RenderX

**Table 5.** Attitude toward COVID-19 tests by characteristics for Chinese Weibo posts.

| Characteristic | Attitude toward individual COVID-19 tests | | | | P value | Attitude toward community-wide COVID-19 tests | | | P value |
|---|---|---|---|---|---|---|---|---|---|
| | Total | Dominant, n (%) | Negotiated, n (%) | Oppositional, n (%) | | Total | Dominant, n (%) | Oppositional, n (%) | |
| **Gender** | | | | | <.001 | | | | <.001 |
| Male | 2844 | 989 (34.8) | 1587 (55.8) | 268 (9.4) | | 660 | 614 (93.0) | 46 (7.0) | |
| Female | 5940 | 1626 (27.4) | 3814 (64.2) | 500 (8.4) | | 910 | 882 (96.9) | 28 (3.1) | |
| Total | 8784 | 2615 (29.8) | 5401 (61.5) | 768 (8.7) | | 1570 | 1496 (95.3) | 74 (4.7) | |
| **Age (years)** | | | | | <.001 | | | | .05 |
| 10-30 | 2705 | 680 (25.1) | 1793 (66.3) | 232 (8.6) | | 356 | 344 (96.6) | 12 (3.4) | |
| 30-90 | 1125 | 420 (37.3) | 602 (53.5) | 103 (9.2) | | 269 | 251 (93.3) | 18 (6.7) | |
| Total | 3830 | 1100 (28.7) | 2395 (62.5) | 335 (8.7) | | 625 | 595 (95.2) | 30 (4.8) | |
| **Education** | | | | | .49 | | | | .90 |
| Bachelor's degree or above | 2297 | 704 (30.6) | 1387 (60.4) | 206 (9.0) | | 472 | 450 (95.3) | 22 (4.7) | |
| High school or below | 6582 | 1945 (29.6) | 4067 (61.8) | 570 (8.7) | | 1122 | 1068 (95.2) | 54 (4.8) | |
| Total | 8879 | 2649 (29.8) | 5454 (61.4) | 776 (8.7) | | 1594 | 1518 (95.2) | 76 (4.8) | |

**Table 6.** Uptake of COVID-19 tests by characteristics for Chinese Weibo posts.

| Characteristic | Uptake of COVID-19 tests | | | P value | Driving force for taking a COVID-19 test | | | | P value |
|---|---|---|---|---|---|---|---|---|---|
| | Total | Yes, n (%) | No, n (%) | | Total | Personal health needs, n (%) | Mandatory testing policies for travel, n (%) | Community-wide tests, n (%) | |
| **Gender** | | | | <.001 | | | | | .44 |
| Male | 2844 | 808 (28.4) | 2036 (71.6) | | 413 | 37 (9.0) | 167 (40.4) | 209 (50.6) | |
| Female | 5940 | 2470 (41.6) | 3470 (58.4) | | 1170 | 126 (10.8) | 487 (41.6) | 557 (47.6) | |
| Total | 8784 | 3278 (37.3) | 5506 (62.7) | | 1583 | 163 (10.3) | 654 (41.3) | 766 (48.4) | |
| **Age (years)** | | | | <.001 | | | | | .09 |
| 10-30 | 2705 | 1187 (43.9) | 1518 (56.1) | | 506 | 60 (11.9) | 229 (45.3) | 217 (42.9) | |
| 30-90 | 1125 | 316 (28.1) | 809 (71.9) | | 193 | 22 (11.4) | 71 (36.8) | 100 (51.8) | |
| Total | 3830 | 1503 (39.2) | 2327 (60.8) | | 699 | 82 (11.7) | 300 (42.9) | 317 (45.4) | |
| **Education** | | | | .001 | | | | | .16 |
| Bachelor's degree or above | 2297 | 922 (40.1) | 1375 (59.9) | | 444 | 45 (10.1) | 166 (37.4) | 233 (52.5) | |
| High school or below | 6582 | 2395 (36.4) | 4187 (63.6) | | 1162 | 118 (10.2) | 493 (42.4) | 551 (47.4) | |
| Total | 8879 | 3317 (37.4) | 5562 (62.6) | | 1606 | 163 (10.1) | 659 (41.0) | 784 (48.8) | |

# Discussion

## Principal Findings

This infodemiology study assessed public attitudes and opinions around COVID-19 testing, including both individual and government-led mass testing, by monitoring and analyzing digital conversations in the United Kingdom (Twitter) and China (Sina Weibo) with a framework of testing hesitancy (confidence, complacency, convenience, and communication). Overall, there was a higher level of support toward individual and mass COVID-19 testing in the United Kingdom and London than in mainland China and Beijing; most opposition originated from the capital cities. Time trend analyses showed that discussions about individual COVID-19 tests were mostly dominant in the United Kingdom, while Weibo posts in China showed a rise of negotiated views over testing. There were much wider reports of public uptake of COVID-19 tests in mainland China than in the United Kingdom. Personal health needs (eg, possible exposure, symptoms, and worry), mandatory testing policies

for work or travel, and government-led mass testing were the main driving forces for people to undergo testing in both countries, with differences in priorities between countries. The Chinese public posted more about uncomfortable experiences and logistical barriers to testing, whereas people in the United Kingdom posted more about prices and priority groups for testing. Perceived risk of the COVID-19 disease was high in both countries. Only 5% or less of the posts discussed test effectiveness, and of them, Chinese users expressed confidence in its effectiveness, whereas British users displayed doubts. Rumors related to COVID-19 test administration and results were identified. In China, females, those under 30 years old, and those with a bachelor's degree or higher were more likely to undergo a COVID-19 test.

Overall, discussions about COVID-19 in both countries showed low complacency (high perceived risk) for the COVID-19 disease and high confidence in testing, which translated into high levels of public support for testing. The cited driving forces for testing (personal health needs in the United Kingdom versus government-led mass testing or mandatory testing policies for travel in China) also reflected the epidemic situation and testing policies implemented in each respective context. Our data showed that, as daily new cases decreased and COVID-19 testing became routine, a negotiated position toward COVID-19 testing became the majority view, leading to an increase in acceptance and uptake behavior when needed. Epidemiologists have argued that widespread dissemination of cheap and rapid tests might be as effective as a vaccine at interrupting coronavirus transmission by identifying and isolating people with the virus when they are most infectious [1,34]. Integrating "complacency (risk perception)" of the disease and "confidence" of testing in messaging by emphasizing the importance of timely testing during an acute epidemic could increase acceptance and uptake.

"Convenience" of testing, including accessibility, frequency, and sample-to-answer time, was a popular topic of digital discussion and also one of the most important factors for effective screening, being an even higher priority than the analytical limits of detection [34]. Inquiries and rumors related to COVID-19 testing pointed to the lack of a frequent and factually correct information campaign. Furthermore, regional analysis showed an association between opposition views toward testing and low case counts, with the exceptions of London and the East Midlands, mostly because of worrying about false negative testing results and worrying that someone could get infected after testing negative, leading to more cases. Inquiries, concerns, and rumors identified during social listening call for rapid communication responses. These findings demonstrate a need for effective emergency risk communication strategies during a public health crisis that are informed by real-time evidence derived from ongoing social listening and tailored to local social and epidemic contexts. These strategies should not only meet immediate public information needs, but also debunk rumors and misinformation as they emerge.

Our data showed how local epidemic situations influenced public attitudes toward COVID-19 testing and highlighted the challenges facing governments when weighing the balance between epidemic control and socioeconomical livelihoods.

This study was performed when the United Kingdom was under its first nationwide lockdown, while China had resumed complete normalcy since late March 2020. Tweets from London and the United Kingdom showed overwhelming support for both strategies, whereas more Weibo users expressed negotiated or oppositional positions of mass community-wide testing. In the United Kingdom, the government has long been criticized for being underprepared for the COVID-19 pandemic, including lacking testing capacity for both the general public and frontline workers [9,35]. Without reliable test results, very limited data were available to develop and introduce an exit strategy for the general lockdown, as health experts had no evidence to inform their decisions. After its implementation, the NHS Test and Trace program was widely criticized over the lack of convenience (ie, pricing and accessibility) and the exacerbation of COVID-19 inequalities, resulting in a new campaign being launched on July 30, 2020 [36] to encourage everyone with symptoms to undergo free testing. In China, comprehensive testing requirements around domestic travel have been in place since March 2020, when the country lifted its nationwide lockdown. Between June 12 and 22, 2020, the Beijing government led mass community-wide testing, with 2.95 million tests completed in 10 days due to a small resurgence of cases. Despite high perceived risk toward COVID-19, some Chinese residents questioned the overall cost-effectiveness of implementing a massive measure against such few cases. Public attitudes and sentiment constantly evolve with local epidemic situations, and as such, public communication about health risks and countermeasures must leverage real-time social listening and disease surveillance data to keep up.

There are geographic differences in public attitudes toward mass community-wide COVID-19 testing, with more people in London and Beijing, and adjacent areas opposing testing. In China, strong evidence indicates gender differences in attitudes toward mass community-wide COVID-19 testing, with more female residents supporting government-led testing. Consistent with the audience reception theory, our data showed that the public is a diverse heterogeneous set of people with varying experiences and needs. They access, process, and react to messages differently based on their individual backgrounds and views. Tailoring engagement strategies to the target community will be critical in increasing acceptance toward COVID-19 testing and other containment measures.

## Limitations

This study has some limitations. First, there is an inherent bias shared among all studies that utilize social media data, where users might present themselves differently online (eg, inflated perception) and/or represent a skewed younger population [37]. Nevertheless, findings from this study had very limited influence by curated perceptions as the investigation mainly focused on how aggregated social media data constituted a dynamic digital environment regarding COVID-19 testing, and how such an information environment affected individuals' acceptance of control measures during a pandemic. Moreover, this study captured routine data from populations that may not be represented in traditional research designs. The opinions gathered via social listening could be less biased than those derived from traditional research methods, such as surveys and

interviews, where unintended errors could be introduced by how questionnaires were presented and implemented (eg, reporting bias and acquiescence bias). Second, we were unable to extract demographic data from all Twitter and some Weibo profiles due to privacy restrictions, and the authenticity of the retrieved data was not directly verifiable. We therefore conducted account analyses of available Weibo profile data to assess data representativeness, which indicated a satisfactory level of comparability to the data in the official Weibo report. Third, data were downloaded daily to avoid the possible interference of comment removal by authorities. Lastly, the findings from this study are mostly exploratory and might not be generalizable due to the small sample size of posts reviewed (approximately 10%). A further investigation employing machine learning algorithms for big data analysis is needed.

## Conclusion

Policy makers tackling factors of COVID-19 testing hesitancy should focus on complacency, confidence, convenience, and communication in relation to testing. There is a need for more comparative studies to identify differences and similarities across populations and experiences with COVID-19 testing. Future infodemiology studies should integrate public and epidemic data (eg, traditional media, social media, polls, and disease surveillance data), both online and offline, and employ machine learning to enable rapid real-time analysis of big data for epidemic preparedness and response.

## Authors' Contributions

ZH and LL conceptualized and designed the study. FYS and XZ collected the data. YS, QW, JP, and YZ analyzed the data. LL, YS, and QW drafted the manuscript. ZH, LL, and HJL contributed to the critical revision of the manuscript for important intellectual content. All authors approved the final manuscript as submitted and agree to be accountable for all aspects of the work.

## Conflicts of Interest

HJL is on the Merck Vaccine Confidence Advisory Board. Her research group, the Vaccine Confidence Project, received research grants from GSK and Merck on vaccine confidence issues. None of those research grants are related to this paper. The other authors have no conflicts to declare.

Multimedia Appendix 1
Workflow of the inclusion and exclusion process of tweets and Weibo posts.
[PNG File , 161 KB - infodemiology_v1i1e26895_app1.png ]

Multimedia Appendix 2
Account analysis of COVID-19 test posts on social media.
[DOCX File , 17 KB - infodemiology_v1i1e26895_app2.docx ]

Multimedia Appendix 3
Coding framework for COVID-19 test posts on social media.
[DOCX File , 24 KB - infodemiology_v1i1e26895_app3.docx ]

## References

1. Mina MJ, Parker R, Larremore DB. Rethinking Covid-19 Test Sensitivity - A Strategy for Containment. N Engl J Med 2020 Nov 26;383(22):e120 [FREE Full text] [doi: 10.1056/NEJMp2025631] [Medline: 32997903]
2. Grassly NC, Pons-Salort M, Parker EPK, White PJ, Ferguson NM, Ainslie K, et al. Comparison of molecular testing strategies for COVID-19 control: a mathematical modelling study. The Lancet Infectious Diseases 2020 Dec;20(12):1381-1389. [doi: 10.1016/s1473-3099(20)30630-7]
3. Sharpe HR, Gilbride C, Allen E, Belij-Rammerstorfer S, Bissett C, Ewer K, et al. The early landscape of coronavirus disease 2019 vaccine development in the UK and rest of the world. Immunology 2020 Jul;160(3):223-232 [FREE Full text] [doi: 10.1111/imm.13222] [Medline: 32460358]

XSL·FO

RenderX

4.  Wong LP, Alias H, Wong P, Lee HY, AbuBakar S. The use of the health belief model to assess predictors of intent to receive the COVID-19 vaccine and willingness to pay. Hum Vaccin Immunother 2020 Sep 01;16(9):2204-2214 [FREE Full text] [doi: 10.1080/21645515.2020.1790279] [Medline: 32730103]

5.  Timeline of new outbreak in Xinfadi, Beijing. Beijing Daily. URL: https://baijiahao.baidu.com/s?id=1674327183099264211&wfr=spider&for=pc [accessed 2020-08-07]

6.  Scally G, Jacobson B, Abbasi K. The UK's public health response to covid-19. BMJ 2020 May 15;369:m1932. [doi: 10.1136/bmj.m1932] [Medline: 32414712]

7.  Iacobucci G. Covid-19: What is the UK's testing strategy? BMJ 2020 Mar 26;368:m1222. [doi: 10.1136/bmj.m1222] [Medline: 32217754]

8.  García LY, Cerda AA. Contingent assessment of the COVID-19 vaccine. Vaccine 2020 Jul 22;38(34):5424-5429 [FREE Full text] [doi: 10.1016/j.vaccine.2020.06.068] [Medline: 32620375]

9.  Vize R. Too slow and fundamentally flawed: why test and trace is a weak and inequitable defence against covid-19. BMJ 2020 Jun 11;369:m2246. [doi: 10.1136/bmj.m2246] [Medline: 32527794]

10. McDermott JH, Newman WG. Refusal of viral testing during the SARS-CoV-2 pandemic. Clin Med (Lond) 2020 Sep 03;20(5):e163-e164 [FREE Full text] [doi: 10.7861/clinmed.2020-0388] [Medline: 32620593]

11. Eysenbach G. Infodemiology: The epidemiology of (mis)information. Am J Med 2002 Dec 15;113(9):763-765. [doi: 10.1016/s0002-9343(02)01473-0] [Medline: 12517369]

12. Eysenbach G. Infodemiology and infoveillance: framework for an emerging set of public health informatics methods to analyze search, communication and publication behavior on the Internet. J Med Internet Res 2009 Mar 27;11(1):e11 [FREE Full text] [doi: 10.2196/jmir.1157] [Medline: 19329408]

13. Tang L, Bie B, Park S, Zhi D. Social media and outbreaks of emerging infectious diseases: A systematic review of literature. Am J Infect Control 2018 Sep;46(9):962-972 [FREE Full text] [doi: 10.1016/j.ajic.2018.02.010] [Medline: 29628293]

14. Choi S, Lee J, Kang M, Min H, Chang Y, Yoon S. Large-scale machine learning of media outlets for understanding public reactions to nation-wide viral infection outbreaks. Methods 2017 Oct 01;129:50-59 [FREE Full text] [doi: 10.1016/j.ymeth.2017.07.027] [Medline: 28813689]

15. Hall S. Encoding, Decoding. In: During S, editor. The Cultural Studies Reader. London, UK: Routledge; 1993:90-103.

16. Hou Z, Tong Y, Du F, Lu L, Zhao S, Yu K, et al. Assessing COVID-19 Vaccine Hesitancy, Confidence, and Public Engagement: A Global Social Listening Study. J Med Internet Res 2021 Jun 11;23(6):e27632 [FREE Full text] [doi: 10.2196/27632] [Medline: 34061757]

17. Barkur G, Vibha, Kamath GB. Sentiment analysis of nationwide lockdown due to COVID 19 outbreak: Evidence from India. Asian J Psychiatr 2020 Jun;51:102089 [FREE Full text] [doi: 10.1016/j.ajp.2020.102089] [Medline: 32305035]

18. Đogaš Z, Lušić Kalcina L, Pavlinac Dodig I, Demirović S, Madirazza K, Valić M, et al. The effect of COVID-19 lockdown on lifestyle and mood in Croatian general population: a cross-sectional study. Croat Med J 2020 Aug;61(4):309-318. [doi: 10.3325/cmj.2020.61.309]

19. Doogan C, Buntine W, Linger H, Brunt S. Public Perceptions and Attitudes Toward COVID-19 Nonpharmaceutical Interventions Across Six Countries: A Topic Modeling Analysis of Twitter Data. J Med Internet Res 2020 Sep 03;22(9):e21419 [FREE Full text] [doi: 10.2196/21419] [Medline: 32784190]

20. Hou Z, Du F, Zhou X, Jiang H, Martin S, Larson H, et al. Cross-Country Comparison of Public Awareness, Rumors, and Behavioral Responses to the COVID-19 Epidemic: Infodemiology Study. J Med Internet Res 2020 Aug 03;22(8):e21143 [FREE Full text] [doi: 10.2196/21143] [Medline: 32701460]

21. Meltwater. URL: https://www.meltwater.com/en [accessed 2021-08-15]

22. Weibo User Development Report 2020. Weibo Data Center. URL: https://hd.weibo.com/article/view/3982 [accessed 2021-03-12]

23. Bacsu J, O'Connell ME, Cammer A, Azizi M, Grewal K, Poole L, et al. Using Twitter to Understand the COVID-19 Experiences of People With Dementia: Infodemiology Study. J Med Internet Res 2021 Feb 03;23(2):e26254 [FREE Full text] [doi: 10.2196/26254] [Medline: 33468449]

24. Chew C, Eysenbach G. Pandemics in the age of Twitter: content analysis of Tweets during the 2009 H1N1 outbreak. PLoS One 2010 Nov 29;5(11):e14118 [FREE Full text] [doi: 10.1371/journal.pone.0014118] [Medline: 21124761]

25. Liao Q, Yuan J, Dong M, Yang L, Fielding R, Lam WWT. Public Engagement and Government Responsiveness in the Communications About COVID-19 During the Early Epidemic Stage in China: Infodemiology Study on Social Media Data. J Med Internet Res 2020 May 26;22(5):e18796 [FREE Full text] [doi: 10.2196/18796] [Medline: 32412414]

26. Daily briefing on novel coronavirus cases in China. National Health Commission of China. URL: http://www.nhc.gov.cn/xcs/yqtb/list_gzbd.shtml [accessed 2021-08-16]

27. Cases tested positive for coronavirus in United Kingdom. Public Health England. URL: https://coronavirus.data.gov.uk/details/cases [accessed 2021-08-16]

28. Report of the SAGE Working Group on Vaccine Hesitancy. World Health Organization. URL: https://www.who.int/immunization/sage/meetings/2014/october/1_Report_WORKING_GROUP_vaccine_hesitancy_final.pdf [accessed 2021-08-16]

29. Braun V, Clarke V. Using thematic analysis in psychology. Qualitative Research in Psychology 2006 Jan;3(2):77-101. [doi: 10.1191/1478088706qp063oa]

30. Baccianella S, Esuli A, Sebastiani F. SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining. In: Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10). 2010 Presented at: Seventh International Conference on Language Resources and Evaluation (LREC'10); 2010; Valletta, Malta.

31. Cohen J. A Coefficient of Agreement for Nominal Scales. Educational and Psychological Measurement 2016 Jul 02;20(1):37-46. [doi: 10.1177/001316446002000104]

32. UK COVID-19 alert level methodology: an overview. UK Department of Health and Social Care. URL: https://www.gov.uk/government/publications/uk-covid-19-alert-level-methodology-an-overview/uk-covid-19-alert-level-methodology-an-overview [accessed 2021-08-06]

33. Lennon R, Small M, Smith R, Van Scoy LJ, Myrick J, Martin M, Data Action Research Group. Unique Predictors of Intended Uptake of a COVID-19 Vaccine in Adults Living in a Rural College Town in the United States. Am J Health Promot 2021 Jul 16:8901171211026132. [doi: 10.1177/08901171211026132] [Medline: 34269077]

34. Larremore DB, Wilder B, Lester E, Shehata S, Burke JM, Hay JA, et al. Test sensitivity is secondary to frequency and turnaround time for COVID-19 screening. Sci Adv 2021 Jan 01;7(1):eabd5393 [FREE Full text] [doi: 10.1126/sciadv.abd5393] [Medline: 33219112]

35. Briggs A, Jenkins D, Fraser C. NHS Test and Trace: the journey so far. The Health Foundation. URL: https://www.health.org.uk/publications/long-reads/nhs-test-and-trace-the-journey-so-far [accessed 2020-09-23]

36. NHS Test and Trace launches campaign to encourage everyone with symptoms to get a free test. UK Department of Health and Social Care. URL: https://www.gov.uk/government/news/nhs-test-and-trace-launches-campaign-to-encourage-everyone-with-symptoms-to-get-a-free-test [accessed 2020-07-30]

37. Golder S, O'Connor K, Hennessy S, Gross R, Gonzalez-Hernandez G. Assessment of Beliefs and Attitudes About Statins Posted on Twitter: A Qualitative Study. JAMA Netw Open 2020 Jun 01;3(6):e208953 [FREE Full text] [doi: 10.1001/jamanetworkopen.2020.8953] [Medline: 32584408]

## Abbreviations

**NHS:** National Health System

Original Paper

# The Impact of the Online COVID-19 Infodemic on French Red Cross Actors' Field Engagement and Protective Behaviors: Mixed Methods Study

Leonardo W Heyerdahl[1], PhD; Benedetta Lana[1], MA; Tamara Giles-Vernick[1], PhD

Department of Global Health, Anthropology and Ecology of Disease Emergence Unit, Institut Pasteur, Paris, France

**Corresponding Author:**
Tamara Giles-Vernick, PhD
Department of Global Health, Anthropology and Ecology of Disease Emergence Unit
Institut Pasteur
25 rue du Docteur Roux
Paris, 75015
France
Phone: 33 0140613982
Email: tamara.giles-vernick@pasteur.fr

## Abstract

**Background:** The COVID-19 pandemic has been widely described as an infodemic, an excess of rapidly circulating information in social and traditional media in which some information may be erroneous, contradictory, or inaccurate. One key theme cutting across many infodemic analyses is that it stymies users' capacities to identify appropriate information and guidelines, encourages them to take inappropriate or even harmful actions, and should be managed through multiple transdisciplinary approaches. Yet, investigations demonstrating how the COVID-19 information ecosystem influences complex public decision making and behavior offline are relatively few.

**Objective:** The aim of this study was to investigate whether information reported through the social media channel Twitter, linked articles and websites, and selected traditional media affected the risk perception, engagement in field activities, and protective behaviors of French Red Cross (FRC) volunteers and health workers in the Paris region of France from June to October 2020.

**Methods:** We used a hybrid approach that blended online and offline data. We tracked daily Twitter discussions and selected traditional media in France for 7 months, qualitatively evaluating COVID-19 claims and debates about nonpharmaceutical protective measures. We conducted 24 semistructured interviews with FRC workers and volunteers.

**Results:** Social and traditional media debates about viral risks and nonpharmaceutical interventions fanned anxieties among FRC volunteers and workers. Decisions to continue conducting FRC field activities and daily protective practices were also influenced by other factors unrelated to the infodemic: familial and social obligations, gender expectations, financial pressures, FRC rules and communications, state regulations, and relationships with coworkers. Some respondents developed strategies for "tuning out" social and traditional media.

**Conclusions:** This study suggests that during the COVID-19 pandemic, the information ecosystem may be just one among multiple influences on one group's offline perceptions and behavior. Measures to address users who have disengaged from online sources of health information and who rely on social relationships to obtain information are needed. Tuning out can potentially lead to less informed decision making, leading to worse health outcomes.

## Introduction

One critical concern emerging during the COVID-19 pandemic has been the infodemic, defined as excessive information that spreads rapidly, may be deliberately or inadvertently misleading, complicates emergency risk communication, and encourages lay publics to engage in harmful actions during this public health emergency [1-5]. Neither this phenomenon during epidemics

XSL•FO
RenderX

nor investigations of it are new. Infodemiology—"the study of the determinants and distribution of health information and misinformation"—emerged in the late 20th century and was shortly thereafter conceptualized as a field of study [6,7]. Infodemics as informational companions to epidemics developed with the first SARS epidemic and continued subsequently during the H1N1, Ebola, and Zika public health emergencies [8-10]. The World Health Organization Director General popularized the term in the context of the COVID-19 pandemic, declaring, "We are not just fighting an epidemic; we're fighting an infodemic" [11].

The current pandemic has catalyzed numerous social media–listening investigations to monitor, evaluate, and respond to circulating misinformation and disinformation during this infodemic [4,12-14]. In identifying problematic narratives and measuring their online spread, one key theme cutting across many such analyses is that the infodemic is a threat; it rapidly overwhelms users with contradictory and misleading information and encourages them to make risky or harmful decisions [5,6,15-19].

How this complex information ecosystem influences offline behavior—real-life choices and practices—is a critically important question during the COVID-19 pandemic, although it remains insufficiently investigated [20]. Media analyses, psychology, and anthropology have addressed this interaction between online information interpretation and offline behavior differently.

Drawing from the fields of communication, marketing, and computer science, social media analyses have detected the emergence and spread of COVID-19 misinformation and characterized narratives and concerns of users [21-24]. Such analyses can shed important light on public concerns about public health measures, yet they offer less insight into real-life practice [25,26]. Islam and colleagues found a correlation between online stigmatization of, and offline violence toward, Asian populations [27]. Although such correlations are compelling, they do not question the specific agents of this violence with the aim to understand how they obtained and interpreted online narratives and decided to act upon them.

Behavioral studies drawing from psychology often measure the influence of social and traditional media on psychological states or on health behaviors [28-34]. During the COVID-19 pandemic, these studies focused primarily on misinformation and disinformation, evaluating the language of media users or employing closed questionnaires to assess subjects' media use and offline emotions and behaviors [35-43]. For these researchers, attending to cultural differences (eg, around mask use) can illuminate divergent psychological states [44]. These studies yield important insights into the psychological impact of misinformation and disinformation. In concentrating on misinformation and disinformation, they characterize and aggregate individual responses, but risk neglecting underlying sociocultural, political, and economic conditions that may inform the emotional and behavioral responses of specific social groups.

Although not well known by other fields, anthropological contributions to evaluating rumors and infodemics are two-fold. First, over the past three decades, anthropologists and ethnohistorians addressing epidemic crises and biomedical research have been less preoccupied with distinguishing truth from rumors than with understanding specific social groups' perceptions and actions, and the factors shaping them [45-50]. Tappan [46] and Graboyes [47], in particular, contend that rumors of blood theft offer rich insight into East Africans' criticisms of biomedical research, in addition to reflections on the political, economic, and social inequalities that late colonialism imposed. This anthropological perspective on rumors, and more broadly an insistence on the valuable insights drawn from evaluating all information, rather than sifting out misinformation and disinformation from a broader body of circulating information, has informed both online ethnographies and research on the current COVID-19 infodemic [8,51-54]. A second contribution of anthropology, as well as other field-based qualitative social sciences, is a preoccupation with situating a social group's specific understandings and practices within their broader sociocultural, political, and economic contexts [55]. This preoccupation necessitates the use of multiple methods that permit anthropologists to capture what informants say, think, and do, and to gain insight into the multiple influences that shape those words, thoughts, and actions.

These two anthropological concerns frame our central question and our approach in this study. Our question focuses on how the COVID-19 infodemic has affected offline public behavior. We investigated whether and how social and traditional media influenced risk perceptions and behaviors of a heterogeneous social group in the Paris region of France: French Red Cross (FRC) volunteers and workers. Specifically, we conducted quantitative and qualitative social listening and analysis of Twitter, its links to other media, and selected print media in France from April to October 2020. This popular microblogging site, its links, and selected print media served as a proxy for public debates around COVID-19. Simultaneously, we evaluated how FRC volunteers and workers experienced and portrayed the influence of social and traditional media on their risk perceptions, their engagement in FRC field actions, and their protective behaviors.

## Methods

### Site and Population Description

This study has been carried out in the Ile-de-France region (Paris region) of France, a region comprised of eight departments that cover the city of Paris and its suburbs, with a population of 12.1 million people.

The population investigated in the study consisted of FRC salaried workers and volunteers carrying out operations in the Ile-de-France region. The FRC is a nongovernmental organization providing critical support for France's public health system during the COVID-19 pandemic, organizing and conducting diagnostic testing and emergency medical services. An estimated 42,800 volunteers in this region are mostly lay people undertaking social assistance activities (eg, food assistance; outreach for elderly, homeless, or other vulnerable populations; and support for school-age children), but a small proportion (less than 5%) serve as rescue workers, physicians, nurses, nursing assistants, and technicians performing

emergency and first aid response and providing health care. The FRC's 4200 salaried workers provide medical, material, and legal assistance to vulnerable populations in specialized centers and manage the logistics of FRC activities, as well as financial and other donations. The FRC's major presence in the humanitarian landscape makes it a useful organization to investigate during the pandemic.

The FRC was also selected because in early March 2020, prior to France's first lockdown, the French Red Cross Foundation contacted our research team, requesting that we investigate how the COVID-19 pandemic affected FRC workers and volunteers. Given that the COVID-19 infodemic was a serious global concern, we eventually decided to explore whether the changing information ecosystem affected workers' and volunteers' motivations: their decisions to participate in FRC field activities, as well as their self-protective measures.

## Data Collection

### Twitter Social Listening and Selected Media Evaluation

We collected data through Twitter and selected traditional media to capture a range of COVID-19 public health debates. Social and traditional media form part of an informational ecosystem: social and traditional media mutually influence one another and, at times, overlap [56].

From April to October 2020, we conducted social listening of COVID-19–related messages on Twitter and tracked COVID-19–related debates in traditional media. Using custom scripts in R (version 4.0.2; The R Foundation) and the rtweet package, we submitted daily keyword-based "coronavirus OR COVID" queries to the Twitter application programming interface. All matching organic tweets—not retweets—in French from accounts declaring a France-based location were collected.

Additional queries were added based on emerging themes from interviews or from traditional print media. These queries addressed lockdowns (keyword: "confinement"), hydroxychloroquine ("hydroxychloroquine"), vaccines ("vaccin"), and masks ("masque"). Such queries characterized narratives about epidemic events and identified scientific and public health debates emerging from interviews or from Twitter data.

To supplement this source, we also viewed linked sources (eg, articles, videos, radio programs, and websites) and followed three of the top five daily newspapers covering a range of political positions—*Le Monde*, *Libération,* and *Le Figaro*—to identify changing debates around COVID-19 [57]. Linked sources provided additional contextual information about the tweet and provided additional content on new debates. We also participated in an international WhatsApp group of researchers sharing diverse media sources from around the world to track new and ongoing pandemic debates. These combined sources served as a proxy for key public debates over public health measures and biomedical investigations; they supplemented our inquiries in Twitter and contributed specific questions about online and media debates concerning the pandemic in our qualitative interviews.

### Semistructured Interviews

From June to October 2020, we conducted 24 semistructured individual interviews with FRC volunteers and workers. The interview guide (Multimedia Appendix 1) addresses informants' training, their activities before the pandemic, and how these activities changed with pandemic emergence, lockdown measures, and deconfinement. Crucially, it explored participants' risk perceptions of COVID-19 and their decisions to participate or not in FRC activities. The interview guide also included questions about traditional and social media and the informant's use of them.

Our analyses of social listening and media tracking enabled us to ask specific interview questions about whether and how specific traditional and social media debates had influenced a participant's emotional responses, decisions, or practices. Even if individual participants did not use Twitter specifically, we nevertheless asked questions about their understanding and interpretation of current public debates over public health measures and biomedical investigations, as well as whether these debates influenced their risk assessments, their decisions to continue field FRC activities, and their adherence to protective measures. All interviews were conducted in French and recorded with informant consent.

### Recruitment

The FRC compiled a randomly selected database of 9000 volunteers and workers in the Ile-de-France region, stratified by proportion of volunteers and workers, gender, department, and age. We recruited interview participants by randomly selecting their names from this database and contacting them via FRC email.

## Data Analysis

### Twitter Social Listening and Media Tracking

We conducted weekly quantitative and qualitative analyses of collected tweets. We evaluated top hashtags, expressing them as a percentage of weekly totals. Our qualitative analysis involved thematic coding of a random sample of all tweets (100 per week), from which we would identify the most frequently mentioned debates to raise in interviews. This random selection enabled us to pick up tweets, including conspiracy-related narratives, that might otherwise have been flagged or filtered out by an algorithm before they could trend.

An initial coding grid for tweets contained four major themes: risk perception, control, interpretations, and key actors and groups. To examine Twitter users' perceptions of epidemic risk, we coded perceived viral origins, transmission and severity, and individual or social group susceptibility to the virus. For the control theme, we evaluated how Twitter users understood the efficacy, safety, and accessibility of preventive, diagnostic, and purportedly curative measures and devices (eg, masks, contact tracing, and distancing). The interpretations theme focused on narratives of viral origins and those profiteering from the pandemic. The key actors and groups theme categorized descriptions of specific actors or social groups.

Linked materials (eg, articles, other tweets, videos, and websites) were also evaluated for content. If the message alone was

insufficiently clear, coders evaluated contextual tweets, as well as titles and descriptions of linked material. A short synthesis of contextual material was noted with the code to support thematic analysis. Two coders reached a coding agreement of Cohen κ=0.65 [58]. Divergent coding occurred when coders evaluated users' perceptions of effectiveness of protective measures; in some cases, coders focused on the tweet's content, and in others, on its broader context and linked content. Another disagreement was related to users' attitudes toward protective measures (ie, those favoring or opposed to control measures and mandatory enforcement of masks or vaccines). Here, divergent coding resulted because control measures and mandatory enforcement (ie, masks, distancing, and limits on group sizes) were closely related. Coders discussed and, when necessary, modified code definitions to reflect consensus.

We consolidated biweekly social listening of pandemic discussions into short reports summarizing COVID-19 discussions. This analysis focuses primarily on the COVID-19 mask queries.

### Interviews

All interviews were transcribed and integrated into NVivo software (2020 version; QSR International). The first author conducted the thematic coding of the interviews, building on the Twitter codebook by adding categories related to FRC activities, FRC actors' motivation and engagement, and protective strategies. Team members wrote memos that formed the basis of our analysis by synthesizing coded content and detailing linkages across codes.

### Ethics

The protocol received ethical approval from the Institut Pasteur Institutional Review Board (IRB 2020-03) and was reviewed and approved by the FRC ethics committee. All study participants received an information notice and provided informed consent.

# Results

We collected and statistically analyzed 9,648,000 tweets, evaluated 1400 tweets qualitatively, and conducted 24 semistructured interviews with 8 FRC workers and 16 volunteers.

## Twitter Social Listening and Supplemental Media Tracking

### Overview

Through Twitter-based social listening and supplemental media tracking, we followed discussions about COVID-19 control, risk perceptions, key actors and groups, and interpretations of purported origins and those profiteering from the pandemic. The supplemental media tracking guided our inquiries into newly emerging concerns on Twitter.

### COVID-19 Risks

Between June and November 2020, many Twitter users discussed disease severity (n=127), most of whom (n=104) depicted COVID-19 as a dangerous disease, whereas a minority (n=23) described it as a mild disease affecting only the elderly;

they compared it to seasonal influenza. The pandemic's social consequences, particularly its social and economic costs and heightened exposures for marginalized social groups and frontline workers, were often emphasized (n=75). Uncertainties about the virus, specifically concerning transmission, mutations, and long-term consequences, also figured in discussions (n=59). A few users also shared diverging interpretations of the virus origins as being zoonotic, laboratory made, or 5G related (n=9).

### Masks as Control Measures

Debates about the utility of masks, their scarcity, and changing policies regarding their use made masks the most discussed control measure for the study period in our social media–listening data.

Concerning the utility of masks, several users (n=14) echoed authorities' and scientists' calls to wear masks. Certain users (n=10) evoked COVID-19 susceptibility and severity as justification for mask wearing (n=10), whereas others (n=14) did so on the grounds that masks could prevent transmission, with a few contending that masks "would not be enough."

Mask opponents, however, contested the legitimacy of political and scientific authorities who insisted on mask wearing (n=5), citing their own or others' clinical experiences of mask use. Some also asserted that masks were ineffective: they "do not protect from COVID," the virus "passes through" the mask, or masks "protect others but not yourself." Two users emphasized the ineffectiveness of masks by observing that other countries successfully managed the pandemic without imposing masks or by insisting that handwashing was more important. Another two users claimed that COVID-19 would have no effect on them, one arguing that it was not a severe disease and the other claiming that the virus "did not exist."

Debates over mask safety also emerged. Those considering masks to be safe (n=3) linked their claims to articles maintaining that masks did not reduce oxygen intake or asserted more generally that safety concerns had not surfaced over extensive mask use in the past. Mask opponents (n=4) countered claims of mask safety; some argued that they reduced oxygen intake—an especially serious concern during the summer, when the French state made them mandatory—whereas others contended that masks provoked skin conditions or were fabricated "in dirty places."

Certain Twitter users criticized obligatory mask wearing as a ploy for political or economic gain. For some (n=5), mask regulations were a means of economic profiteering; the French government could "sell masks" or "collect fines," retailers could become "rich" from mask proceeds, or "Chinese made the virus and sell the masks." For others (n=11), masks were one tool in the arsenal of interests by politicians who sought to impose a "sanitary dictatorship," to surveil citizens by embedding a "tracking device" in the masks or literally "muzzling" them.

Unequal access to masks also emerged as a subject of debate. Some attributed the mask shortage to poor governance (n=7) and others argued that masks should be distributed free of charge (n=11).

XSL·FO

RenderX

How social media users did or did not integrate masks into their daily lives was an additional debate. For some (n=8), masks were inconvenient in the workplace, on sunny days, and for makeup wearers. Others (n=4) found mask policies that recommended frequently changing or washing masks unrealistic; one offered tactics to evade mask use, and another claimed to wear masks to "avoid the fine." Some users (n=10) argued about appropriate contexts for mask wearing, suggesting that they be worn only in closed spaces and by at-risk individuals and deploring constant reminders to wear masks. These arguments were countered by those who sought to make the best of mask wearing (n=3), maintaining that masks were a small price to pay for safety or offering tips for more comfortable mask wearing.

Some pro-mask users criticized inappropriate mask wearing (n=35), which included "no-maskers," those wearing masks "under the nose" or "in the pocket," and those who removed masks altogether to talk or smoke. They publicly denounced "no-maskers" (n=13), communicating descriptions of offenders to named authorities, recounting scenes of "no-maskers" driven

out of public spaces, and deploring violence and resistance to calls for mask wearing.

Finally, state and regional changes in mask-wearing policies catalyzed many users to communicate new rules (n=44); some users interpreted these changes as governmental incoherence (n=8), whereas others (n=4) complained that traditional media outlets "talked too much" about masks (see Multimedia Appendix 2 for mask-related coding tree).

### In-depth Interviews

#### Interview Population

A total of 427 volunteers and workers were randomly invited from a contact database provided by the FRC to share their experiences anonymously and to provide recommendations for the organization. The response rate was 5.6% (24/427); 24 FRC workers and volunteers between the ages of 31 and 70 years participated in qualitative interviews (Table 1). Interviews did not collect data on the educational level of participants (see Multimedia Appendix 1 for the interview guide).

**Table 1.** Gender of interviewed French Red Cross workers and volunteers.

| Gender | Participants (N=24), n (%) | | |
|---|---|---|---|
| | Workers | Volunteers | Total |
| Women | 5 (21) | 10 (42) | 15 (63) |
| Men | 2 (8) | 7 (29) | 9 (38) |
| Total | 7 (29) | 17 (71) | 24 (100) |

### Social and Traditional Media Provoked Anxiety, Uncertainty, and Disdain

All interviewees agreed that COVID-19 was dangerous for themselves, for the elderly, and for those with comorbidities. One-half contended that the persistence, volume, and contradictions in social and traditional media coverage cultivated deep uncertainty and/or anxiety. A few participants disdained claims of certainty expressed by social media users and commentators. Although some informants valued continuous media coverage encouraging public adherence to preventive measures or, in the case of social media, promoting FRC visibility or alternatives to traditional media coverage, most informants were critical. Social media networks, including Twitter and Facebook, came under particular criticism, reproached by FRC participants as uncontrolled sources and disseminators of rumors that undermined effective public health strategies. Social media users, they claimed, uncritically accepted rumors and, as a result, perceived themselves as an authority. Another critic, a retired industrial engineer and now FRC volunteer, lamented the following:

> You know there is something that is hurting society badly nowadays, and that is social networks. It disintegrates society at the speed of lightning. So as soon as someone on a social network says, "The vaccine causes this, it causes that, there are risks, there are things," it is seen, 10, 20, 30 thousand times and there you have it: it ignites.

Responses varied when we encouraged informants to identify how social and traditional media influenced their work offline during the pandemic. Some participants noted that internet-borne misinformation did not pose real-life problems for their daily field activities. Indeed, the retired engineer could not identify specific ways that social media had influenced FRC work, but nonetheless insisted that it did. He acknowledged, "I cannot link it to that...I could not tell you that there is an influence of social media or not...[although] I personally believe it." Several other respondents, however, found that traditional and social media production fanned their disquiet. Media production was "anxiety-producing" ("anxiogène"), and several expressed disdain for the relentless criticism of experts, who seemed not to be concerned with the consequences of their remarks for listeners. One former hospital worker and volunteer observed the following:

> I have never seen as many professors [as I have on TV] than during COVID. Every evening, there was a new one, saying, "We should not have done it this way, we should do it that way," or "Not at all, we should have used x treatment," [or] "It was not this, it was not that." That was incredibly anxiety-producing for the average person...The epicenter was the unknown.

Other informants echoed this uncertainty and anxiety, bemoaning traditional media outlets' desires to attract viewers and readers and to elicit their fuller engagement with disseminated information. One worker described media outlets

as "horrendous," always in "search for a buzz." He opined, "Me, I find that horrendous...Because it puts people into a terrible state of mind, the term that summarizes everything: it's anxiety-producing."

Hence, several informants identified traditional and social media as exacerbating their uncertainty and anxiety, but also recognized that it remained difficult to pinpoint how social media and media productions influenced specific offline actions, either their own or those of other FRC workers.

### "Tuning Out" as a Coping Strategy

One strategy reported by participants was to "tune out" traditional and social media, by "shutting down" their televisions or "logging off of Facebook." For certain informants, turning off social and traditional media feeds resulted from an effort to reduce information-provoking anxiety and uncertainty. For others, turning off social media feeds stemmed from a desire to avoid excessive misleading information, or from a need to shut out unduly negative, repetitive claims. Hence, one volunteer claimed, "I cut off my Facebook and Twitter...I just stopped all that idiocy," whereas another said of traditional media, "they [only] recounted the deaths, the hospitalizations, etc. I'd had enough...I didn't want to listen anymore." In such cases, workers and volunteers relied on the expertise of knowledgeable family members and the FRC to provide information and advice.

Turning off media feeds was not foolproof. One FRC worker who avoided social media feeds inadvertently found in her Facebook feed a news report suggesting a link between children with COVID-19 and Kawasaki syndrome, just as schools were about to reopen in France. At that point, she decided not to seek further information about the subject; she sent her children back to school and returned to volunteering with the FRC.

### Traditional and Social Media: One of Multiple Influences on Decisions to Participate in Field Activities

The COVID-19 pandemic exacerbated needs among marginalized populations in the Paris region and imposed heightened demands on responding organizations, including the FRC, but multiple factors shaped individual workers' and volunteers' decisions to participate in FRC field activities that potentially brought them into contact with people suffering from the disease.

Social and traditional media coverage and commentary of ongoing FRC activities was one salient influence and, in some cases, galvanized our informants' decisions to participate in field activities. Informants discussed television and social media depicting FRC support to transfer COVID-19 patients from overwhelmed hospitals to less busy ones, as well as the new *Croix-Rouge chez vous* (Red Cross at your home) platform, a call center–based service to deliver food and medicine to at-risk people confined to their homes and to provide psychological counseling. Workers and volunteers signaled that this media coverage heightened their desire to participate in field activities. Our informants underscored the pandemic as a singular, urgent, "historical" moment and their pride in contributing to FRC interventions. One volunteer noted, "It's the first time for any volunteer who is alive today, at the Red Cross...at least in France we had never lived a crisis of such magnitude. We had all been

trained since day 1 at the FRC for this kind of catastrophe." A worker proclaimed,

> We have all been very proud of what the Red Cross has done during this crisis. I think that whatever the domain, we have seen our colleagues in trains [patient transfers via high-speed railway, depicted on national television], in social outreach.

Nonetheless, decisions to return to FRC volunteering and work resulted from multiple factors. The influence of social and traditional media conjugated with perceptions of personal and familial health risks from exposure to COVID-19, family obligations concerning childcare and gendered expectations, income, employer pressures, and FRC regulations.

Our informants reported that FRC colleagues and family members relentlessly called and shared information over the phone and through messaging apps, such as WhatsApp, to influence their decisions and to remind them of familial obligations. Several participants faced active discouragement to participate in field activities by family members worrying about their health or that of other family members. One volunteer reported WhatsApp discussions of familial pressures among her fellow volunteers:

> The most experienced [volunteers] asked themselves, "COVID, we don't know what it is. I have a wife, I have kids, do I go into the field or not?" The fathers started a discussion...I think there were a couple of people under pressure from their wives, who were saying, "You shouldn't go into the field." So they would come to the crisis center [to volunteer] and then return home...to be there for snacks, baths, etc. They would "do the Red Cross" behind the scenes, and afterwards, in the evening, they'd be [at home] in "daddy mode."

In addition to these familial obligations were gendered expectations. Several women volunteers and workers discontinued field activities because they had primary responsibility for childcare during school closures, although one man assumed childcare activities so that his wife could continue her field activities, even though this decision reduced family income.

Employer and related financial pressures were also important. A volunteer withdrew from FRC field activities at the request of her employer, who worried about workplace COVID-19 transmission. One worker received a request from the FRC not to work remotely, but instead to be paid for part-time unemployment because she had children to care for. Because the family could not cope financially with her partial unemployment, the worker and her husband found a family member to help care for the children.

Finally, certain volunteers did not decide at all: the FRC forbade volunteers aged 70 years and older, as well as volunteers and workers with existing comorbidities, to participate in any field-based activities. Although we heard reports of a few volunteers and workers circumventing these restrictions or participating through distant support, this measure made the decision for the volunteers and workers themselves.

XSL•FO

RenderX

### Informational and Institutional Influences on Daily Protective Practices

For workers and volunteers engaging in FRC activities during the pandemic, how to carry out their work safely remained an important, daily consideration. Widely circulating, contradictory information about the virus and nonpharmaceutical prevention measures was just one factor shaping their everyday practices on the job or in the field. Other factors included state regulations, FRC institutional measures, communications and ethos, as well as coworker relationships and practices.

Several study participants noted that social media rumors and misinformation circulated among some FRC workers and volunteers, but they were divided about their effects on offline activities. Some contended that misinformation about mask effectiveness and its associated offline behaviors (eg, consequent refusals to wear masks) could pose problems for implementing certain activities. They also perceived misinformation attributing the origin or voluntary spread of the virus to a specific group as dangerous, for such claims contradicted the International Federation of the Red Cross principles of "humanity," "unity," and "universality." Others considered this misinformation "ridiculous" and maintained that it would not affect volunteer or worker practice.

Multiple FRC institutional influences hindered the circulation of inaccurate information or diminished its effects on workers and volunteers and encouraged good protective practices to limit COVID-19 transmission (eg, mask use, physical distancing, and frequent handwashing). The FRC implemented an internal communication strategy, issuing bulletins and holding webinars to summarize information concerning protective equipment, other protective measures, and changing knowledge about SARS-CoV-2. It developed an internal "Frequently Asked Questions" page on its internal website to respond to worker and volunteer questions. Although FRC communications personnel recognized that workers and volunteers developed their own Facebook and WhatsApp groups, they built internal social networking tools and produced engaging and humorous social media content on the FRC action.

Several informants contended that the FRC institutional ethos, materialized in its uniform (Multimedia Appendix 3), compelled adherence to specific behaviors, including mask wearing. One male volunteer, when asked about how social media rumors about COVID-19 and protective measures might affect volunteers, reflected the following:

> When you wear a uniform, I think that you execute the instructions that were given. There is a chain of command, and it is there to make sure that things will be respected. A uniform, I think that the moment we put it on, we put our personal opinions aside to focus on the mission.

Coworkers, notably those who had fallen ill with COVID-19 and returned to work, also influenced adherence to protective practices. An older volunteer noted the following:

> It's simple: everyone at the local unit knows that I was contaminated. I arrive in the morning with my mask on, so they see me arriving and say, "right, we

> have to wear the mask" [laughs]...but it does not come naturally; if I am not there, they don't wear it...Because they did not have an experience of COVID, they don't feel that the mask is of any use. The problem is that for a while, they were told that masks were useless, then they were told that they are useful but only in public transport, and, finally, they were told that they are useful everywhere. You see, it's not easy for people to understand.

Yet not all our informants complied with FRC measures or coworker influences, reporting that normative measures could be negotiated in practice. Some, for instance, complained that FRC communication strategies contributed to the infodemic and exacerbated anxiety. One worker opined the following:

> I think...that...[FRC management] should communicate about things that really happen, and not those that may not happen. All of those things that generate a huge amount of anguish and are never going to happen...In our service, for certain people, that generates a lot of anxiety...Instead, they should give themselves time to see what happens...to reflect and to put into place in collaboration with the teams, not to put the cart before the horse.

Echoing online and traditional media controversies, certain participants complained that some protective measures—floor distancing marks, one-way corridors, and disinfection routines around everyday objects—were not scientifically justifiable or realistic. Mask wearing preoccupied many informants, nearly half of whom expressed confusion about changing official discourses around masks. Actual mask wearing, they reported, varied considerably. One worker negotiated her mask-wearing practice through her daily interactions with coworkers. She noted, "I share the office with a colleague...I ask, 'What do we do? Mask or no mask?' I ask as if we needed to ensure that everyone present would be on the same page." When asked if this was "a form of sanitary consent," she responded, "Yes, that's right."

## Discussion

### Overview

Anthropological investigations of social media narratives are relatively few [8,59,60]; however, they fit into a much longer tradition of situating such narratives into social, political, economic, and historical relationships and understandings of rumors, risks, and practice [50,61-63]. This mixed methods study sought to determine whether the COVID-19 infodemic, particularly online and media debates about viral risks and protective measures, affected FRC workers' and volunteers' decisions to return to work and their protective practices in the Paris region. It analyzed two distinct data sets: (1) social media (Twitter) and selected traditional media and (2) qualitative interviews with volunteers and workers in one of France's most important nongovernmental humanitarian assistance organizations, the FRC. Its contributions lie in an anthropological analysis of the influence of online debates around the virus, its risks, and protective masks on one group's risk perceptions, decisions, and daily practices. We show that

although online debates did affect FRC workers' and volunteers' emotions, decisions to return to work, and protective practices, other influences also played a role on their responses to the

COVID-19 pandemic. Figure 1 summarizes the interactions between online COVID-19 debates and offline responses among FRC workers and volunteers in our study.

**Figure 1.** Factors shaping French Red Cross (FRC) volunteer and worker activity decisions and workplace masking practices.



## Emotional Responses, Risk Perception, and Decision Making

We found that the COVID-19 infodemic incited anxiety, uncertainty, and, in some cases, disdain for expert opinions among interviewees. This response is coherent with other studies that found that social media can elicit emotional responses among users [28,29], including during the COVID-19 pandemic [35,38,39,43,64]. In some cases, FRC volunteers and workers protected themselves from the infodemic by shutting off social and traditional media or by relying on the FRC or on knowledgeable family members, colleagues, and friends. Although an adaptive response, tuning out and relying on social relationships can lead to less informed decision making and possibly worse health outcomes. New measures to reach such populations should be developed.

How the infodemic shaped FRC volunteers' and workers' risk perceptions and, ultimately, their decisions to return to field activities appears more complex. Although the information ecosystem provoked anxiety and uncertainty among our informants, other factors shaped their decisions: family, friends, colleagues and employers, financial concerns, and gendered expectations and norms influenced decisions to participate in FRC field activities. Risk perceptions and decisions to conduct field work did not simply entail epidemiological risks, but also social risks (ie, alienating family members encouraging one to return to the field or not), financial risks (ie, losing income), professional risks (ie, countering the wishes of one's employer), and even a risk of losing one's own sense of self-worth [63,65,66]. Our interviews suggested that the linkage between COVID-19 media debates and offline perceptions and decision

making is far from straightforward. These decisions were contingent on the personal, familial, sociopolitical, and economic relationships in which volunteers and workers were embedded. Our anthropological lens thus contributes to prior studies of online influences on offline perceptions and behavior by accounting for these multiple factors in shaping decisions [8,54].

## Daily Protective Practices

Online social and traditional media debates around protection from COVID-19 comprised one factor among several that affected the daily protective practices of FRC volunteers and workers. Mask wearing could be inconsistent, which was explained by our informants as the consequence of a volatile, dynamic informational environment that, at times, discounted the effectiveness of masks. State regulations and FRC messaging and enforcement, in particular, did much to reinforce mask wearing among workers and volunteers, although our field evidence also suggests that masking could be a socially situated practice; an individual's past history (eg, COVID-19 infection) and social relationships with coworkers or volunteers also played into mask-wearing practices.

## Anthropological Analyses of Infodemic Narratives

In contrast to many COVID-19 infodemic studies that conducted network or sentiment analyses [67-70], we employed an anthropological analysis of our evidence: we analyzed *all* claims, rather than triaging data as "true" or "false" and focusing solely on "false" information [19,21,71,72]. This approach has been useful in a pandemic, when biomedical and public health uncertainties about the virus have persisted and knowledge and policy have changed rapidly. Moreover, in evaluating epidemic

narratives regardless of their truthfulness, we followed a useful anthropological tradition of exploring local understandings of risk, misfortune, sorcery, and the occult in their political, economic, and sociocultural contexts [45,73-77]. Examining the claims and circulation of epidemic narratives as neither true nor false helped us to better understand online narratives as well as their offline influences.

Evaluating only "false" narratives about masks would have neglected the changing narratives of masks in traditional and online media, the confusion that such changes precipitated, and their consequent erosion of public trust in health authorities and political leaders. Early in the pandemic, many authorities, including those in France, claimed that surgical or cloth masks would not prevent COVID-19 transmission and that lay publics would not use masks correctly. Moreover, multiple questions about transmission remained unresolved, including those about fomites, sexual contact, and aerosolized transmission [78-81]. In July 2020, French regulations around masks were implemented. Online debates questioning the effectiveness and necessity of masks reflect this changing knowledge, echoing Eysenbach's observation that the early pandemic period must work with the "best evidence at the time," not immutable claims to truth [6]. Claims that masks would usher in a "sanitary dictatorship" and surveil populations are not simply "false." From an anthropological perspective, they yield rich insight into how certain French lay publics experienced current state public health measures for pandemic control as oppressive. Suspicions of state surveillance or ambitions of dictatorship did not spring forth suddenly in 2020 but have built on earlier political tensions. From 2017, the government has faced protests and massive strikes over retirement reforms, tax policy, police violence, budget cuts, and the climate crisis, among other concerns, as well as the rise of the *Gilets Jaunes* ("Yellow Vests") as a new political opposition group [24,82,83]. Further examination of social and traditional media debates through an anthropological lens will be useful to open up specific online debates, to situate them in longer-term political, economic, and social tensions.

## Principal Results

This study found that social and traditional media were just one of many influences on FRC workers' and volunteers' decisions to work in the field and on their daily protective practices. FRC informants reported that social and traditional media provoked anxiety, uncertainty, and disdain for commentators' claims to expertise. They also sought to "tune out" traditional and social media as a means of coping emotionally with persistent COVID-19 pandemic coverage.

## Limitations

This study has several limitations. First, French Twitter users and their beliefs, claims, or preoccupations are not representative of all populations in France or in the Paris region. For this reason, we undertook an iterative approach to media tracking, social listening, and interviews: selective media tracking influenced our social listening queries, and we used both as a proxy to identify major pandemic debates and a general typology or narratives about which we could ask our FRC informants during their interviews. Moreover, although captured tweets

numbered in the millions, our qualitative analyses of a random sample of those tweets could only evaluate a small proportion of their online narratives. We were, thus, unable to address all online narratives through this analysis. Nevertheless, 1400 tweets are a substantial number for qualitative analysis, fitting into a range of similar thematic analyses of tweets [84-86].

Second, only a small proportion of FRC actors use social media, including Twitter, making it difficult to track how participants engaged with online information. We coupled our Twitter analysis, however, with selective media tracking to ensure that we had a proxy for major debates. This weakness is simultaneously an advantage, in that our informants experienced the informational environment as a complex, multi-sourced, contradictory onslaught of information, and not through the framework of a single social media platform.

Third, we were unable to conduct numerous interviews, although 24 is generally acceptable for publishable qualitative research. Our interviews do not reflect the perspectives of all FRC workers and volunteers in the Paris region. We initially hypothesized that our informants, because of their participation in a nongovernmental organization assisting people in humanitarian emergencies, were less likely to experience the influence of the infodemic, or perhaps less likely to admit to this influence. Recent literature, however, shows that frontline workers are highly likely to suffer from the infodemic and that US nurses were uncertain about or opposed to receiving COVID-19 vaccines [64,87]. Our qualitative approach helped to mitigate these limitations. The small sample size and the flexibility of our qualitative interviews allowed us to pursue lengthy conversations, cultivate nonjudgmental interactions, and build trust with informants during the interviews.

Anthropological interviews can yield rich data concerning how informants perceive risks, describe decision-making processes, and explain their protective behaviors, and they can situate these narratives and practices in a broad context of social, political, and economic relationships. They cannot, however, shed light on what people do in practice. Two lockdowns and the FRC's heavy workload during the pandemic have hampered our efforts to undertake field observations of FRC activities. We supplemented our insights by meeting with FRC field actors and by the first author's observations of FRC outreach actions.

## Conclusions

This study found that the social and traditional media narratives about COVID-19 and protective practices had an important emotional influence on interviewed FRC workers and volunteers. Excessive, rapidly circulating, and misleading information produced by the social and traditional media was only one of several factors, however, that affected FRC workers' and volunteers' decisions to contribute to field activities and to pursue daily protective practices, namely masking. Additional investigation of online narratives and expanded qualitative investigation, including observations of their offline influences among larger population samples, will be crucial to develop further insights. Moreover, measures to address users who have disengaged from online sources of health information and who rely on social relationships to obtain information are necessary.

Tuning out can potentially lead to less informed decision   making, leading to worse health outcomes.

## Conflicts of Interest

None declared.

Multimedia Appendix 1
Interview guide.
[DOCX File , 20 KB - infodemiology_v1i1e27472_app1.docx ]

Multimedia Appendix 2
Mask codes and subcodes.
[XLSX File (Microsoft Excel File), 15 KB - infodemiology_v1i1e27472_app2.xlsx ]

Multimedia Appendix 3
French Red Cross uniform.
[PNG File , 2000 KB - infodemiology_v1i1e27472_app3.png ]

## References

1. Risk Communication and Community Engagement Readiness and Response to Coronavirus Disease (COVID-19): Interim Guidance. Geneva, Switzerland: World Health Organization; 2020 Mar 19. URL: https://www.who.int/publications/i/item/risk-communication-and-community-engagement-readiness-and-initial-response-for-novel-coronaviruses [accessed 2021-01-31]

2. Ratzan SC, Sommariva S, Rauh L. Enhancing global health communication during a crisis: Lessons from the COVID-19 pandemic. Public Health Res Pract 2020 Jun 30;30(2):1-6 [FREE Full text] [doi: 10.17061/phrp3022010] [Medline: 32601655]

3. Zarocostas J. How to fight an infodemic. Lancet 2020 Feb 29;395(10225):676 [FREE Full text] [doi: 10.1016/S0140-6736(20)30461-X] [Medline: 32113495]

4. Orso D, Federici N, Copetti R, Vetrugno L, Bove T. Infodemic and the spread of fake news in the COVID-19-era. Eur J Emerg Med 2020 Apr 23:327-328 [FREE Full text] [doi: 10.1097/MEJ.0000000000000713] [Medline: 32332201]

5. Okan O, Bollweg TM, Berens E, Hurrelmann K, Bauer U, Schaeffer D. Coronavirus-related health literacy: A cross-sectional study in adults during the COVID-19 infodemic in Germany. Int J Environ Res Public Health 2020 Jul 30;17(15):5503 [FREE Full text] [doi: 10.3390/ijerph17155503] [Medline: 32751484]

6. Eysenbach G. How to fight an infodemic: The four pillars of infodemic management. J Med Internet Res 2020 Jun 29;22(6):e21820 [FREE Full text] [doi: 10.2196/21820] [Medline: 32589589]

7. Eysenbach G. Infodemiology: The epidemiology of (mis)information. Am J Med 2002 Dec;113(9):763-765. [doi: 10.1016/S0002-9343(02)01473-0]

8. Stalcup M. The invention of infodemics: On the outbreak of Zika and rumors. Somatosphere. 2020 Mar 16. URL: http://somatosphere.net/2020/infodemics-zika.html/ [accessed 2021-09-28]

9. Safarnejad L, Xu Q, Ge Y, Bagavathi A, Krishnan S, Chen S. Identifying influential factors in the discussion dynamics of emerging health issues on social media: Computational study. JMIR Public Health Surveill 2020 Jul 28;6(3):e17175 [FREE Full text] [doi: 10.2196/17175] [Medline: 32348275]

10. Chew C, Eysenbach G. Pandemics in the age of Twitter: Content analysis of tweets during the 2009 H1N1 outbreak. PLoS One 2010;5(11):e14118 [FREE Full text] [doi: 10.1371/journal.pone.0014118] [Medline: 21124761]

11. Munich Security Conference. World Health Organization. 2020 Feb 15. URL: https://www.who.int/director-general/speeches/detail/munich-security-conference [accessed 2021-01-30]

12. Rovetta A, Bhagavathula AS. COVID-19-related web search behaviors and infodemic attitudes in Italy: Infodemiological study. JMIR Public Health Surveill 2020 May 05;6(2):e19374 [FREE Full text] [doi: 10.2196/19374] [Medline: 32338613]

13. 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]

XSL•FO
RenderX

14. Xue J, Chen J, Hu R, Chen C, Zheng C, Su Y, et al. Twitter discussions and emotions about the COVID-19 pandemic: Machine learning approach. J Med Internet Res 2020 Nov 25;22(11):e20550 [FREE Full text] [doi: 10.2196/20550] [Medline: 33119535]

15. Savoia E, Lin L, Gamhewage GM. A conceptual framework for the evaluation of emergency risk communications. Am J Public Health 2017 Sep;107(S2):S208-S214. [doi: 10.2105/AJPH.2017.304040] [Medline: 28892436]

16. Chong YY, Cheng HY, Chan HYL, Chien WT, Wong SYS. COVID-19 pandemic, infodemic and the role of eHealth literacy. Int J Nurs Stud 2020 Aug;108:103644 [FREE Full text] [doi: 10.1016/j.ijnurstu.2020.103644] [Medline: 32447127]

17. Pobiruchin M, Zowalla R, Wiesner M. Temporal and location variations, and link categories for the dissemination of COVID-19-related information on Twitter during the SARS-CoV-2 outbreak in Europe: Infoveillance study. J Med Internet Res 2020 Aug 28;22(8):e19629 [FREE Full text] [doi: 10.2196/19629] [Medline: 32790641]

18. Paakkari L, Okan O. COVID-19: Health literacy is an underestimated problem. Lancet Public Health 2020 May;5(5):e249-e250 [FREE Full text] [doi: 10.1016/S2468-2667(20)30086-4] [Medline: 32302535]

19. Naeem SB, Bhatti R, Khan A. An exploration of how fake news is taking over social media and putting public health at risk. Health Info Libr J 2020 Jul 12:143-149 [FREE Full text] [doi: 10.1111/hir.12320] [Medline: 32657000]

20. Correia RB, Wood IB, Bollen J, Rocha LM. Mining social media data for biomedical signals and health-related behavior. Annu Rev Biomed Data Sci 2020 Jul;3:433-458 [FREE Full text] [doi: 10.1146/annurev-biodatasci-030320-040844] [Medline: 32550337]

21. Vijjali R, Potluri P, Kumar S, Teki S. Two stage transformer model for COVID-19 fake news detection and fact checking. ArXiv Preprint posted online on November 26, 2020 [FREE Full text]

22. Medina Serrano JC, Papakyriakopoulos O, Hegelich S. NLP-based feature extraction for the detection of COVID-19 misinformation videos on YouTube. In: Proceedings of the 1st Workshop on NLP for COVID-19 at the 58th Annual Meeting of the Association for Computational Linguistics. 2020 Presented at: 1st Workshop on NLP for COVID-19 at the 58th Annual Meeting of the Association for Computational Linguistics; July 5-10, 2020; Virtual URL: https://aclanthology.org/2020.nlpcovid19-acl.17.pdf

23. Ahmed W, López Seguí F, Vidal-Alaball J, Katz MS. COVID-19 and the "film your hospital" conspiracy theory: Social network analysis of Twitter data. J Med Internet Res 2020 Oct 05;22(10):e22374 [FREE Full text] [doi: 10.2196/22374] [Medline: 32936771]

24. Smith R, Cubbon S, Wardle C. Under the surface: Covid-19 vaccine narratives, misinformation and data deficits on social media. First Draft. 2020 Nov 12. URL: https://firstdraftnews.org/long-form-article/under-the-surface-covid-19-vaccine-narratives-misinformation-and-data-deficits-on-social-media/ [accessed 2020-12-18]

25. Aiello AE, Renson A, Zivich PN. Social media- and internet-based disease surveillance for public health. Annu Rev Public Health 2020 Apr 02;41:101-118. [doi: 10.1146/annurev-publhealth-040119-094402] [Medline: 31905322]

26. Lazarus JV, Ratzan SC, Palayew A, Gostin LO, Larson HJ, Rabin K, et al. A global survey of potential acceptance of a COVID-19 vaccine. Nat Med 2021 Feb;27(2):225-228 [FREE Full text] [doi: 10.1038/s41591-020-1124-9] [Medline: 33082575]

27. Islam MS, Sarkar T, Khan SH, Mostofa Kamal A, Hasan SMM, Kabir A, et al. COVID-19-related infodemic and its impact on public health: A global social media analysis. Am J Trop Med Hyg 2020 Oct;103(4):1621-1629 [FREE Full text] [doi: 10.4269/ajtmh.20-0812] [Medline: 32783794]

28. Martínez-Ferrer B, Moreno D, Musitu G. Are adolescents engaged in the problematic use of social networking sites more involved in peer aggression and victimization? Front Psychol 2018;9:801 [FREE Full text] [doi: 10.3389/fpsyg.2018.00801] [Medline: 29896139]

29. Marino C, Gini G, Vieno A, Spada MM. The associations between problematic Facebook use, psychological distress and well-being among adolescents and young adults: A systematic review and meta-analysis. J Affect Disord 2018 Jan 15;226:274-281. [doi: 10.1016/j.jad.2017.10.007] [Medline: 29024900]

30. Doornwaard SM, ter Bogt TFM, Reitz E, van den Eijnden RJJM. Sex-related online behaviors, perceived peer norms and adolescents' experience with sexual behavior: Testing an integrative model. PLoS One 2015;10(6):1-18 [FREE Full text] [doi: 10.1371/journal.pone.0127787] [Medline: 26086606]

31. Moreno MA, Whitehill JM. Influence of social media on alcohol use in adolescents and young adults. Alcohol Res 2014;36(1):91-100 [FREE Full text] [Medline: 26259003]

32. Tran BX, Huong LT, Hinh ND, Nguyen LH, Le BN, Nong VM, et al. A study on the influence of internet addiction and online interpersonal influences on health-related quality of life in young Vietnamese. BMC Public Health 2017 Jan 31;17(1):138 [FREE Full text] [doi: 10.1186/s12889-016-3983-z] [Medline: 28143462]

33. Saran I, Fink G, McConnell M. How does anonymous online peer communication affect prevention behavior? Evidence from a laboratory experiment. PLoS One 2018;13(11):1-16 [FREE Full text] [doi: 10.1371/journal.pone.0207679] [Medline: 30462718]

34. Centola D. The spread of behavior in an online social network experiment. Science 2010 Sep 3;329(5996):1194-1197. [doi: 10.1126/science.1185231] [Medline: 20813952]

35.  Siebenhaar KU, Köther AK, Alpers GW. Dealing With the COVID-19 infodemic: Distress by information, information avoidance, and compliance with preventive measures. Front Psychol 2020;11:567905 [FREE Full text] [doi: 10.3389/fpsyg.2020.567905] [Medline: 33224060]

36.  Lep Ž, Babnik K, Hacin Beyazoglu K. Emotional responses and self-protective behavior within days of the COVID-19 outbreak: The promoting role of information credibility. Front Psychol 2020;11:1846 [FREE Full text] [doi: 10.3389/fpsyg.2020.01846] [Medline: 32849087]

37.  Lin Y, Hu Z, Alias H, Wong LP. Influence of mass and social media on psychobehavioral responses among medical students during the downward trend of COVID-19 in Fujian, China: Cross-sectional study. J Med Internet Res 2020 Jul 20;22(7):e19982 [FREE Full text] [doi: 10.2196/19982] [Medline: 32584779]

38.  Lee JJ, Kang K, Wang MP, Zhao SZ, Wong JYH, O'Connor S, et al. Associations between COVID-19 misinformation exposure and belief with COVID-19 knowledge and preventive behaviors: Cross-sectional online study. J Med Internet Res 2020 Nov 13;22(11):e22205 [FREE Full text] [doi: 10.2196/22205] [Medline: 33048825]

39.  Pahayahay A, Khalili-Mahani N. What media helps, what media hurts: A mixed methods survey study of coping with COVID-19 using the media repertoire framework and the appraisal theory of stress. J Med Internet Res 2020 Aug 06;22(8):e20186 [FREE Full text] [doi: 10.2196/20186] [Medline: 32701459]

40.  Low DM, Rumker L, Talkar T, Torous J, Cecchi G, Ghosh SS. Natural language processing reveals vulnerable mental health support groups and heightened health anxiety on reddit during COVID-19: Observational study. J Med Internet Res 2020 Oct 12;22(10):e22635 [FREE Full text] [doi: 10.2196/22635] [Medline: 32936777]

41.  Skoda E, Spura A, De Bock F, Schweda A, Dörrie N, Fink M, et al. Change in psychological burden during the COVID-19 pandemic in Germany: Fears, individual behavior, and the relevance of information and trust in governmental institutions [Article in German]. Bundesgesundheitsblatt Gesundheitsforschung Gesundheitsschutz 2021 Mar;64(3):322-333 [FREE Full text] [doi: 10.1007/s00103-021-03278-0] [Medline: 33481055]

42.  Bala R, Srivastava A, Ningthoujam GD, Potsangbam T, Oinam A, Anal CL. An observational study in Manipur State, India on preventive behavior influenced by social media during the COVID-19 pandemic mediated by cyberchondria and information overload. J Prev Med Public Health 2021 Jan;54(1):22-30 [FREE Full text] [doi: 10.3961/jpmph.20.465] [Medline: 33618496]

43.  Wang C, Pan R, Wan X, Tan Y, Xu L, McIntyre RS, et al. A longitudinal study on the mental health of general population during the COVID-19 epidemic in China. Brain Behav Immun 2020 Jul;87:40-48 [FREE Full text] [doi: 10.1016/j.bbi.2020.04.028] [Medline: 32298802]

44.  Wang C, Chudzicka-Czupała A, Grabowski D, Pan R, Adamus K, Wan X, et al. The association between physical and mental health and face mask use during the COVID-19 pandemic: A comparison of two countries with different views and practices. Front Psychiatry 2020;11:569981 [FREE Full text] [doi: 10.3389/fpsyt.2020.569981] [Medline: 33033485]

45.  White L. Speaking With Vampires: Rumor and History in Colonial Africa. Berkeley, CA: University of California Press; 2000.

46.  Tappan J. Blood work and "rumors" of blood: Nutritional research and insurrection in Buganda, 1935–1970. Int J Afr Hist Stud 2014;47(3):473-494.

47.  Graboyes M. The Experiment Must Continue: Medical Research and Ethics in East Africa, 1940-2014. Athens, OH: Ohio University Press; 2015.

48.  Fairhead J, Leach M, Small M. Where techno-science meets poverty: Medical research and the economy of blood in The Gambia, West Africa. Soc Sci Med 2006 Aug;63(4):1109-1120. [doi: 10.1016/j.socscimed.2006.02.018] [Medline: 16630676]

49.  Scheper-Hughes N. Commodity fetishism in organs trafficking. Body Soc 2001 Sep;7(2-3):31-62. [doi: 10.1177/1357034x0100700203]

50.  Geissler PW, Pool R. Editorial: Popular concerns about medical research projects in sub-Saharan Africa--A critical voice in debates about medical research ethics. Trop Med Int Health 2006 Jul;11(7):975-982 [FREE Full text] [doi: 10.1111/j.1365-3156.2006.01682.x] [Medline: 16827698]

51.  Kozinets R. Netnography: Redefined. 2nd edition. London, UK: SAGE Publications; Jul 24, 2015.

52.  Larson HJ. Stuck: How Vaccine Rumors Start─and Why They Don't Go Away. Oxford, UK: Oxford University Press; 2020.

53.  Durand JY, Cunha MI. 'To all the anti-vaxxers out there…': Ethnography of the public controversy about vaccination in the time of COVID-19. Soc Anthropol 2020 May 18:1-2 [FREE Full text] [doi: 10.1111/1469-8676.12805] [Medline: 32836939]

54.  Krieg LJ, Berning M, Hardon A. Anthropology with algorithms? An exploration of online drug knowledge using digital methods. Med Anthropol Theory 2017 Sep 28;4(3):21-52 [FREE Full text] [doi: 10.17157/mat.4.3.458]

55.  Pool R, Geissler W. Medical Anthropology. Maidenhead, UK: Open University Press; 2005.

56.  Carlson M. Embedded links, embedded meanings. Journal Stud 2016 Apr 20;17(7):915-924. [doi: 10.1080/1461670x.2016.1169210]

57.  Guadaloupe F. Audiences presse : "Le Figaro" en tête devant "Le Parisien" et "Le Monde", "20 Minutes" solide. Pure Médias. 2020 Jun 24. URL: https://www.ozap.com/actu/audiences-presse-le-figaro-en-tete-devant-le-parisien-et-le-monde-20-minutes-solide/594406 [accessed 2021-03-05]

58. Landis JR, Koch GG. The measurement of observer agreement for categorical data. Biometrics 1977 Mar;33(1):159-174. [Medline: 843571]

59. Fogarty S, Elmir R, Hay P, Schmied V. The experience of women with an eating disorder in the perinatal period: A meta-ethnographic study. BMC Pregnancy Childbirth 2018 May 02;18(1):121 [FREE Full text] [doi: 10.1186/s12884-018-1762-9] [Medline: 29720107]

60. Güzel H. Pain as performance: Re-virginisation in Turkey. Med Humanit 2018 Jun;44(2):89-95. [doi: 10.1136/medhum-2017-011414] [Medline: 29724778]

61. Bourdieu P. Esquisse d'une Théorie de la Pratique: Précédé de Trois Études d'Ethnologie Kabyle. Geneva, Switzerland: Librairie Droz; 1972.

62. Fine GA, Campion-Vencent V, Heath C, editors. Rumor Mills: The Social Impact of Rumor and Legend. Piscataway, NJ: Transaction Publishers; 2009.

63. Douglas M. Risk and Blame: Essays in Cultural Theory. London, UK: Routledge; 1992.

64. Dubey S, Biswas P, Ghosh R, Chatterjee S, Dubey MJ, Chatterjee S, et al. Psychosocial impact of COVID-19. Diabetes Metab Syndr 2020;14(5):779-788 [FREE Full text] [doi: 10.1016/j.dsx.2020.05.035] [Medline: 32526627]

65. Launiala A, Honkasalo M. Malaria, danger, and risk perceptions among the Yao in rural Malawi. Med Anthropol Q 2010 Sep;24(3):399-420. [doi: 10.1111/j.1548-1387.2010.01111.x] [Medline: 20949843]

66. Agbo IE, Johnson RC, Sopoh GE, Nichter M. The gendered impact of Buruli ulcer on the household production of health and social support networks: Why decentralization favors women. PLoS Negl Trop Dis 2019 Apr;13(4):1-20 [FREE Full text] [doi: 10.1371/journal.pntd.0007317] [Medline: 30986205]

67. Gruzd A, Mai P. Going viral: How a single tweet spawned a COVID-19 conspiracy theory on Twitter. Big Data Soc 2020 Jul 20;7(2):1-9 [FREE Full text] [doi: 10.1177/2053951720938405]

68. Bruns A, Harrington S, Hurcombe E. 'Corona? 5G? or both?': The dynamics of COVID-19/5G conspiracy theories on Facebook. Media Int Aust 2020 Aug 04;177(1):12-29. [doi: 10.1177/1329878x20946113]

69. Hung M, Lauren E, Hon ES, Birmingham WC, Xu J, Su S, et al. Social network analysis of COVID-19 sentiments: Application of artificial intelligence. J Med Internet Res 2020 Aug 18;22(8):e22590 [FREE Full text] [doi: 10.2196/22590] [Medline: 32750001]

70. Alamoodi A, Zaidan B, Zaidan A, Albahri O, Mohammed K, Malik R, et al. Sentiment analysis and its applications in fighting COVID-19 and infectious diseases: A systematic review. Expert Syst Appl 2021 Apr 01;167:114155 [FREE Full text] [doi: 10.1016/j.eswa.2020.114155] [Medline: 33139966]

71. Donovan J. Social-media companies must flatten the curve of misinformation. Nature 2020 Apr 14. [doi: 10.1038/d41586-020-01107-z] [Medline: 32291410]

72. Sharma K, Seo S, Meng C, Rambhatla S, Liu Y. COVID-19 on social media: Analyzing misinformation in Twitter conversations. ArXiv Preprint posted online on October 22, 2020. [FREE Full text]

73. Venturini T. Diving in magma: How to explore controversies with actor-network theory. Public Underst Sci 2009 May 29;19(3):258-273. [doi: 10.1177/0963662509102694]

74. Evans-Pritchard EE. Witchcraft, Oracles and Magic among the Azande. Oxford, UK: Clarendon Press; 1937.

75. Favret-Saada J. Deadly Words: Witchcraft in the Bocage. New York, NY: Cambridge University Press; 1980.

76. Wenzel Geissler P. 'Kachinja are coming!': Encounters around medical research work in a Kenyan village. Africa 2005 May;75(2):173-202. [doi: 10.3366/afr.2005.75.2.173]

77. Peeters Grietens K, Ribera J, Erhart A, Hoibak S, Ravinetto R, Gryseels C, et al. Doctors and vampires in sub-Saharan Africa: Ethical challenges in clinical trial research. Am J Trop Med Hyg 2014 Aug;91(2):213-215 [FREE Full text] [doi: 10.4269/ajtmh.13-0630] [Medline: 24821846]

78. Eikenberry S, Mancuso M, Iboi E, Phan T, Eikenberry K, Kuang Y, et al. To mask or not to mask: Modeling the potential for face mask use by the general public to curtail the COVID-19 pandemic. Infect Dis Model 2020;5:293-308 [FREE Full text] [doi: 10.1016/j.idm.2020.04.001] [Medline: 32355904]

79. Leung NHL, Chu DKW, Shiu EYC, Chan K, McDevitt JJ, Hau BJP, et al. Respiratory virus shedding in exhaled breath and efficacy of face masks. Nat Med 2020 May;26(5):676-680 [FREE Full text] [doi: 10.1038/s41591-020-0843-2] [Medline: 32371934]

80. Makison Booth C, Clayton M, Crook B, Gawn J. Effectiveness of surgical masks against influenza bioaerosols. J Hosp Infect 2013 May;84(1):22-26. [doi: 10.1016/j.jhin.2013.02.007] [Medline: 23498357]

81. Asadi S, Cappa CD, Barreda S, Wexler AS, Bouvier NM, Ristenpart WD. Efficacy of masks and face coverings in controlling outward aerosol particle emission from expiratory activities. Sci Rep 2020 Sep 24;10(1):15665 [FREE Full text] [doi: 10.1038/s41598-020-72798-7] [Medline: 32973285]

82. Taylor P. Yellow Jackets hit Macron's blind spot. Politico. 2018 Nov 28. URL: https://www.politico.eu/article/yellow-jackets-protests-emmanuel-macron-blind-spot/ [accessed 2020-12-18]

83. Breeden AJ, Breeden A. Macron Faces first big street protests, a challenge to his labor overhaul. The New York Times. 2017 Sep 12. URL: https://www.nytimes.com/2017/09/12/world/europe/france-labor-law-protests.html [accessed 2020-12-18]

84.  Berry N, Lobban F, Belousov M, Emsley R, Nenadic G, Bucci S. #WhyWeTweetMH: Understanding why people use Twitter to discuss mental health problems. J Med Internet Res 2017 Apr 05;19(4):e107 [FREE Full text] [doi: 10.2196/jmir.6173] [Medline: 28381392]

85.  Golder S, Bach M, O'Connor K, Gross R, Hennessy S, Gonzalez Hernandez G. Public perspectives on anti-diabetic drugs: Exploratory analysis of Twitter posts. JMIR Diabetes 2021 Jan 26;6(1):e24681 [FREE Full text] [doi: 10.2196/24681] [Medline: 33496671]

86.  Slavik CE, Buttle C, Sturrock SL, Darlington JC, Yiannakoulias N. Examining Tweet content and engagement of Canadian public health agencies and decision makers during COVID-19: Mixed methods analysis. J Med Internet Res 2021 Mar 11;23(3):e24883 [FREE Full text] [doi: 10.2196/24883] [Medline: 33651705]

87.  McClendon S, Proctor K. New survey of 13k US nurses: Findings indicate urgent need to educate nurses about COVID-19 vaccines. American Nurses Association. 2020 Oct 29. URL: https://www.nursingworld.org/news/news-releases/2020/new-survey-of-13k-u.s.-nurses-findings-indicate-urgent-need-to-educate-nurses-about-covid-19-vaccines/ [accessed 2021-01-17]

## Abbreviations

**FRC:** French Red Cross

XSL•FO
**RenderX**

<u>Original Paper</u>

# Health Care Providers' Trusted Sources for Information About COVID-19 Vaccines: Mixed Methods Study

Eden Brauer[1,2], PhD, RN; Kristen Choi[1,2], PhD, RN; John Chang[3], MPH; Yi Luo[3], PhD; Bruno Lewin[4], MD, DTMH; Corrine Munoz-Plaza[3], MPH; David Bronstein[4], MD; Katia Bruxvoort[3,4,5], MPH, PhD

[1]School of Nursing, University of California, Los Angeles, Los Angeles, CA, United States

[2]Department of Health Policy and Management, Fielding School of Public Health, University of California, Los Angeles, Los Angeles, CA, United States

[3]Department of Research & Evaluation, Kaiser Permanente Southern California, Pasadena, CA, United States

[4]Southern California Permanente Medical Group, Kaiser Permanente Southern California, Pasadena, CA, United States

[5]Department of Epidemiology, School of Public Health, University of Alabama at Birmingham, Birmingham, AL, United States

**Corresponding Author:**
Kristen Choi, PhD, RN
School of Nursing
University of California, Los Angeles
700 Tiverton Ave
Los Angeles, CA
United States
Phone: 1 3107947493
Email: krchoi@ucla.edu

## *Abstract*

**Background:**   Information and opinions shared by health care providers can affect patient vaccination decisions, but little is known about who health care providers themselves trust for information in the context of new COVID-19 vaccines.

**Objective:**   The purpose of this study is to investigate which sources of information about COVID-19 vaccines are trusted by health care providers and how they communicate this information to patients.

**Methods:**   This mixed methods study involved a one-time, web-based survey of health care providers and qualitative interviews with a subset of survey respondents. Health care providers (physicians, advanced practice providers, pharmacists, nurses) were recruited from an integrated health system in Southern California using voluntary response sampling, with follow-up interviews with providers who either accepted or declined a COVID-19 vaccine. The outcome was the type of information sources that respondents reported trusting for information about COVID-19 vaccines. Bivariate tests were used to compare trusted information sources by provider type; thematic analysis was used to explore perspectives about vaccine information and communicating with patients about vaccines.

**Results:**   The survey was completed by 2948 providers, of whom 91% (n=2683) responded that they had received ≥1 dose of a COVID-19 vaccine. The most frequently trusted source of COVID-19 vaccine information was government agencies (n=2513, 84.2%); the least frequently trusted source was social media (n=691, 9.5%). More physicians trusted government agencies (n=1226, 93%) than nurses (n=927, 78%) or pharmacists (n=203, 78%; *P*<.001), and more physicians trusted their employer (n=1115, 84%) than advanced practice providers (n=95, 67%) and nurses (n=759, 64%; *P*=.002). Qualitative themes (n=32 participants) about trusted sources of COVID-19 vaccine information were identified: processing new COVID-19 information in a health care work context likened to a "war zone" during the pandemic and communicating information to patients. Some providers were hesitant to recommend vaccines to pregnant people and groups they perceived to be at low risk for COVID-19.

**Conclusions:**   Physicians have stronger trust in government sources and their employers for information about COVID-19 vaccines compared with nurses, pharmacists, and advanced practice providers. Strategies such as role modeling, tailored messaging, or talking points with standard language may help providers to communicate accurate COVID-19 vaccine information to patients, and these strategies may also be used with providers with lower levels of trust in reputable information sources.

XSL•FO
**RenderX**

## Introduction

### Background

The rapid onset of the COVID-19 pandemic has created a secondary "infodemic" of health information challenges globally [1,2]. Health information about COVID-19 has proliferated in news media and social media (ie, web-based applications for creating or sharing content and social networking), and has rapidly evolved as scientists and public health professionals learned new information about the transmission and management of SARS-CoV-2 [3,4]. The real-time availability of new scientific and health information on COVID-19 has undoubtedly aided pandemic response but has also created information challenges for health care providers and the public in navigating misinformation, contradictions, and complexity [5]. Understanding how to effectively navigate a complex health information environment is an essential component of pandemic response for health care providers, who must apply changing information about the COVID-19 pandemic to practice.

### Prior Work

Despite growing reliance on the internet as a source of health information, many individuals still rely upon health care providers to learn new health information [6,7]. There is strong evidence that physicians, nurses, and other health care providers are among the most trusted entities for health information [8,9]. Although having up-to-date pandemic knowledge is essential for health care providers to educate the public, in the COVID-19 pandemic, health care providers are challenged to keep pace with ever-growing health information on SARS-CoV-2 and COVID-19 [5]. Although there is much literature on health care professionals as a trusted entity for health information among the public, including information about COVID-19 vaccines [8], less is known about who health care professionals themselves trust for health information.

How health care providers learn new COVID-19 information and convey that information to patients is especially important in regard to COVID-19 vaccines. Evidence suggests that health care provider opinions about vaccines and vaccine recommendations can affect patient decisions about vaccines [10-12]. Though nearly 70% of the US adult public has received at least one dose of a COVID-19 vaccine as of July 2021 [13] and more than 80% of health care providers have received a COVID-19 vaccine [14,15], vaccination levels vary substantially by locale, and there are still sizeable populations of adults that are unvaccinated. Health care providers have the potential to address barriers to COVID-19 vaccination and increase vaccine confidence as the US vaccination strategy shifts from mass vaccination to more traditional clinic-based administration of vaccines [16-18].

### Study Purpose

Given the high level of public trust in health care professionals for health information, the health information literacy of providers is essential for appropriate patient education and communication about COVID-19 vaccines. However, to date, there have been few studies about the specific sources that health care providers rely on to find trusted health information and how these sources affect their discussions about COVID-19 vaccines with patients. The purpose of this mixed methods study was to investigate which sources of information about COVID-19 vaccines are trusted by health care providers and how they communicate COVID-19 vaccine information with patients.

## Methods

### Design

This study used an explanatory-sequential mixed methods design with data from a web-based survey followed by qualitative interviews [19]. The study took place from March to May 2021 at Kaiser Permanente Southern California (KPSC), an integrated health system with approximately 15 hospitals, 235 clinics, and over 20,000 clinical employees. A one-time survey was sent to KPSC health care providers to assess COVID-19 experiences, COVID-19 perceptions including trusted sources of information, and demographics characteristics. We also conducted semistructured interviews using Rapid Assessment Procedures (RAP) for qualitative research [20,21] to further investigate health care provider perspectives on trusted sources of COVID-19 vaccine information. The study was approved by the KPSC Institutional Review Board, and all participants gave informed consent.

### Survey Procedures

KPSC health care providers were eligible to participate in the survey if they were actively practicing in the KPSC health system at the time of the survey and had access to a web-enabled device to complete the survey (phone, tablet, computer). We engaged leadership in medicine, nursing, and pharmacy to email the survey opportunity to their staff. Two reminder emails were sent from clinical leadership, and they were also provided with study flyers to post at hospitals and clinics.

### Survey Measures

#### Outcome Measures

The primary outcome was a survey item asking providers to select which sources of information they trusted for learning about COVID-19 vaccines among the following: government entities (local, state, or federal), their health system employer, mainstream news (television, print, radio), social media, personal contacts, physicians, and other, where participants could specify other sources with free text. The categories were not mutually exclusive, allowing respondents to select multiple sources.

#### Exposure Measures

The primary exposure was self-identified provider type (physician, advanced practice provider [Physician Assistant or Advanced Practice Registered Nurse], nurse [Registered Nurse or Licensed Vocational Nurse], pharmacist, and other). We also examined demographic and health history characteristics of the

sample, including gender, age, race/ethnicity, and history of testing positive for COVID-19.

## Rapid Qualitative Assessment Procedures

Qualitative data were collected to further elucidate perspectives on COVID-19 vaccines among health care workers. As part of the survey, respondents had the option to indicate their willingness to be contacted for a follow-up interview. From those participants who volunteered to be contacted for an interview, we stratified potential interviewees by provider type (physician, pharmacist, nurse) and whether they had received a COVID-19 vaccine (yes/no). We then contacted 10 participants in each of these six groups (physician-acceptor, physician-decliner, pharmacist-acceptor, pharmacist-decliner, nurse-acceptor, nurse-decliner, 60 potential participants total) to ensure that interviews reflected a range of experiences and perspectives regarding COVID-19 vaccine confidence and hesitancy among providers. Interview participants were offered a small gift as an incentive for their time.

Interviews were conducted by authors KC and JC, who are experienced researchers with a background in conducting qualitative research and using semistructured interview guides. KC has a background in nursing and health services research, and as such, she conducted all interviews with nurses. JC has a background in public health and health services research and conducted all interviews with physicians and pharmacists.

Semistructured interviews with providers who either accepted or declined the COVID-19 vaccine were conducted using RAP [20,21]. An interview guide was developed with open-ended questions and probes about providers' experiences with information about COVID-19; the vaccines; and how they receive, gather, and appraise various information sources. Perspectives on educational resources or other interventions that could be used to support vaccine confidence were also explored. Interviews were conducted by a member of the research team with experience in qualitative research, took place by telephone, and lasted approximately 15 to 30 minutes each. Interview data were digitally recorded and then transcribed for analysis and triangulation with survey data.

## Analysis

For quantitative survey data, we used chi-square tests to compare health care providers (physicians, physician assistants and nurse practitioners, pharmacists, nurses, others) by which sources of information they had indicated that they trusted. Analyses were conducted using R version 4.0.3 (R Foundation for Statistical Computing). We systematically analyzed the qualitative data using inductive thematic analysis [22-24]. A member of the research team reviewed the interview transcripts for data familiarization and generated codes with attached segments of data that were relevant to the research question. These codes were reviewed by study investigators, collapsed or broadened to ensure good fit with the data, and organized into themes and subthemes. To enhance credibility, the technique of member-checking was used, where stakeholder representatives from each provider group (nurses, physicians, pharmacists) reviewed and provided feedback about preliminary analyses. Themes were further refined to capture the most salient patterns in the data and then triangulated with quantitative data to gain deeper insight about providers' experiences with COVID-19 and vaccine information.

## *Results*

### Sample Description

A total of 3164 potential participants opened the survey, 3052 verified eligibility and consented to the survey, and 2948 went on to complete the survey. The sample comprised 45.0% (n=1326) physicians, 40.2% (n=1184) nurses, 8.8% (n=259) pharmacists, and 5.7% (n=169) advanced practice providers (Table 1). The majority of respondents were female (n=2051, 69.6%) and White (n=1087, 36.9%) or Asian (n=1153, 39.1%). About 8% (n=240) of respondents reported a history of testing positive for COVID-19, and 1.9% (n=55) of the sample reported being currently pregnant. Among the total sample, 91.3% (n=2683) reported receiving at least one dose of a COVID-19 vaccine.

**Table 1.** Sample description.

| Variable | Participants (N=2948), n (%) |
| --- | --- |
| **Age (years)[a]** | |
| 18-30 | 65 (2.2) |
| 31-40 | 811 (27.5) |
| 41-50 | 1027 (34.8) |
| 51-60 | 697 (23.7) |
| 61-70 | 313 (10.6) |
| >70 | 34 (1.2) |
| **Gender** | |
| Female | 2051 (69.6) |
| Male | 891 (30.2) |
| Other | 6 (0.2) |
| **Provider type** | |
| Physician | 1326 (45.0) |
| Advanced practice provider | 169 (5.7) |
| Pharmacist | 259 (8.8) |
| Nurse | 1184 (40.2) |
| **Race/ethnicity** | |
| White | 1087 (36.9) |
| African American/Black | 98 (3.3) |
| Hispanic/Latinx | 340 (11.5) |
| Asian | 1153 (39.1) |
| Native American/Alaskan/Hawaiian | 18 (6.1) |
| Multiple | 167 (5.7) |
| Other | 85 (2.9) |
| **Ever had COVID-19** | |
| Yes | 240 (8.1) |
| No | 2536 (86) |
| Unsure | 172 (5.8) |
| Received at least one dose of a COVID-19 vaccine[b] | 2683 (91.3) |
| **Plans to recommend COVID-19 vaccines to patients** | |
| Will recommend | 2203 (74.9) |
| Will recommend if asked | 593 (20.2) |
| Unsure | 99 (3.4) |
| Will not recommend | 47 (1.6) |

[a]One participant did not provide information on age.

[b]Nine participants skipped this question.

## Survey Results: Comparison of Trusted Information Sources by Provider Type

The most trusted source of COVID-19 vaccine information across all health care provider types in our sample was government agencies (n=2513, 84.2% of the sample), followed by KPSC (n=2191, 74.3%). The least frequently trusted source of COVID-19 information by health care providers in our sample across all provider types was social media (n=691, 9.5%). When comparing information sources by provider type, there were significant differences for three information sources: government agencies, employer, and news media. More physicians trusted government agencies (n=1226, 93%) than nurses (n=927, 78%) or pharmacists (n=203, 78%; $P<.001$). For trust in one's

employer, there were differences for physicians compared with nurses and advanced practice providers. Although advanced practice providers trusted their employer at a frequency of 67% (n=95) and nurses at 64% (n=759), 84% (n=1115) of physicians reported trusting their employer for information about COVID-19 (P=.002). Overall trust in news as a source of information was lower for all provider groups (P=.003), but physicians (n=66, 27%) and pharmacists (n=351, 25%) more frequently reported trusting news media than advanced practice providers (n=29, 21%) or nurses (n=231, 20%). When compared to other provider groups, nurses generally reported lower levels of trust in nearly all information sources.

## Interview Results

A total of 32 interviews were conducted across all provider/vaccine groups (15 nurses, 8 pharmacists, 9 physicians). Of these, 17 interviewees indicated that they had declined the vaccine (10 nurses, 4 pharmacists, 3 physicians). In this analysis on experiences with information about COVID-19 and COVID-19 vaccines, we report on two overarching themes among provider vaccine acceptors and decliners: processing information in a health care work context likened to a "war zone" during the pandemic and communicating information to patients.

### Theme: Processing Information in a Health Care Work Context Likened to a "War Zone"

The first theme reflects provider accounts of navigating the constant influx of new information during the COVID-19 pandemic while also managing fluctuating work demands and protecting their own health and safety in a workplace, described by several participants as a "war zone." As one nurse-decliner stated:

> *It was absolutely horrible. Patients were dying every day. Lots and lots of death that I witnessed there, lots of strain on staff. Physically, mentally, it was hard.*

The war zone work environment was characterized by unpredictability, with one nurse-decliner recalling, "it didn't really seem like anyone knew what was going on," while another nurse-decliner described work as "different every day." Several providers recalled being unexpectedly "deployed" to COVID-19 units and having to adapt to rapidly changing information, patient volume and acuity, and work responsibilities. One nurse-decliner explained, "there was no warning, this was pandemic world." Providers also observed the impact of these conditions on quality of care. One nurse-decliner stated, "there was no choice. [...] We couldn't provide the same level of care."

### Subtheme: Valuing Transparency

Participants described how they evaluated COVID-19 and vaccine information in these circumstances. Many providers emphasized the need for transparency and "more balanced information," particularly in the context of government and corporation-led dissemination, with one nurse-acceptor stating:

> *It would be helpful if [...] people knew that it wasn't just these two [pharmaceutical] companies or the government that was supporting it.*

Another nurse-decliner shared:

> *It's very one-sided, the information that's being given out. People have a false sense of security thinking they're vaccinated because they don't think they can still get COVID. It fits the narrative.*

### Subtheme: Acknowledging Ambiguity

Providers also questioned the oversimplification of COVID-19 information and vaccination decisions, with one nurse-decliner explaining, "It didn't answer our doubts." Part of this questioning stemmed from their firsthand experiences with unfolding information early in the pandemic. A nurse-decliner remembered:

> *Trying to preserve PPE [personal protective equipment], when we weren't really sure how [the virus] was transmitted.*

Many participants felt that oversimplification and lack of transparency contributed to feelings of hesitancy, distrust, or questioning. Instead, there was a preference for open acknowledgment of the complexities and limitations of available information, and respect for multiple points of view. As one pharmacist-acceptor pointed out, little attention was paid to "figuring out what those issues are [related to hesitancy] and addressing those issues."

In making sense of COVID-19 information, participants also described the need to recognize the biases in their professional experiences. One nurse-decliner shared:

> *As healthcare workers sometimes our perspectives can be skewed, toward really bad. We're not going to see people that have mild cases. I have to remind myself, "This isn't what everyone is going through that has COVID."*

### Subtheme: Appraising Various Sources of COVID-19 Information

Participants shared their perceptions of various sources of information related to COVID-19 vaccines. As displayed in Figure 1, nearly all providers identified major governmental entities such as the Centers for Disease Control and Prevention (CDC) as their primary source of trusted COVID-19 vaccine information, and this was generally consistent in interviews as well. However, a small subgroup of providers—most often, nurses—expressed misgivings about government sources during the interviews. As one nurse-decliner noted of information from the CDC:

> *It changes all the time so it's really scary. It's a lot of changes. It's kind of hard to rely on data when data is practically new all the time.*

Another nurse-decliner perceived changes in information and discrepancies with other organizations as reasons to distrust the CDC, stating:

> *I'm really not trusting what the CDC is saying, just because they have just been going back and forth...They're contradicting what the World Health Organization is saying. I really question the FDA [Food and Drug Administration], I question the CDC.*

**Figure 1.** This figure shows the frequency of which sources of information health care providers (N=2948) reported trusting in a survey conducted from March to May 2021 in Southern California. *Group differences were significant at the .005 level in a chi-square test. APRN: advanced practice nurse; LVN: licensed vocational nurse; PA: physician assistant; RN: registered nurse.



One nurse-decliner appreciated the visual information provided at the county level by local government officials, explaining, "I'm very visual, I need to see the graphs, I need to see the trends."

Participants also reported strong levels of trust in the information provided by their health system employer. Some described using updates from their employer as a reference in their clinical practice, but others noted the challenge of keeping up with the constant barrage of information from management. As one pharmacist-acceptor stated, "After a while, you keep posting things on the wall and it just ends up being wallpaper."

Trust in mainstream news as an information source was low across all provider groups, with one nurse-decliner sharing, "I stopped watching the news." Although overall trust in social media was comparable to mainstream news, some providers emphasized the credibility of personal testimonies, or what one nurse-decliner called "real life experiences, real life realities," shared on such platforms. Another nurse-decliner used social media as an information starting point, explaining:

> I definitely get most of my news from social media, from Instagram. Then I go research it for myself to make sure it's true.

The public social media accounts of frontline physicians were also mentioned as trustworthy information sources.

In addition to the major sources of information listed, many participants described how they relied on their own personal experiences with COVID-19 as an information source in their perspectives about the vaccine. For some, the experience of personally becoming ill or caring for ill family members provided information about COVID-19 that was distinct from other conventional information sources, and these experiences influenced their perspectives on vaccine decision-making. One nurse-decliner said:

> We actually ended up having COVID. I still can't taste very well, I still can't smell very well, I'm not 100% back to where my energy level was and that's part of why I'm still hesitant to get the vaccine.

Another nurse-decliner shared:

> I actually had COVID a few weeks ago, and my views on the vaccine have truly changed. It's been rough. [...] I still feel congested, still have a mild cough, don't have 100% energy.

In sharing about loved ones who had COVID-19, a physician-acceptor stated, "most everyone survived thankfully but I do have friends who still have symptoms." Finally, some providers explained that no single source of information was sufficient in the context of rapidly evolving information, with one physician-acceptor stating, "for me, the key is to have multiple sources of information."

### Theme: Communicating Information About COVID-19 and Vaccines to Patients

The second overarching theme relates to how health care providers communicate COVID-19 and COVID-19 vaccine information to patients. Many participants described the impact of a changed work environment, specifically a shift to telehealth, which resulted in "limited face-to-face encounters with patients" and required new approaches to sharing information.

**Subtheme: "Being a Role Model Matters"**

Many participants who had received the COVID-19 vaccine believed it was their professional responsibility to serve as an example to patients. They described the impact of disclosing their own personal vaccine experiences in conveying information to patients. One nurse-acceptor explained:

> When you see me, and I'm like, "Hey, I'm three months out from my second dose and I'm doing fine," I'm a witness that it's OK.

Similarly, a physician-acceptor stated:

> I think it does help to say as a physician that I've been vaccinated and that it was fine for me and that I believe in it. Being a role model matters.

**Subtheme: "Tailoring the Message"**

Many participants recognized a need to "tailor the message" in communication with patients to reflect individual preferences and values. In some instances, this meant framing the risks of vaccines in the context of benefits; for example, focusing the discussion on serious risks such as hospitalization, death, and

other long-term consequences. Tailoring the message also involved consideration of incentives that might resonate with an individual patient. In some cases, participants discussed the vaccine as a path back toward normalcy, with one physician-acceptor stating, "it can allow us to get back to normal life again, and that's exciting." For others, the health and safety of others—loved ones or the broader community—was used to invoke collective responsibility and an opportunity to help, particularly with patients who did not perceive COVID-19 as a serious threat to their own health. Tailoring the message was also important in preserving patients' sense of autonomy in vaccine decision-making. One physician-acceptor explained, "these are the options, these are the pros and cons, take your pick."

In tailoring vaccine messaging, providers discussed prioritizing some patients over others considering available vaccine safety data and perceived patient risk. One pharmacist-acceptor reflected on how they would discuss the vaccine with young women and stated:

> I wouldn't encourage them as much, especially to females who are of childbearing age, because I don't want to recommend something that prevents them from having a child.

Another pharmacist-acceptor said, "for the young healthy crowd, I wouldn't push it as much as the older group."

### Subtheme: Recognizing Social, Political, and Historical Factors

Recognizing broader contextual factors in COVID-19 vaccine communication was an important consideration for participants. Several providers emphasized the technical nature of COVID-19 information and challenges in communication with people who lacked foundational science knowledge. For example, a nurse-acceptor said:

> I don't think there's enough information out there to explain to the medical staff, EVS [environmental services], housekeepers, people who aren't knowledgeable in the science aspect. Even nurses, some nurses, they don't understand what mRNA does.

Another nurse-acceptor reflected:

> You have these highly educated physicians, but then you have people who aren't as educated who don't have as much resources to get the education. It should be fair and equal.

Concerns about political and historical reasons for not getting vaccinated were also raised by participants. One physician-acceptor shared:

> Nurses that I've spoken to and tried to encourage vaccination, what I'm aware of is, there's actually a history in the Philippines where Sanofi rolled out a kind of mandatory dengue vaccine, and I think the government profited off of it but many children died. And, so there's a lot of pressure, people tell me, from their family or others that are still living in the Philippines not to be vaccinated.

There was recognition that groups with specific social, political, and historical meaning around vaccination would benefit from tailored communication approaches. Ultimately, many participants found it "very difficult to convince someone to do it if they truly do not believe in it," in the words of a pharmacist-acceptor.

## Discussion

### Principal Results

In this mixed methods study of COVID-19 vaccine experiences and perceptions, we examined the information sources that health care providers use and trust, and how they have navigated the COVID-19 *infodemic* [1]. Providers generally trusted government health sources—specifically the CDC, noted in qualitative interviews—and the health system where they practiced. They had less trust in news media and social media. Though these patterns were consistent across provider types, we found small differences in trust by provider type. Nurses, pharmacists, and advanced practice providers had less trust in information from government sources, their employer, and the news compared with physicians. Qualitative interviews suggested that this mistrust stemmed from frequently changing and at times conflicting information about COVID-19 from the government, challenging and even traumatic pandemic working conditions, and perceiving COVID-19 vaccine information to be "one-sided" such that it did not fully resolve providers' questions and doubts. These experiences and perceptions may reflect differences in pandemic working conditions by provider type, leading to differences between physicians and other providers. For example, some physician specialties were able to provide care via telehealth during the pandemic, while nurses had a direct patient-facing role and may have found changing or conflicting information difficult to integrate with a traumatic or stressful pandemic clinical context. Health care providers have been significantly challenged by keeping abreast of the latest understanding and guidance on COVID-19 clinical practice in the midst of misinformation, a high volume of new scientific information, and errors in or misunderstanding of the latest science [25,26]. Providers have faced the difficult task of integrating evolving, incomplete information into their practice while also needing to take immediate action for their patients and manage potential implications of information changes for their own personal health and safety [27].

Providers who had received a COVID-19 vaccine shared strategies for how they communicated information about the vaccines to patients but also recognized that convincing patients who did not believe in the vaccine was challenging. These strategies included role modeling the benefits and safety of the vaccine by disclosing their vaccination status as providers, tailoring the messaging to patient concerns, and recognizing structural forces that might contribute to vaccine hesitancy in specific demographic subgroups. Health care providers described the challenge of making sense of and sharing technical data with diverse groups of patients while avoiding oversimplification and confronting misinformation about vaccines in the public. Other strategies for vaccine messaging have been proposed based on principles of social, communication, and behavioral

science, such as prosocial appeals, framing recommendations positively, and making strong or presumptive recommendations for vaccination [28]. Our findings suggest that health care providers are weathering the challenge of providing patients with accurate information about COVID-19 vaccines but also that additional support for clinicians may be needed from public health entities and health systems so that they are fully prepared with messaging and educational tools. This may include standard messaging strategies and patient educational tools that providers can tailor. Additionally, there may be a need for interventions to reinforce health care provider trust in reputable information sources to ensure that providers are prepared to give accurate, quality information to patients.

## Limitations

This study had strengths and limitations that should be considered in interpreting its findings. The study used mixed methods, which allowed us to explore health care provider perspectives on COVID-19 vaccine information in greater depth than a survey alone would allow. The sample was large and diverse, representing multiple provider types and race/ethnicities. Study limitations were the cross-sectional and self-report nature of the survey. The study used voluntary response sampling, which did not allow for determination of an exact response rate or number of potential participants reached, which may have oversampled providers with favorable views about COVID-19 vaccines. However, levels of vaccination reported by providers in our study are consistent with other similar surveys and national averages, suggesting that the sample was reasonably representative [14,15,29]. The sources of trusted information assessed in our survey were not necessarily exhaustive of all sources providers may rely upon, although we provided an option for providers to write in an *other* response for sources not included in the survey response list. Finally, qualitative results were intended to explore quantitative findings in greater depth within our sample and thus are not necessarily generalizable. With qualitative research, there is risk for interviewer bias in data collection and coding/analysis, but we attempted to mitigate these risks by using a consistent interview protocol and using a team approach coming to consensus about codes and themes.

## Conclusion

Scientific evidence on the prevention and treatment of COVID-19 has been changing rapidly since the onset of the pandemic in early 2020. Early in the pandemic, the World Health Organization (WHO) passed a resolution on the COVID-19 response that included the importance of managing the *infodemic* in controlling the COVID-19 pandemic [30]. The WHO called for the provision of reliable content and science-based data to the public, measures to counter misinformation, and prevention of information activities that undermined public health response. As the uncertainty of the pandemic and the politicization of vaccines continue, there is a need to, first, ensure that all health care providers receive accurate information from reputable information sources that they can trust and, second, to ensure that health care providers have informational tools available to give quality information and recommendations to patients about vaccines.

## References

1. Zarocostas J. How to fight an infodemic. Lancet 2020 Feb 29;395(10225):676 [FREE Full text] [doi: 10.1016/S0140-6736(20)30461-X] [Medline: 32113495]
2. Calleja N, AbdAllah A, Abad N, Ahmed N, Albarracin D, Altieri E, et al. A public health research agenda for managing infodemics: methods and results of the First WHO Infodemiology Conference. JMIR Infodemiology 2021;1(1):e30979 [FREE Full text] [doi: 10.2196/30979] [Medline: 34604708]
3. Fauci AS, Lane HC, Redfield RR. Covid-19 - navigating the uncharted. N Engl J Med 2020 Mar 26;382(13):1268-1269 [FREE Full text] [doi: 10.1056/NEJMe2002387] [Medline: 32109011]
4. Song P, Karako T. COVID-19: real-time dissemination of scientific information to fight a public health emergency of international concern. Biosci Trends 2020 Mar 16;14(1):1-2. [doi: 10.5582/bst.2020.01056] [Medline: 32092748]
5. Mokhtari H, Mirzaei A. The tsunami of misinformation on COVID-19 challenged the health information literacy of the general public and the readability of educational material: a commentary. Public Health 2020 Oct;187:109-110 [FREE Full text] [doi: 10.1016/j.puhe.2020.08.011] [Medline: 32942170]
6. Smith D. Health care consumer's use and trust of health information sources. J Commun Healthcare 2013 Jul 18;4(3):200-210. [doi: 10.1179/1753807611y.0000000010]

XSL•FO
RenderX

7.  Hou J, Shim M. The role of provider-patient communication and trust in online sources in Internet use for health-related activities. J Health Commun 2010;15 Suppl 3:186-199. [doi: 10.1080/10810730.2010.522691] [Medline: 21154093]

8.  Rozek LS, Jones P, Menon A, Hicken A, Apsley S, King EJ. Understanding vaccine hesitancy in the context of COVID-19: the role of trust and confidence in a seventeen-country survey. Int J Public Health 2021;66:636255 [FREE Full text] [doi: 10.3389/ijph.2021.636255] [Medline: 34744589]

9.  Hong H, Oh HJ. The effects of patient-centered communication: exploring the mediating role of trust in healthcare providers. Health Commun 2020 Apr;35(4):502-511. [doi: 10.1080/10410236.2019.1570427] [Medline: 30706741]

10. Paterson P, Meurice F, Stanberry LR, Glismann S, Rosenthal SL, Larson HJ. Vaccine hesitancy and healthcare providers. Vaccine 2016 Dec 20;34(52):6700-6706 [FREE Full text] [doi: 10.1016/j.vaccine.2016.10.042] [Medline: 27810314]

11. Smith PJ, Kennedy AM, Wooten K, Gust DA, Pickering LK. Association between health care providers' influence on parents who have concerns about vaccine safety and vaccination coverage. Pediatrics 2006 Nov;118(5):e1287-e1292. [doi: 10.1542/peds.2006-0923] [Medline: 17079529]

12. Ylitalo KR, Lee H, Mehta NK. Health care provider recommendation, human papillomavirus vaccination, and race/ethnicity in the US National Immunization Survey. Am J Public Health 2013 Jan;103(1):164-169. [doi: 10.2105/AJPH.2011.300600] [Medline: 22698055]

13. COVID data tracker. Centers for Disease Control and Prevention. URL: https://covid.cdc.gov/covid-data-tracker/#datatracker-home [accessed 2021-08-31]

14. Moniz MH, Townsel C, Wagner AL, Zikmund-Fisher BJ, Hawley S, Jiang C, et al. COVID-19 vaccine acceptance among healthcare workers in a United States medical center. medRxiv (forthcoming) Preprint posted online on April 30, 2021. [doi: 10.1101/2021.04.29.21256186]

15. COVID vaccine facts for nurses survey. COVID Vaccine Facts for Nurses. 2021. URL: https://covidvaccinefacts4nurses.org/covid-19-survey [accessed 2021-10-10]

16. Verger P, Dubé E. Restoring confidence in vaccines in the COVID-19 era. Expert Rev Vaccines 2020 Nov;19(11):991-993. [doi: 10.1080/14760584.2020.1825945] [Medline: 32940574]

17. Goralnick E, Kaufmann C, Gawande AA. Mass-vaccination sites - an essential innovation to curb the Covid-19 pandemic. N Engl J Med 2021 May 06;384(18):e67. [doi: 10.1056/NEJMp2102535] [Medline: 33691058]

18. Gianfredi V, Pennisi F, Lume A, Ricciardi GE, Minerva M, Riccò M, et al. Challenges and opportunities of mass vaccination centers in COVID-19 times: a rapid review of literature. Vaccines (Basel) 2021 Jun 01;9(6):574 [FREE Full text] [doi: 10.3390/vaccines9060574] [Medline: 34205891]

19. Creswell J, Clark V. Designing and Conducting Mixed Methods Research. 3rd edition. Thousand Oaks, CA: Sage Publications; 2011:209-256.

20. Palinkas LA, Zatzick D. Rapid assessment procedure informed clinical ethnography (RAPICE) in pragmatic clinical trials of mental health services implementation: methods and applied case study. Adm Policy Ment Health 2019 Mar;46(2):255-270 [FREE Full text] [doi: 10.1007/s10488-018-0909-3] [Medline: 30488143]

21. Scrimshaw S, Hurtado E. Rapid assessment procedures for nutrition and primary health care: anthropological approaches to improving programme effectiveness. Latin American Center Reference Series. URL: https://agris.fao.org/agris-search/search.do?recordID=XF19900021532 [accessed 2021-10-10]

22. Carmichael T, Cunningham N. Theoretical data collection and data analysis with gerunds in a constructivist grounded theory study. Electronic J Business Res Methods 2017;15(2):73.

23. Clarke V, Braun V. Thematic analysis. In: Teo T, editor. Encyclopedia of Critical Psychology. New York, NY: Springer; 2012:57-71.

24. Braun V, Clarke V. Using thematic analysis in psychology. Qualitative Res Psychol 2006 Jan;3(2):77-101. [doi: 10.1191/1478088706qp063oa]

25. Moradi S, Abdi S. Pandemic publication: correction and erratum in COVID-19 publications. Scientometrics 2020 Nov 21:1-9 [FREE Full text] [doi: 10.1007/s11192-020-03787-w] [Medline: 33250543]

26. Zdravkovic M, Berger-Estilita J, Zdravkovic B, Berger D. Scientific quality of COVID-19 and SARS CoV-2 publications in the highest impact medical journals during the early phase of the pandemic: a case control study. PLoS One 2020;15(11):e0241826 [FREE Full text] [doi: 10.1371/journal.pone.0241826] [Medline: 33152034]

27. Delgado D, Wyss Quintana F, Perez G, Sosa Liprandi A, Ponte-Negretti C, Mendoza I, et al. Personal safety during the COVID-19 pandemic: realities and perspectives of healthcare workers in Latin America. Int J Environ Res Public Health 2020 Apr 18;17(8):2798 [FREE Full text] [doi: 10.3390/ijerph17082798] [Medline: 32325718]

28. Finney Rutten LJ, Zhu X, Leppin AL, Ridgeway JL, Swift MD, Griffin JM, et al. Evidence-based strategies for clinical organizations to address COVID-19 vaccine hesitancy. Mayo Clin Proc 2021 Mar;96(3):699-707 [FREE Full text] [doi: 10.1016/j.mayocp.2020.12.024] [Medline: 33673921]

29. Grumbach K, Judson T, Desai M, Jain V, Lindan C, Doernberg SB, et al. Association of race/ethnicity with likeliness of COVID-19 vaccine uptake among health workers and the general population in the San Francisco Bay Area. JAMA Intern Med 2021 Jul 01;181(7):1008-1011 [FREE Full text] [doi: 10.1001/jamainternmed.2021.1445] [Medline: 33783471]

30. COVID-19 response. World Health Organization. URL: https://www.who.int/health-cluster/news-and-events/news/COVID19/en/ [accessed 2021-10-09]

XSL•FO
RenderX

## Abbreviations

**CDC:** Centers for Disease Control and Prevention
**CIRT:** Care Improvement Research Team
**EVS:** environmental services
**FDA:** Food and Drug Administration
**KPSC:** Kaiser Permanente Southern California
**PPE:** personal protective equipment
**RAP:** Rapid Assessment Procedures
**WHO:** World Health Organization

XSL·FO
**RenderX**

<u>Original Paper</u>

# Impact of the World Inflammatory Bowel Disease Day and Crohn's and Colitis Awareness Week on Population Interest Between 2016 and 2020: Google Trends Analysis

Krixie Silangcruz[1], MBA, MD; Yoshito Nishimura[1,2], MD, MPH, PhD; Torrey Czech[1], MD; Nobuhiko Kimura[1], MD; Hideharu Hagiya[2], MD, PhD; Toshihiro Koyama[2], PhD; Fumio Otsuka[2], MD, PhD

[1]University of Hawaii, Honolulu, HI, United States

[2]Okayama University, Okayama, Japan

**Corresponding Author:**
Yoshito Nishimura, MD, MPH, PhD
University of Hawaii
1356 Lusitana St
Honolulu, HI
United States
Phone: 1 808 586 2910
Email: nishimura-yoshito@okayama-u.ac.jp

## Abstract

**Background:** More than 6 million people are affected by inflammatory bowel disease (IBD) globally. The World IBD Day (WID, May 19) and Crohn's and Colitis Awareness Week (CCAW, December 1-7) occur yearly as national health observances to raise public awareness of IBD, but their effects are unclear.

**Objective:** The aim of this study was to analyze the relationship between WID or CCAW and the public health awareness on IBD represented by the Google search engine query data.

**Methods:** This study evaluates the impact of WID and CCAW on the public awareness of IBD in the United States and worldwide from 2016 to 2020 by using the relative search volume of "IBD," "ulcerative colitis," and "Crohn's disease" in Google Trends. To identify significant time points of trend changes (joinpoints), we performed joinpoint regression analysis.

**Results:** No joinpoints were noted around the time of WID or CCAW during the study period in the search results of the United States. Worldwide, joinpoints were noted around WID in 2020 with the search for "IBD" and around CCAW in 2017 and 2019 with the search for "ulcerative colitis." However, the extents of trend changes were modest without statistically significant increases.

**Conclusions:** These results posed a question that WID and CCAW might not have worked as expected to raise public awareness of IBD. Additional studies are needed to precisely estimate the impact of health observances to raise the awareness of IBD.

**KEYWORDS**

inflammatory bowel disease; ulcerative colitis; Crohn disease; google trends; trend analysis; online health information; awareness; chronic disease; gastrointestinal; trend; impact; public health; United States

## Introduction

Inflammatory bowel disease (IBD) is a global disease with an increasing prevalence in newly industrialized countries, and rising cases have been documented in every continent [1]. Recent systematic reviews have demonstrated that IBD is increasing in such countries [2]. Globally in 2017, there were 6.8 million cases of IBD with an increased age-standardized prevalence rate from 1990 to 2017 [2]. Within the United States,

it was estimated that more than a million adult Americans had IBD [3,4].

Research into IBD, however, is largely underrepresented despite its prevalence owing to the multifactorial nature of the disease [5]. Global efforts have been made to raise the awareness of IBD. In 2010, the World IBD Day (WID) was created by the European Federation of Crohn's and Ulcerative Colitis Association and patient organizations to increase IBD awareness and to provide education about IBD to the public [6,7].

XSL•FO
RenderX

Similarly, Crohn's and Colitis Awareness Week (CCAW) was created by a US Senate resolution in 2011, with goals of encouraging all people in the United States to engage in activities aimed at raising awareness of IBD among the general public [8].

Disease awareness and health promotion campaigns are created to increase public health education, and awareness, and ultimately change behavior [9]. Approximately 200 health awareness days, weeks, or months are on the US National Health Observances calendar [10], and nearly 70% of these health awareness occasions have been introduced after 2005. Despite the increasing number of awareness initiatives, there is a lack of data regarding evidence of their effectiveness and impact [11]. This lack of data highlights the need for greater evaluation and quantifiable metrics to determine the impact of health behaviors on a global scale.

Because web-based searches are a predominant source of access to health awareness–related information, internet searches are a reflection of engagement between the public and resources, which increase disease awareness. Searches are individual proxies for public disease awareness and provide insight into the effect of dissemination of information via global public health days and weeks. Google Trends (GT) is a novel, open-source, freely accessible resource that allows researchers to analyze Google search query data [12]. An analysis regarding the efficacy and public health behaviors that resulted from IBD awareness initiatives has not been done previously. We aimed to perform a hypothesis generation if the WID or CCAW effectively increased the public health awareness for IBD through GT data by using joinpoint trend analysis.

## Methods

### Data Source

GT is a data source generated from the total Google search data [13]. These data are available to the public, and GT has been used in multiple social, public health, or global health research to dig into the public attention [14-25]. Surrogate of the public attention in GT is the relative popularity of specific search terms or topics in a certain category (eg, health), place, and time range. The relative popularity is defined as a relative search volume (RSV) with a scale of 0-100 (0 being the lowest popularity) [14,19-21]. The RSV correlates with how popular the terms are at a certain time point.

### Search Input

We followed protocols noted by previous studies [17,19,21]. Briefly, we accessed data between July 11 and 13, 2021, and chose [Inflammatory bowel disease], [Ulcerative colitis], and

[Crohn's disease] as search inputs. The location of the search included United States and worldwide.

### Search Variables

To specifically obtain the popularity of the disease-related search inputs, all searches were done with a "disease" option in the Health category (with a "disease" option, search volumes of subtopics or relevant themes are included). We chose each full year from 2016 to 2020 as search scales to visualize weekly trends of the RSVs (each year contains 52 or 53 weeks; the WID occurred in the 20th week of 2016-2019 and in the 19th week of 2020; CCAW occurred in the 48th-49th week in 2016-2019 and in the 49th-50th week in 2020).

### Statistical Analyses

We used a joinpoint regression model with the Joinpoint Regression Program version 4.9.0.0, March 2021 [26] to analyze the RSV data and their time trend. This software enables us to identify time points called joinpoints, where a temporal trend significantly changes. We defined the analysis criteria to look for up to 3 joinpoints. The weekly percentage changes between trend change points were determined with 95% CIs. The threshold for statistical significance was defined as a *P* value <.05, suggesting the level at which the slope differed from zero.

### Ethical Considerations

The publicly available data published by GT are utilized in the project [13]. This study was approved by the institutional review board of the Okayama University Hospital with a waiver for informed consent since the study intended to retrospectively analyze open data (1910-009). All research methods were performed in accordance with relevant guidelines and regulations.

## Results

### Trends in the Search Volume of Inflammatory Bowel Disease

Table 1 and Figure 1 describe the trends and trend changes of the weekly RSVs for "inflammatory bowel disease" in each full year from 2016 to 2020. With respect to the search results in the United States, no joinpoints were observed throughout the period. Regarding the search results worldwide, there was a joinpoint at the 45th week in 2019 before which a significant increase in the weekly percentage change of 0.2% (95% CI 0.1-0.4) was observed. In 2020, a joinpoint was noted in the 17th week (3 weeks before WID), after which there was a significant weekly increase in the RSV by 0.3% (*P*<.001). Further, the third joinpoint was observed in the 48th week (a week prior to CCAW). However, no joinpoints were noted from 2015 to 2018 around the time of WID or CCAW.

**Table 1.** Trend changes in the relative search volumes of inflammatory bowel disease in 2016-2020.[a]

| Country, year | Period 1 | | Period 2 | | Period 3 | | Period 4 | |
|---|---|---|---|---|---|---|---|---|
| | Weeks | Weekly percentage change (%) (95% CI) | Weeks | Weekly percentage change (%) (95% CI) | Weeks | Weekly percentage change (%) (95% CI) | Weeks | Weekly percentage change (%) (95% CI) |
| United States, 2016 | 1-52 | –0.1 (–0.4 to 0.2) | N/A[b] | N/A | N/A | N/A | N/A | N/A |
| United States, 2017 | 1-53 | –0.2 (–0.5 to 0.1) | N/A | N/A | N/A | N/A | N/A | N/A |
| United States, 2018 | 1-52 | –0.2 (–0.5 to 0.2) | N/A | N/A | N/A | N/A | N/A | N/A |
| United States, 2019 | 1-52 | –0.3[c] (–0.6 to 0) | N/A | N/A | N/A | N/A | N/A | N/A |
| United States, 2020 | 1-52 | 0.1 (–0.2 to 0.4) | N/A | N/A | N/A | N/A | N/A | N/A |
| Worldwide, 2016 | 1-52 | 0 (–0.2 to 0.2) | N/A | N/A | N/A | N/A | N/A | N/A |
| Worldwide, 2017 | 1-53 | 0 (–0.2 to 0.1) | N/A | N/A | N/A | N/A | N/A | N/A |
| Worldwide, 2018 | 1-52 | 0 (–0.2 to 0.2) | N/A | N/A | N/A | N/A | N/A | N/A |
| Worldwide, 2019 | 1-45 | 0.2[c] (0.1 to 0.4) | 45-52 | –2.6 (–5.4 to 0.3) | N/A | N/A | N/A | N/A |
| Worldwide, 2020 | 1-13 | –2.5[c] (–3.6 to –1.4) | 13-17 | 4.6 (–4.9 to 15.0) | 17-48 | 0.3[c] (0 to 0.6) | 48-52 | –5.4 (–10.9 to 0.5) |

[a]Periods were separated as Period 1-4, when the trend changes were statistically detected in the joinpoint regression analysis during the study period.

[b]N/A: not applicable.

[c]Significantly different from zero (*P*<.05).

**Figure 1.** Trends in the relative search volume of inflammatory bowel disease during 2016-2020. Weekly relative search volumes for the search term "inflammatory bowel disease" are described. World Inflammatory Bowel Disease day occurred in the 20th week of 2016-2019 and the 19th week of 2020; Crohn's & Colitis Awareness Week occurred in the 48th to 49th week in 2016-2019 and the 49th to 50th week in 2020. The number of slopes is determined by the number of joinpoints identified by the analysis. Joinpoints are the time points when statistically significant changes in the linear slopes are noted. RSV: relative search volume.



## Trends in the Search Volume of Ulcerative Colitis

Table 2 and Figure 2 describe the trends and trend changes of the weekly RSVs for ulcerative colitis in the designated period. In the search results of the United States and worldwide, a big surge was observed in the 3rd week in 2016. In 2020, a joinpoint was noted in the 16th week (4 weeks before WID), after which a nonstatistically significant but considerable weekly RSV increase by 3.7% (*P*<.001) was observed until the 24th week. For worldwide results, there was a prominent joinpoint in the 49th week (CCAW) in 2017. No other joinpoints were observed around the time of WID or CCAW in 2016 or 2018 to 2020.

**Table 2.** Trend changes in the relative search volumes of ulcerative colitis in 2016-2020.[a]

| Country, year | Period 1 | | Period 2 | | Period 3 | |
|---|---|---|---|---|---|---|
| | Weeks | Weekly percentage change (%) (95% CI) | Weeks | Weekly percentage change (%) (95% CI) | Weeks | Weekly percentage change (%) (95% CI) |
| United States, 2016 | 1-3 | 80.2[b] (40.1 to 131.8) | 3-6 | –19.5 (–37.4 to 3.6) | 6-52 | –0.2[b] (–0.4 to 0) |
| United States, 2017 | 1-53 | 0.2 (0 to 0.4) | N/A[c] | N/A | N/A | N/A |
| United States, 2018 | 1-52 | 0.1 (–0.1 to 0.3) | N/A | N/A | N/A | N/A |
| United States, 2019 | 1-52 | –0.3[b] (–0.4 to –0.1) | N/A | N/A | N/A | N/A |
| United States, 2020 | 1-16 | –2.8[b] (–4.0 to –1.6) | 16-24 | 3.7 (–0.3 to 7.8) | 24-52 | –0.5 (–0.9 to 0) |
| Worldwide, 2016 | 1-52 | –0.3[b] (–0.6 to –0.1) | N/A | N/A | N/A | N/A |
| Worldwide, 2017 | 1-46 | 0.1 (0 to 0.2) | 46-49 | 5.7 (–6.5 to 19.6) | 49-53 | –7.9[b] (–11.5 to –4.3) |
| Worldwide, 2018 | 1-52 | 0 (–0.1 to 0.1) | N/A | N/A | N/A | N/A |
| Worldwide, 2019 | 1-46 | 0 (–0.1 to 0.1) | 46-52 | –2.1 (–4.5 to 0.3) | N/A | N/A |
| Worldwide, 2020 | 1-52 | 0 (–0.4 to 0.3) | N/A | N/A | N/A | N/A |

[a]Periods were separated as Period 1-4, when the trend changes were statistically detected in the joinpoint regression analysis during the study period.

[b]Significantly different from zero ($P<.05$).

[c]N/A: not applicable.

**Figure 2.** Trends in the relative search volume of ulcerative colitis during 2016-2020. Weekly relative search volumes for the search term "ulcerative colitis." Except for the 16th week (4 weeks before World Inflammatory Bowel Disease Day) in the United States and the 49th week (Crohn's and Colitis Awareness Week) worldwide in 2020, no other joinpoints were noted around the time of World Inflammatory Bowel Disease Day or Crohn's and Colitis Awareness Week during the designated period. RSV: relative search volume.



## Trends in the Search Volume of Crohn disease

Table 3 and Figure 3 describe the trends and trend changes in the weekly RSVs for Crohn disease in the designated period. Between 2017 and 2019, there was no remarkable trend change in both the United States and worldwide. In 2020, joinpoints were observed in the 8th week, the 16th week, and the 24th week in the United States. For worldwide, joinpoints were observed in the 10th week, the 14th week, and the 24th week. However, there were no joinpoints around the time of WID or CCAW throughout the period.

**Table 3.** Trend changes in the relative search volumes of Crohn disease during 2016-2020.[a]

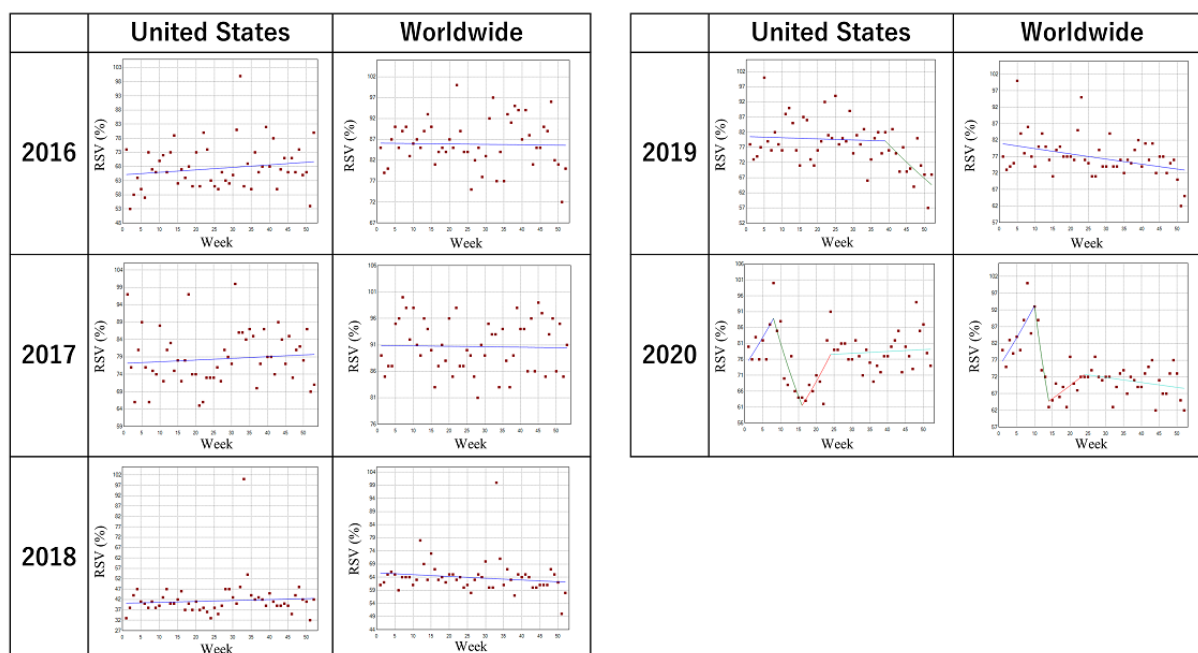| Country, year | Period 1 | | Period 2 | | Period 3 | | Period 4 | |
|---|---|---|---|---|---|---|---|---|
| | Weeks | Weekly percentage change (%) (95% CI) | Weeks | Weekly percentage change (%) (95% CI) | Weeks | Weekly percentage change (%) (95% CI) | Weeks | Weekly percentage change (%) (95% CI) |
| United States, 2016 | 1-52 | 0.1 (–0.1 to 0.3) | N/A[b] | N/A | N/A | N/A | N/A | N/A |
| United States, 2017 | 1-53 | 0.1 (–0.1 to 0.2) | N/A | N/A | N/A | N/A | N/A | N/A |
| United States, 2018 | 1-52 | 0.1 (–0.2 to 0.4) | N/A | N/A | N/A | N/A | N/A | N/A |
| United States, 2019 | 1-39 | 0 (–0.3 to 0.2) | 39-52 | –1.5[c] (–2.8 to –0.3) | N/A | N/A | N/A | N/A |
| United States, 2020 | 1-8 | 2.3 (–0.4 to 5.1) | 8-16 | –4.5[c] (–7.0 to –1.9) | 16-24 | 2.9[c] (0.2 to 5.7) | 24-52 | 0.1 (–0.3 to 0.4) |
| Worldwide, 2016 | 1-52 | 0 (–0.1 to 0.1) | N/A | N/A | N/A | N/A | N/A | N/A |
| Worldwide, 2017 | 1-53 | 0 (–0.1 to 0.1) | N/A | N/A | N/A | N/A | N/A | N/A |
| Worldwide, 2018 | 1-52 | –0.1 (–0.3 to 0.1) | N/A | N/A | N/A | N/A | N/A | N/A |
| Worldwide, 2019 | 1-52 | –0.2 (–0.3 to –0.1) | N/A | N/A | N/A | N/A | N/A | N/A |
| Worldwide, 2020 | 1-10 | 2.2[c] (0.6 to 3.8) | 10-14 | –8.7[c] (–16.4 to –0.4) | 14-24 | 1.2 (–0.4 to 2.8) | 24-52 | –0.2 (–0.5 to 0.1) |

[a]Periods were separated as Period 1-4, when the trend changes were statistically detected in the joinpoint regression analysis during the study period.

[b]N/A: not applicable.

[c]Significantly different from zero ($P<.05$).

**Figure 3.** Trends in the relative search volume of Crohn disease during 2016-2020. Weekly relative search volumes for the search term "Crohn's disease." No joinpoints were noted around the time of World Inflammatory Bowel Disease Day or Crohn's and Colitis Awareness Week throughout the period. RSV: relative search volume.



## Discussion

This study evaluated how the global campaigns for promoting IBD, such as WID and CCAW, affected public awareness by using the RSVs of GT data as a surrogate. Although there were several significant joinpoints for IBD, Crohn disease, and ulcerative colitis, overall, the results in this study posed a hypothesis that WID and CCAW might not have affected the public interest in the United States and worldwide. Rather, the RSVs seem to have been affected by timely topics. For example, in the 3rd week of 2016, when there was a significant increase in the RSV of ulcerative colitis in the United States, a famous US singer-songwriter reportedly passed away due to the disease. In March 2020, when significant trend changes were observed in the United States and worldwide, a well-known American comedian revealed that he had Crohn disease. Similarly, when the same comedian was featured in a film about a man who has Crohn disease in June 2020, there were trend changes both in the United States and worldwide (24th week). Only in 2017 worldwide, considerable trend changes in the RSVs for

ulcerative colitis were noted around CCAW, although the weekly percentage change was not statistically significant. Given the rapid increase in the global prevalence of IBD with increasing health care costs, raising public awareness of IBD is a pressing global health issue. While one would think that people may be more aware of IBD, given the rising number of IBD cases worldwide, more efforts are needed to rigorously evaluate if public awareness of IBD has trended up or not.

Since 2020, the dramatic challenges of the COVID-19 pandemic have greatly affected our lives, which might have affected the public interests for IBD as well. Because immunocompromised patients may be vulnerable to COVID-19, there were concerns about whether patients with IBD might be more susceptible to COVID-19 and have worser outcomes [27]. In a cross-sectional questionnaire, patients with IBD were apprehensive about the COVID-19 pandemic, as they felt more vulnerable to COVID-19 owing to their condition and their immunosuppressive therapies, including biologics. Many patients also felt disturbed, depressed, and tense when thinking of the infection [28]. To provide solid supports for patients with IBD during the pandemic, further efforts to increase public awareness of the entity are crucial. One example of a successful public health awareness campaign is the annual breast cancer awareness campaign [29], which achieved appropriate identification of targets, early involvement of the key stakeholders such as celebrities with the condition, and utilization of smartphone apps or eHealth platforms even during the current pandemic.

This study's strength is that this is the first hypothesis-generating study to see the extent of public awareness of IBD in the United States and worldwide by using the GT database. Using the open data, we could quantify the current trends of general interest in IBD. However, several limitations need to be addressed. First, owing to the nature of GT, the results of this study only included results from those who had internet access and sought health-related information via Google search. Given the high internet penetration rates—approximately 90.4% in North America and 60.1% worldwide [30]—and high US Google search market share of approximately 83% [31], GT is considered a good surrogate of public awareness. Second, GTs are proxies for engagement. Sentinel surveillance such as surveys may be needed to clarify the findings. Third, the potential effect on RSVs may lag the intervention by weeks, and it is uncertain how long the effects of the intervention would last, making it challenging to assess the impact of the intervention as RSVs. Fourth, there are confounders such as separate media coverage of the disease, which are difficult to identify and account for the uncertainty about how to attribute to the independent variable of interest. Further, incorporating the analysis of Google search query data of the actual public awareness campaigns (in this case, WID and CCAW) might be a preferred approach to reduce confounding factors and directly evaluate the effects of these health observances to garner an audience. However, the RSVs for WID and CCAW were too few to conduct a joinpoint analysis during the study period (Multimedia Appendix 1). No joinpoints were noted around the time of WID or CCAW throughout the period. Regarding WID, there were spikes in the RSV in the week of WID on May 19 (worldwide search in 2016 to 2020 and the US search in 2018 and 2020). Otherwise, the RSVs were consistently zero. For CCAW, GT could not return queries since there were too few Google searches using the term. Despite these limitations, our approach is interestingly novel to generate hypotheses on campaign effectiveness or ineffectiveness in the public awareness of IBD.

In conclusion, using the GT data as a surrogate, our study posed a possibility that WID and CCAW might not have successfully improved public awareness toward IBD. There is a need to look deeper into how to precisely assess the public awareness and improve public awareness using these health observances based on good examples.

## Authors' Contributions

KS, TC, and NK wrote the manuscript. YN proposed the study concept, designed the study, wrote the manuscript, and analyzed the data. HH and TK revised the manuscript critically. FO supervised the research.

## Conflicts of Interest

None declared.

Multimedia Appendix 1
Weekly relative search volumes of World Inflammatory Bowel Disease Day and Crohn's and Colitis Awareness Week in 2016-2020.
[PNG File , 71 KB - infodemiology_v1i1e32856_app1.png ]

## References

1. Kaplan GG. The global burden of IBD: from 2015 to 2025. Nat Rev Gastroenterol Hepatol 2015 Dec;12(12):720-727. [doi: 10.1038/nrgastro.2015.150] [Medline: 26323879]
2. Ouyang G, Pan G, Liu Q, Wu Y, Liu Z, Lu W, et al. The global, regional, and national burden of pancreatitis in 195 countries and territories, 1990-2017: a systematic analysis for the Global Burden of Disease Study 2017. BMC Med 2020 Dec 10;18(1):388 [FREE Full text] [doi: 10.1186/s12916-020-01859-5] [Medline: 33298026]

XSL•FO
RenderX

3.   Ye Y, Manne S, Bennett D. Prevalence of Inflammatory Bowel Disease in the U.S. Adult Population: Recent Estimates from Large Population-Based National Databases. American Journal of Gastroenterology 2018;113(Supplement):S373-S374. [doi: 10.14309/00000434-201810001-00654]

4.   Xu F, Carlson SA, Liu Y, Greenlund KJ. Prevalence of Inflammatory Bowel Disease Among Medicare Fee-For-Service Beneficiaries - United States, 2001-2018. MMWR Morb Mortal Wkly Rep 2021 May 14;70(19):698-701 [FREE Full text] [doi: 10.15585/mmwr.mm7019a2] [Medline: 33983913]

5.   Ananthakrishnan AN, Bernstein CN, Iliopoulos D, Macpherson A, Neurath MF, Ali RAR, et al. Environmental triggers in IBD: a review of progress and evidence. Nat Rev Gastroenterol Hepatol 2018 Jan;15(1):39-49. [doi: 10.1038/nrgastro.2017.136] [Medline: 29018271]

6.   No authors L. Announcement: World IBD Day - May 19, 2017. MMWR Morb Mortal Wkly Rep 2017 May 19;66(19):516 [FREE Full text] [doi: 10.15585/mmwr.mm6619a9] [Medline: 28520712]

7.   Take action for world IBD day!. Crohn's & Colitis Foundation. URL: https://www.crohnscolitisfoundation.org/WorldIBDDay [accessed 2021-07-20]

8.   S. Res. 199 Supporting the goals and ideals of "Crohn's and Colitis Awareness Week". The United States Congress. URL: https://www.congress.gov/bill/112th-congress/senate-resolution/199/text [accessed 2021-07-20]

9.   Purtle J, Roman LA. Health Awareness Days: Sufficient Evidence to Support the Craze? Am J Public Health 2015 Jun;105(6):1061-1065. [doi: 10.2105/ajph.2015.302621]

10.  National health observances 2021. U.S. Department of Health and Human Services. URL: https://health.gov/news/category/national-health-observances [accessed 2021-07-20]

11.  Wakefield MA, Loken B, Hornik RC. Use of mass media campaigns to change health behaviour. The Lancet 2010 Oct;376(9748):1261-1271. [doi: 10.1016/s0140-6736(10)60809-4]

12.  Nuti SV, Wayda B, Ranasinghe I, Wang S, Dreyer RP, Chen SI, et al. The use of google trends in health care research: a systematic review. PLoS One 2014;9(10):e109583 [FREE Full text] [doi: 10.1371/journal.pone.0109583] [Medline: 25337815]

13.  Google trends 2021. Google. URL: https://trends.google.com/trends/ [accessed 2021-07-20]

14.  Motosko C, Zakhem G, Ho R, Saadeh P, Hazen A. Using Google to Trend Patient Interest in Botulinum Toxin and Hyaluronic Acid Fillers. J Drugs Dermatol 2018 Nov 01;17(11):1245-1246. [Medline: 30500150]

15.  Frauenfeld L, Nann D, Sulyok Z, Feng Y, Sulyok M. Forecasting tuberculosis using diabetes-related google trends data. Pathog Glob Health 2020 Jul;114(5):236-241 [FREE Full text] [doi: 10.1080/20477724.2020.1767854] [Medline: 32453658]

16.  Patel JC, Khurana P, Sharma YK, Kumar B, Ragumani S. Chronic lifestyle diseases display seasonal sensitive comorbid trend in human population evidence from Google Trends. PLoS One 2018;13(12):e0207359 [FREE Full text] [doi: 10.1371/journal.pone.0207359] [Medline: 30540756]

17.  Tabuchi T, Fukui K, Gallus S. Tobacco Price Increases and Population Interest in Smoking Cessation in Japan Between 2004 and 2016: A Google Trends Analysis. Nicotine Tob Res 2019 Mar 30;21(4):475-480. [doi: 10.1093/ntr/nty020] [Medline: 29394419]

18.  Cacciamani GE, Bassi S, Sebben M, Marcer A, Russo GI, Cocci A, et al. Consulting "Dr. Google" for Prostate Cancer Treatment Options: A Contemporary Worldwide Trend Analysis. Eur Urol Oncol 2020 Aug;3(4):481-488. [doi: 10.1016/j.euo.2019.07.002] [Medline: 31375427]

19.  Havelka E, Mallen C, Shepherd T. Using Google Trends to assess the impact of global public health days on online health information seeking behaviour in Central and South America. J Glob Health 2020 Jun;10(1):010403 [FREE Full text] [doi: 10.7189/jogh.10.010403] [Medline: 32373327]

20.  Patel JC, Khurana P, Sharma YK, Kumar B, Sugadev R. Google trend analysis of climatic zone based Indian severe seasonal sensitive population. BMC Public Health 2020 Mar 12;20(1):306 [FREE Full text] [doi: 10.1186/s12889-020-8399-0] [Medline: 32164654]

21.  Peng Y, Li C, Rong Y, Chen X, Chen H. Retrospective analysis of the accuracy of predicting the alert level of COVID-19 in 202 countries using Google Trends and machine learning. J Glob Health 2020 Dec;10(2):020511 [FREE Full text] [doi: 10.7189/jogh.10.020511] [Medline: 33110594]

22.  Russo GI, di Mauro M, Cocci A, Cacciamani G, Cimino S, Serefoglu EC, EAU-YAU Men's Health Working Group. Consulting "Dr Google" for sexual dysfunction: a contemporary worldwide trend analysis. Int J Impot Res 2020 Jul;32(4):455-461. [doi: 10.1038/s41443-019-0203-2] [Medline: 31591474]

23.  Sharma M, Sharma S. The Rising Number of COVID-19 Cases Reflecting Growing Search Trend and Concern of People: A Google Trend Analysis of Eight Major Countries. J Med Syst 2020 May 20;44(7):117 [FREE Full text] [doi: 10.1007/s10916-020-01588-5] [Medline: 32430650]

24.  Brodeur A, Clark AE, Fleche S, Powdthavee N. COVID-19, lockdowns and well-being: Evidence from Google Trends. J Public Econ 2021 Jan;193:104346 [FREE Full text] [doi: 10.1016/j.jpubeco.2020.104346] [Medline: 33281237]

25.  Zitting K, Lammers-van der Holst HM, Yuan RK, Wang W, Quan SF, Duffy JF. Google Trends reveals increases in internet searches for insomnia during the 2019 coronavirus disease (COVID-19) global pandemic. J Clin Sleep Med 2021 Feb 01;17(2):177-184. [doi: 10.5664/jcsm.8810] [Medline: 32975191]

XSL·FO
RenderX

26.    Joinpoint trend analysis software 2021. National Cancer Institute. URL: https://surveillance.cancer.gov/joinpoint/ [accessed 2021-07-20]

27.    Nakase H, Matsumoto T, Matsuura M, Iijima H, Matsuoka K, Ohmiya N, et al. Expert Opinions on the Current Therapeutic Management of Inflammatory Bowel Disease during the COVID-19 Pandemic: Japan IBD COVID-19 Taskforce, Intractable Diseases, the Health and Labor Sciences Research. Digestion 2021;102(5):814-822 [FREE Full text] [doi: 10.1159/000510502] [Medline: 32892197]

28.    Zingone F, Siniscalchi M, Savarino EV, Barberio B, Cingolani L, D'Incà R, et al. Perception of the COVID-19 Pandemic Among Patients With Inflammatory Bowel Disease in the Time of Telemedicine: Cross-Sectional Questionnaire Study. J Med Internet Res 2020 Nov 02;22(11):e19574 [FREE Full text] [doi: 10.2196/19574] [Medline: 33006945]

29.    Glynn RW, Kelly JC, Coffey N, Sweeney KJ, Kerin MJ. The effect of breast cancer awareness month on internet search activity--a comparison with awareness campaigns for lung and prostate cancer. BMC Cancer 2011 Oct 12;11:442 [FREE Full text] [doi: 10.1186/1471-2407-11-442] [Medline: 21993136]

30.    Global internet penetration rate as of April 2021, by region 2021. Statista. URL: https://www.statista.com/statistics/269329/penetration-rate-of-the-internet-by-region/ [accessed 2021-07-20]

31.    Google: search engine market share in selected countries 2021. Statista. URL: https://www.statista.com/statistics/220534/googles-share-of-search-market-in-selected-countries/ [accessed 2021-07-20]

## Abbreviations

**CCAW:** Crohn's and Colitis Awareness Week
**GT:** Google trends
**IBD:** inflammatory bowel disease
**RSV:** relative search volume
**WID:** World IBD Day

---

XSL•FO
**RenderX**

<u>Original Paper</u>

# Difficulty Regulating Social Media Content of Age-Restricted Products: Comparing JUUL's Official Twitter Timeline and Social Media Content About JUUL

Danny Valdez[1], PhD; Jennifer B Unger[2], PhD

[1]Department of Applied Health Science, Indiana University School of Public Health, Bloomington, IN, United States
[2]Keck School of Medicine, University of Southern California, Los Angeles, CA, United States

**Corresponding Author:**
Danny Valdez, PhD
Department of Applied Health Science
Indiana University School of Public Health
1025 E 7th Street, #111
Bloomington, IN, 47405
United States
Phone: 1 8128551561
Email: danvald@iu.edu

## *Abstract*

**Background:**   In 2018, JUUL Labs Inc, a popular e-cigarette manufacturer, announced it would substantially limit its social media presence in compliance with the Food and Drug Administration's (FDA) call to curb underage e-cigarette use. However, shortly after the announcement, a series of JUUL-related hashtags emerged on various social media platforms, calling the effectiveness of the FDA's regulations into question.

**Objective:**   The purpose of this study is to determine whether hashtags remain a common venue to market age-restricted products on social media.

**Methods:**   We used Twitter's standard application programming interface to download the 3200 most-recent tweets originating from JUUL Labs Inc's official Twitter Account (@JUULVapor), and a series of tweets (n=28,989) from other Twitter users containing either #JUUL or mentioned JUUL in the tweet text. We ran exploratory (10×10) and iterative Latent Dirichlet Allocation (LDA) topic models to compare @JUULVapor's content versus our hashtag corpus. We qualitatively deliberated topic meanings and substantiated our interpretations with tweets from either corpus.

**Results:**   The topic models generated for @JUULVapor's timeline seemingly alluded to compliance with the FDA's call to prohibit marketing of age-restricted products on social media. However, the topic models generated for the hashtag corpus of tweets from other Twitter users contained several references to flavors, vaping paraphernalia, and illicit drugs, which may be appealing to younger audiences.

**Conclusions:**   Our findings underscore the complicated nature of social media regulation. Although JUUL Labs Inc seemingly complied with the FDA to limit its social media presence, JUUL and other e-cigarette manufacturers are still discussed openly in social media spaces. Much discourse about JUUL and e-cigarettes is spread via hashtags, which allow messages to reach a wide audience quickly. This suggests that social media regulations on manufacturers cannot prevent e-cigarette users, influencers, or marketers from spreading information about e-cigarette attributes that appeal to the youth, such as flavors. Stricter protocols are needed to regulate discourse about age-restricted products on social media.

XSL·FO
**RenderX**

## Introduction

Following the Food and Drug Administration's (FDA's) call to curb underage e-cigarette use and increasing criticism of JUUL's youth-oriented "Vaporized" campaign [1], JUUL Labs Inc announced it would limit its social media presence. As part of the FDA agreement, JUUL deleted its official Facebook and Instagram accounts, reduced its Twitter activity, and removed older Twitter posts that could be attractive to youth or interpreted as marketing to youth [1].

FDA can regulate what JUUL and other e-cigarette manufacturers can post on official social media platforms [2]. However, the FDA cannot regulate posts about JUUL by customers or influencers, who can identify their posts as JUUL-related by using hashtags—short words or phrases preceded by the '#' symbol that label the content of a social media post and cause the post to appear in users' keyword searches. Hashtags spread social media content rapidly [3] and are therefore used for branding and marketing of certain products for mainstream appeal [4]. Any social media user (including paid or unpaid social media influencers, retailers, or enthusiastic consumers) can use hashtags to spread content about any topic, including age-restricted products subject to federal regulations. Alcohol, for example, is heavily marketed through hashtags [5-7], though much social media content about alcohol does not originate from official corporate accounts.

Marketing research further suggests that hashtags are used as branding or marketing ploys to promote age-restricted products including alcohol [7] and tobacco [8] on social networking websites. Using hashtags for age-restricted products may help circumvent age-gates, which are already proven to be ineffective at deterring underage engagement with age-restricted products [9]. Thus, any effort by JUUL Labs Inc to curb marketing to underage users may be stunted by the presence (and popularity) of vaping-related hashtags not subject to regulations imposed on manufacturers.

Indeed, the prevalence of JUUL-related hashtags on Instagram increased after JUUL reduced its own social media presence [1]. This suggests a limitation to FDA regulations wherein age-restricted products can still be marketed separately from official company platforms. By consequence, age-restricted items that are popular among youth including alcohol, tobacco, and e-cigarettes remain overtly visible and marketable to this audience, despite official corporate positions that denounce such use.

Regulation of harmful social media content is a critical public health issue [10]. To our knowledge, however, no study has compared verified corporate accounts versus similarly related hashtags from noncorporate posters to examine the effectiveness of social media regulation efforts. This study uses an inductive approach and natural language processing (NLP) modeling to examine differences in JUUL's official, regulated Twitter account @JUULVapor, JUUL-related content posted on social media. Our study is guided by three research questions:

1. Do people's social media posts about JUUL provide evidence of greater reach and visibility of JUUL Labs Inc's official Twitter account @JUULVapor?
2. Can we leverage Latent Dirichlet Allocation (LDA) topic models to dissect JUUL-related corpora?
3. What are salient content differences between @JUULVapor and tweets published by social media users containing either #JUUL or "JUUL" within the tweet's text?

Collectively, findings from our study will contribute to discourse about information diffusion via social media. We hypothesize that despite JUUL's efforts to scrub their social media platforms of youth-oriented content, hashtags about JUUL remain pervasive and highly visible to youth.

## Methods

### Data

Data for this study were procured by leveraging Twitter's application programming interface (API). From the API, we collected two corpora unique to this study: (1) @JUULVapor's Twitter timeline (n=3200 tweets, the maximum number of most recent tweets posted by a single user allotted for download through the standard API), hereafter referred to as the @JUULVapor corpus (January 1 to May 31, 2021) and (2) a 1-month collection of tweets containing #JUUL or "JUUL" (n=29,989 tweets), referred to as the #JUUL corpus (May 1 to June 1, 2021). For the #JUUL corpus, specifically, we performed a Bot analysis [11] to remove tweets that originated from nonhuman accounts (n=135). No Bot analysis was required for the @JUULVapor corpus considering those tweets were pulled from JUUL Lab Inc's official Twitter account. We performed this procedure to ensure that discourse captured in subsequent analyses originated from humans and not an automated program. Upon removing bot accounts, 2 raters independently reviewed the text of the #JUUL tweets and removed any from the corpus that were not expressly about e-cigarettes or vaping (n=23). An author of this study also cross-checked tweet IDs in either corpus to ensure there was no accidental overlap in tweets (ie, the same tweet appearing in both the @JUUL and #JUUL corpora). Note that our total sample inclusive of both corpora (N=33,189 tweets) exceeds the mean observed sample size of collected tweets in a meta-analysis of public health social media studies (n=10,000) [12].

### Analysis

Our research questions are exploratory. Thus, we chose to use LDA topic models, a Bayes-driven, unsupervised NLP method, to examine differences in themes for the @JUULVapor and #JUUL corpora. LDA and related topic modeling analyses have been similarly leveraged in other health contexts, including studying discourse about the COVID-19 pandemic [13] and map themes among corpora of age-restricted products [14].
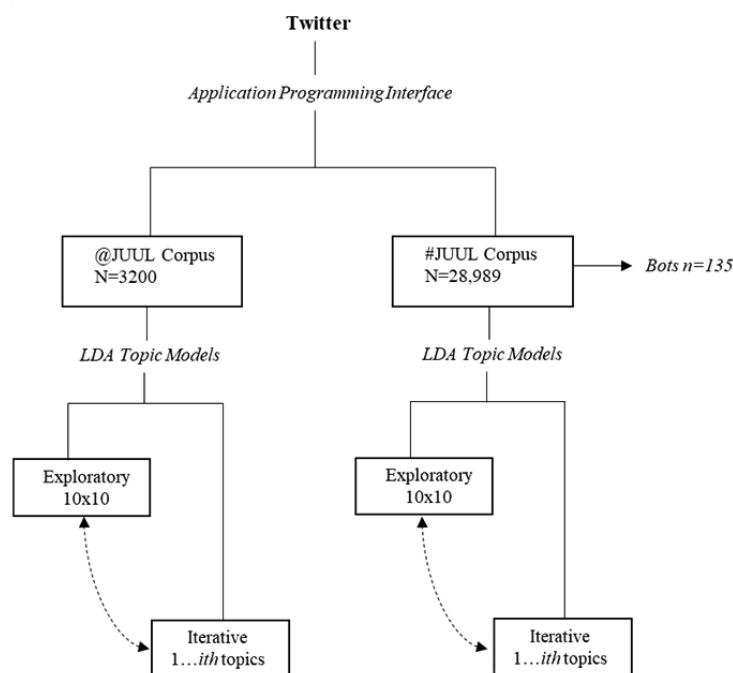
While previous studies generate topic models for differing corpora, qualitatively review differences between corpora, and discuss the meaning of those differences, our study takes a 2-step approach. The first step broadly examines themes for a fixed set of topics or words per topic (ie, 10 topics and 10 words per topic). Valdez et al [15] have provided examples of

exploratory topic models in practice. The second step uses an iterative topic model analysis that meta-analytically generates models with an increasing number of topics per corpus (ie, 1 topic, 2 topics…20 topics) [16]. This analysis generates a coherence score for each iteration, such that higher scores are equated with better model fit and interpretability. We used this second analysis to identify the optimal number of topics per corpus and further refine the models (ie, eliminate redundancy and noise) for maximum interpretability. To ensure the validity of our coherence scores, we selected a random sample of 50 tweets per corpus and matched each tweet's content to a respective theme identified by the topic model. We successfully placed each tweet within a topic, suggesting our topic models were both coherent and precise.

## Procedure

Our workflow is detailed in Figure 1. Upon downloading and cleaning the @JUULVapor and #JUUL corpora, we performed the following. First, we calculated standard descriptive statistics for each corpus, including the average number of likes, retweets, and number of tweets that originated from Verified accounts, or accounts reviewed to ensure they are owned and operated by a specific person (research question 1). Second, we performed an exploratory 10×10 topic model for the #JUUL and @JUULVapor corpora and qualitatively compared differences between them. Lastly, we performed an iterative analysis to identify the optimal number of topics and again qualitatively reviewed the topic model for each corpus for differences (research questions 2 and 3).

**Figure 1.** Conceptual framework guiding our study. LDA: Latent Dirichlet Allocation.



## Ethical Use of Data

All procedures and analyses undertaken in this study conform to the Twitter's terms for data use agreement. Our study was exempt from institutional review board review, given the secondary nature of this data collection and analysis.

## *Results*

### Descriptive Differences

We identified differences in total retweets and favorites per corpus. On average, content in the @JUULVapor corpus, JUUL Lab Inc's official Twitter handle, was retweeted 1.29 times (SD 16.77 times) and favorited 0.25 times (SD 4.49 times). For the #JUUL corpus, tweets were on average favorited 0.41 times (SD 4.74 times) and retweeted 4.53 times (SD 52.07 times). Exactly 237 tweets in the #JUUL corpus originated from verified twitter accounts. Given that such a marginal number of tweets originated from a verified account, we did not perform statistical tests to determine whether scope and reach were significantly different between verified and nonverified accounts.

### Exploratory Topic Models

Table 1 outlines an exploratory topic model for the @JUULVapor corpus. This topic model represents a condensed version of JUUL Lab Inc's 3200 most recent tweets delineated by 10 topics and 10 words per topic. The themes in the @JUULVapor's topic model were generally interpretable. Five of the topics in the @JUULVapor topic model contained references related to customer support or product warrant-related queries—which we interpreted as responses to complaints about JUUL associated products. Words recurrent among this body of topics include *please*, *dm* [direct message], *sorry*, *thank*, *contact*, and *customer* for support-oriented topics; and *JUUL*, *device*, *limited*, *warranty*, *information* for warranty-related topics. This model also referenced adult product use (ie, *legal*, *adults*, and *age*) and acknowledgement of underage use and underage use prevention (ie, *underage*, *prevention*, *minors*, *market*, and *seriously*)—which we interpreted as JUUL Lab Inc's forthrightly attempt to address controversies long associated with its brand. Notably, there were very few references to controversial topics such as flavors, addiction, and other drugs such as cannabis—which we interpreted as JUUL

Lab Inc's seeming attempt to distance itself from controversial aspects also associated with its brand. Other topics that emerged in this model including "Recycling" and "Warranty." Recycling-associated tweets generally referenced the importance of recycling used JUUL cartridges (which are disposable). We interpreted warranty as a topic related to customer support—namely ways in which customers can secure refunds if products are defective.

Table 2 outlines an exploratory topic model for the #JUUL corpus, which is a random collection of tweets discussing JUUL but not originating from JUUL's official Twitter timeline. This topic model represents a condensed version of a months' worth of tweets about JUUL, which were identified by either using #JUUL or containing the word "JUUL" in each tweet's text. The themes in the #JUUL topic model were somewhat interpretable, though less so than the @JUULVapor corpus reflecting greater content diversity. For example, topics that were somewhat vague, yet still referenced vaping, were labeled as a "General Vape" topic, which comprised seemingly unrelated words related to various aspects of vaping but not necessarily related to JUUL as a brand. Words recurrent among these topics include *vape*, *JUUL*, *hit*, *smoke*, *smoking*, *take*, and others. Beyond vague references to vaping, several clearer topics also emerged; these include a topic about marijuana and cannabis, the intersection of vaping and cigarettes, and a topic about nicotine, which we collectively interpreted as youth-appealing narratives about vaping. Note, topics consisting of "nicotine" or "flavors" are entirely absent from the @JUUL 10×10 topic model, which may be indicative of JUUL Lab Inc's attempt to distance itself from web-based controversies and present a cleaner image.

**Table 1.** A 10×10 topic model of JUUL Lab Inc's official Twitter timeline (n=3200 tweets).

| Product warranty | Customer support | Purchase | Customer support | Customer support | Smoking | Underage use | Unclear | Recycle | Warranty |
|---|---|---|---|---|---|---|---|---|---|
| always | please | hi | team | please | adult | Use | hear | reaching | JUUL |
| ready | dm | JUUL | know | team | JUUL | Underage | switch | thanks | device |
| year | hi | JUUL Pod | working | care | smokers | JUUL | hi | currently | limited |
| limited | hey | age | definitely | sorry | labels | take | JUUL | products | year |
| warranty | assist | price | product | hear | lives | prevention | making | program | hey |
| hi | number | retail | want | help | world | seriously | hey | team | warranty |
| customer | lock | legal | frustrating | customer | billion | hi | helped | place | access |
| hey | information | packs | must | contact | improving | minors | thanks | always | products |
| replace | case | adults | thanks | hi | mission | product | congrats | recycle | hi |
| customer | additional | available | thank | hey | high | market | leaking | quality | submit |

**Table 2.** A 10×10 topic model of a collection of tweets referencing JUUL Lab Inc and vaping. Redacted tweets refer to specific mentions which cannot be published per Twitter's data use agreement.

| General vape | Marijuana | General vape | Cigarettes | Unclear | Flavor | Nicotine | General vape | General vape | General vape |
|---|---|---|---|---|---|---|---|---|---|
| vape | vape | vape | vape | vape | JUUL | vape | vape | vape | vape |
| pen | cbd | JUUL | lung | JUUL | pods | vaping | vaping | police | JUUL |
| kid | cannabis | hit | start | mom | pod | nicotine | ublo | cops | though |
| city | mg | time | use | actually | hit | smoking | lungs | vaping | smoke |
| police | juice | think | help | two | vape | JUUL | covid | f*****g | smoking |
| ocean | weed | s**t | cigarettes | care | mango | shop | ml | people | take |
| cops | edibles | smoke | enough | vaping | days | tobacco | already | smoke | hookah |
| REDACTED | cbd oil | pen | cause | baby | lost | cigarettes | ecig | white | ur |
| vaping | liquid | f**k | says | freaking | mint | quit | vapehop | want | free |
| take | thc | cana | safe | stop | spring | lungs | vapelife | smoking | oh |

## Iterative Topic Model

Iterative analysis revealed the optimal number of topics given the total number of words in each corpus. Figure 2 plots the coherence score, which measures the semantic similarity of words in each topic, per corpus [17]. Peaks in the graphs denote the optimal number of topics for each corpus.

For the @JUULVapor corpus, there were 2 optimal topics (coherence score=0.36) (Textbox 1). Both topics were interpreted as referring to either responses to customer's

concerns or complaints about JUUL products. Topics from the general topic model that centered on underage use, purchasing, and recycling were absent. This may suggest, at least partially, that the renewed purpose of JUUL Lab Inc's official Twitter account is to field customer complaints and comments.

For the #JUUL corpus, there were 4 optimal topics (Coherence score=0.50) (Textbox 2). These topics were more diverse than the @JUULVapor corpus; containing topics related to marijuana, vape/smoking, general vaping, and vaping-related damage. Here, there is more emphasis on the elicit side of vaping/smoking, and youth appealing narratives. These topics stand in sharp contrast with the #JUUL Optimal Topic Model (Textbox 1), which only revealed customer support–related topics.

**Figure 2.** Coherence score plot by corpus. The X axis represents the total number of topics; the Y axis represents coherence score per iteration.



**Textbox 1.** Topics sharply contrasting the #JUUL Optimal Topic Model (coherence score=0.36). The topics in bold represent the theme; bulleted words represent the words per topic.

**Customer complaints**

- JUUL
- hi
- use
- hey
- frustrating
- products
- must
- products
- working
- definitely

**Customer support**

- please
- team
- sorry
- hear
- care
- hi
- hey
- customer
- help
- contact

**Textbox 2.** Iterative topic model for #JUUL (coherence score=0.5). The topics in bold represent the theme; bulleted words represent the words per topic.

**Marijuana**

- vape
- cbd
- juice
- vaping
- pen
- vapejuice
- shope
- vapelife
- weed
- cannabis

**General vape**

- vape
- vaping
- people
- smoke
- pen
- new
- lungs
- covid
- already
- day

**Vaping/smoking**

- JUUL
- vape
- pods
- hit
- nicotine
- smoke
- back
- day
- think
- hitting
- mango

**Vaping damage**

- vape
- use
- lungs
- cigarettes
- cause
- enough

- safe

- damage

- irreversible

- vanishvaping

- clear

## *Discussion*

### Principal Findings

This study examined the use of hashtags to indirectly market age-restricted products on social media. We leveraged the Twitter API to archive and compare 2 corpora specific to e-cigarette use with LDA topic models. The first corpus (ie, @JUULVapor) contained 3200 tweets derived from JUUL Lab Inc's official Twitter account. The second corpus (ie, #JUUL) contained a month's worth of tweets (May 1 to June 1, 2021) that contained #JUUL or mentioned JUUL within the tweet text (n=28,989). When the corpora were compared, we identified several telling observations within each corpus, which showcase disparate uses in @JUULVapor vs #JUUL. These partially include the @JUULVapor corpus seeming compliance to prevent underage marketing versus an array of random, youth-appealing content in the #JUUL corpus. Below we discuss these differences within the scope of their current literature delineated by each research question.

### RQ1: Evidence of Greater Reach in the #JUUL Corpus

Our first research question asked whether content about JUUL contained evidence of greater reach and visibility than content posted on JUUL's official Twitter account. Reach and visibility, here, was measured by the average number of likes and retweets per tweet in each corpus. Our findings suggest that, overall, content about JUUL, and e-cigarettes more broadly, are clearly visible on social media spaces via hashtags. This finding corroborates a large body of work that suggests hashtags are often used to quickly distribute branding and product marketing information [4,18,19].

Regarding content, tweets in the #JUUL corpus were, on average, retweeted and liked with greater frequency than content posted by @JUULVapor. That content in the #JUUL corpus was retweeted more often than @JUUL is perhaps not entirely surprising. As mentioned previously (and throughout the remainder of the discussion), content in the #JUUL corpus contained topics of discussion that are inherently appealing to youth, versus content in the @JUUL corpus that seemed to unilaterally focus on customer complaints. For example, in the #JUUL corpus, we observed mentions of flavors (ie, mint, mango, and cucumber), cannabis vaping (ie, vape cartridges), and meme/joke-sharing, all of which are inherently conducive appealing to youth and higher post engagement. Additionally, as hashtags are used for rapid content organization of content, it is likely any social media user (agnostic of age differences) can see content posted by #JUUL, including those who did not expressly seek this information themselves.

### RQ2: LDA Topic Models as Tools to Contrast Corporate Corpora With an Assortment of Related Tweets

Beyond the reach and scope of tweets, we also investigated whether LDA could be leveraged to identify content differences in corporate versus lay user social media accounts. LDA topic models have been historically leveraged in an exploratory capacity to consolidate an overwhelming amount of text data into manageable chunks (ie, themes) that represent the most salient components of that text data [20,21]. For example, prior studies have used topic models to explore the underlying thematic structures across a broad range of corpora, including studying discourse about societal events [13], identifying alcohol branding strategies [14], and mapping publication histories of leading Health journals [22]. As topic models become increasingly used in the social and medical sciences, it remains debatable how these models can be used to test applied, rather than exploratory, hypotheses [22]. This includes ample discussion how topic models could theoretically be used to inform possible digital e-health interventions [23,24] and to construct bots from topic modeling data that meaningfully identify mental health distress [25].

To our knowledge, LDA topic models have not been used in either exploratory or applied capacities to compare social media content originating from a specific corporation and a collection of tweets about that product (though not necessarily originating from the corporate account). Our findings show that such models can be leveraged for this purpose, evidenced by our findings that identified qualitative differences in content between the @JUULVaporVapor and #JUUL corpora. We used exploratory models as a standardized metric to generate the same number of topics and words per topic for each corpus. We then ran iterative models to identify the optimal number of topics within each corpus (ie, improve granularity and precision of the models). Gethers and Poshyvanyk [26] provide more insight into granularity and relational topic models.

We contend the combined use of exploratory and iterative model may provide a conceptual framework for future topic modeling studies. For example, exploratory models may uncover broad themes in a corpus. Iterative models will then only identify highly salient (or themes of highest priority) given a corpus. The range of topics uncovered by the iterative models may highlight how broad or narrow the corpus is in scope—more themes equate to broad content in a corpus, few themes indicate narrow scope or focus. For our study @JUULVapor's two optimal topics, contrasted with four in the #JUUL corpus, suggests the content in the @JUULVapor corpus was much narrower and more defined; for @JUULVApor, that is customer support. More topics in the #JUUL corpus suggest the content

was more diverse, containing a wider array of underlying themes; that is, more youth-appealing narrative. More research is needed to identify optimal use of exploratory and precision topic models in a research context. However, we encourage the use of both exploratory and iterative models when comparing corpora of vastly different sizes.

## RQ3: Implications for Content differences between @JUULVapor and #JUUL

Our final research question posited whether content differences identified between corpora were meaningful. Across each analysis, we identified differences that clearly distinguished each corpus, including vastly different ways in which e-cigarettes were mentioned and discussed between @JUULVapor and #JUUL. This includes, as mentioned, a narrow scope of content in the @JUULVapor corpus, versus more diverse, often youth-appealing content in the #JUUL corpus. This finding, coupled with increased engagement in the #JUUL corpus supports extant research that hashtags are effective means of disseminating age-inappropriate content rapidly [3,27]. Nonetheless, deeper insights into topic nuance are needed.

First, cursory insights into JUUL Lab Inc's corporate Twitter account show a seeming attempt to comply with FDA regulations barring youth marketing. In 2018, JUUL Labs Inc had been accused of using corporate social media accounts to market to youth and, in compliance with court orders and regulations, scrubbed their social media histories of youth-appealing content. Our inability to collect *any* deleted tweets suggests a natural limitation to social media research; namely that deleted content is truly removed from archives and cannot be accessed. However, remaining tweets posted by @JUULVapor—that is, those analyzed in this study—showcase a semiactive Twitter account almost entirely devoid of marketing content. Indeed, both exploratory and iterative topic models, the majority of topics and words per topic for @JUULVapor were customer support oriented. A review of individual tweets further revealed that the majority of posts were corporate response to complaints about JUUL products (eg, TWEET *Hi there, we're sorry to hear that. You can access troubleshooting tips for your JUUL device at…*). This shift in content may indicate that JUUL is trying to position itself as a responsible company, similar to the corporate responsibility advertising campaigns used by Big Tobacco companies to present a respectable image while selling a dangerous product (eg, TWEET *Minors should not use any nicotine product and we take the prevention of underage use of JUUL very seriously*) [28].

By contrast, themes in the #JUUL exploratory and iterative models were more diverse and contained several references that may be appealing to youth. For example, the #JUUL corpus contained references to cartridge flavors, which have been banned in the United States because they are attractive to youth but are still legal in disposable JUUL-like products [29]. Although the JUUL company is no longer actively promoting flavors, it appears users continue to associate the JUUL product with flavors, including mango, cucumber, mint, and others. Beyond flavors, we also observed a high co-occurrence of

flavors with "marijuana." Marijuana was prominent in the exploratory #JUUL model (ie, topic 2) and retained its prominence during the iterative model. This suggests a significant portion of the #JUUL corpus contained references to cannabis. Interestingly, few tweets or topics directly mentioned JUUL (the company). JUUL not being expressly mentioned topics indicates few tweets expressly mentioned JUUL and marijuana together in the same post. However, despite not mentioning the brand directly, the web-based conversation regularly discussed the use of vape products for marijuana, which may be at least partially explained by JUUL's evolution in mainstream vernacular form a noun (ie, JUUL products) to a verb (ie, JUUL'ing, a specific and colloquial term for "vaping"). We also observed profanity in the #JUUL topic models, which was entirely absent in @JUULVapor (eg, TWEET *Bro, we are all [expletive] high on this vape*). Profanity may indicate the presence of younger social media users [30].

Although the #JUUL corpus contained youth-appealing content (ie, profanity, high mentions of marijuana, and flavors among other indicators), we also observed topics in both @JUULVapor and #JUUL corpora detailing antivaping-related advocacy. For the @JUULVapor corpus, perhaps unsurprisingly, this may provide evidence that JUUL Labs Inc complied with court orders to stop marketing to underage users (eg, TWEET *today we're implementing a series of new measures that build upon or existing efforts to reduce underage use*). For the #JUUL corpus, this may also suggest a substantive body of antivaping-related advocacy that adopted a hashtag strategy to spread their messages more effectively (eg, #vanishvaping). However, evidence of antivape advocacy in either corpus does not suggest that the messaging effectively deters youth use or substantively changes the wider web-based conversation. Rather, in some cases, there seem to be additional sarcastic comments that offset antivape messaging (eg, TWEET *All these anti-vape adds make me want to snort meth*). Thus, such policies and court orders are ineffective in regulating the totality of messages received by underage users, particularly given that nonofficial tweets were more likely to be shared/favorited than official @JUULVapor messages. Antimessaging campaigns of other age-restricting products have also shown to have wide reach but inconclusive results [31], which may suggest that despite the best efforts of antivape advocates, their behavior change attempts may fail.

## Implications for Social Media Messaging About Vaping

Together, these findings further demonstrate the overall lack of control the official @JUULVaporVapor account has in directing web-based conversations about vape products. Despite JUUL presenting a "clean image," their brand remains associated with a dangerous and addictive product that is naturally appealing to youth. However, it is also clear that more research on who is tweeting about JUUL and vaping, and how hashtags facilitate marketing illicit behavior, is needed. Future research should consider adding deep learning models to partition tweets about vaping by demographic variables to, among other matters, predict the likelihood an account posted about JUUL was that of an underage user.

From a public health/medical/interventionist perspective, our findings also compel us to ponder how mining communication patterns (ie, tracking discourse about JUUL) can be further leveraged to identify intervention targets promoting antivape messaging. In this study, the sharp divide in content between @JUUL and #JUUL suggest that provape messaging did not end after JUUL Lab Inc's court order to curb marketing efforts. Rather, marketing manifested through shared user content about products that remain popular while not necessarily referencing the elicit product. It is indeed possible that some influencers are paid by manufacturers such as JUUL to surreptitiously market products without seeming association with the brand. This was a strategy used by Big Tobacco to continue marketing indirectly while seemingly complying with antimarking efforts. Additionally, on Twitter, it is difficult to determine which posts from celebrity and other verified accounts are paid advertisements. Other platforms, including Instagram and Facebook, expressly designate ads with a "#ad" notice; however, this is less common with Twitter. Policy efforts should also center on clearer guidelines for designating paid or sponsored posts versus regular posts on Twitter.

## Limitations

Our study is subject to limitations we hope to address in future work. First, we acknowledge a likely demographic bias inherent to social media studies. This includes a sample that likely skews younger, male, wealthier, and whiter than the general population [32]. Given that our study was exploratory, we also did not control spatial and geographic patterns in social media data, which affect how and what users post within a given time (ie, rural vs urban settings, or older users posting earlier in the morning than younger users) [33,34]. Regarding data analysis, we acknowledge that we did not perform a formal qualitative analysis with these data. Topic models were used to, instead, consolidate each corpus and allow us to draw inferences about the data from those topics. Because of sample size constraints, we were also unable to draw meaningful comparisons between verified and nonverified users in the #JUUL corpus. However, despite these limitations, we believe gaps in our study present opportunities for future research related to social media discourse on age restricted products. This includes performing other NLP methodologies with similar data (ie, sentiment analysis) to understand polarity in discourse or applying classifiers to accounts tweeting about age-restricted products to predict age and gender among other demographic traits. Dai et al [35] have provided further information about the M3 classifier in a health context. Such studies may provide a deeper and more nuanced landscape of social media discourse related to age-restricted products.

## Conclusions

@JUUVapor may be compliant with web-based marketing restrictions and promoting antivaping messaging. However, JUUL Labs Inc is powerless to control the larger narrative about vaping on social media. Indeed, hashtags about vaping and JUUL contain much of the youth-directed content that led to the initial impositions placed on JUUL by the FDA. Our results underscore the difficulty of regulating social media content despite federal impositions that ban marketing of age-restricted products in web-based spaces. Given the limited ability of social media to restrict what underage users see on their websites, companies can bypass marketing regulations by allowing users to freely share related hashtags or by paying social media influencers to disseminate these hashtags [36]. Although it is impossible to regulate free speech on the internet, perhaps public health advocates could harness the power of hashtags to deliver antivaping messages, though the effectiveness of such campaigns is not guaranteed.

## Conflicts of Interest

None declared.

## References

1. Czaplicki L, Tulsiani S, Kostygina G, Feng M, Kim Y, Perks SN, et al. #toolittletoolate: JUUL-related content on Instagram before and after self-regulatory action. PLoS One 2020;15(5):e0233419 [FREE Full text] [doi: 10.1371/journal.pone.0233419] [Medline: 32437397]
2. Jackler RK, Li VY, Cardiff RAL, Ramamurthi D. Promotion of tobacco products on Facebook: policy versus practice. Tob Control 2019 Jan;28(1):67-73. [doi: 10.1136/tobaccocontrol-2017-054175] [Medline: 29622602]
3. Rauschnabel P, Sheldon P, Herzfeldt E. What motivates users to hashtag on social media? Psychol Mark 2019 Jan 15;36(5):473-488 [FREE Full text] [doi: 10.1002/mar.21191]
4. Bernard A. Theory of the Hashtag. Hoboken, NJ: John Wiley & Sons; 2019.
5. Moreno MA, Whitehill JM. Influence of Social Media on Alcohol Use in Adolescents and Young Adults. Alcohol Res 2014;36(1):91-100 [FREE Full text] [Medline: 26259003]
6. Burton S, Dadich A, Soboleva A. Competing Voices: Marketing and Counter-Marketing Alcohol on Twitter. J Nonprofit Public Sect 2013 Apr;25(2):186-209. [doi: 10.1080/10495142.2013.787836]
7. Nicholls J. Everyday, everywhere: alcohol marketing and social media--current trends. Alcohol Alcohol 2012;47(4):486-493. [doi: 10.1093/alcalc/ags043] [Medline: 22532575]
8. Astuti PAS, Assunta M, Freeman B. Raising generation 'A': a case study of millennial tobacco company marketing in Indonesia. Tob Control 2018 Jul;27(e1):e41-e49. [doi: 10.1136/tobaccocontrol-2017-054131] [Medline: 30042229]
9. Barry A, Primm K, Russell H, Russell A. Characteristics and Effectiveness of Alcohol Website Age Gates Preventing Underage User Access. Alcohol Alcohol 2021 Jan 04;56(1):82-88. [doi: 10.1093/alcalc/agaa090] [Medline: 33098290]

10. Hoek J, Jones SC. Regulation, public health and social marketing: a behaviour change trinity. J Soc Mark 2011 Feb 11;1(1):32-44. [doi: 10.1108/20426761111104419]

11. Gilani Z, Wang L, Crowcroft J, Almeida M, Farahbakhsh R. Stweeler: A Framework for Twitter Bot Analysis. 2016 Presented at: WWW '16: 25th International World Wide Web Conference; April 11-15, 2016; Montréal, QC. [doi: 10.1145/2872518.2889360]

12. Edo-Osagie O, De La Iglesia B, Lake I, Edeghere O. A scoping review of the use of Twitter for public health research. Comput Biol Med 2020 Jul;122:103770 [FREE Full text] [doi: 10.1016/j.compbiomed.2020.103770] [Medline: 32502758]

13. Valdez D, Ten Thij M, Bathina K, Rutter LA, Bollen J. Social Media Insights Into US Mental Health During the COVID-19 Pandemic: Longitudinal Analysis of Twitter Data. J Med Internet Res 2020 Dec 14;22(12):e21418 [FREE Full text] [doi: 10.2196/21418] [Medline: 33284783]

14. Barry AE, Valdez D, Padon AA, Russell AM. Alcohol Advertising on Twitter—A Topic Model. Am J Health Educ 2018 Jun 29;49(4):256-263. [doi: 10.1080/19325037.2018.1473180]

15. Valdez D, Picket AC, Young B, Golden S. On Mining Words: The Utility of Topic Models in Health Education Research and Practice. Health Promot Pract 2021 May;22(3):309-312. [doi: 10.1177/1524839921999050] [Medline: 33759597]

16. Anoop VS, Prem Sankar C, Asharaf S, Zonin A. Generating and visualizing topic hierarchies from microblogs: An iterative latent dirichllocation approach. 2015 Presented at: 2015 International Conference on Advances in Computing, Communications and Informatics (ICACCI); August 10-13, 2015; Kochi. [doi: 10.1109/icacci.2015.7275712]

17. Lau J, Newman D, Baldwin T. Machine Reading Tea Leaves: Automatically Evaluating Topic Coherence and Topic Model Quality. 2014 Presented at: 14th Conference of the European Chapter of the Association for Computational Linguistics; 2014; Gothenburg. [doi: 10.3115/v1/e14-1056]

18. Zappavigna M. Searchable Talk: Hashtags and Social Media Metadiscourse. London: Bloomsbury Publishing; 2018.

19. Knapp L, Baum N. Hashtags and How to Use Them on Social Media. J Med Pract Manage 2015;31(2):131-133. [Medline: 26665486]

20. Blei D, Ng A, Jordan M. Latent Dirichlet Allocation. J Mach Learn Res 2003;3:993-1022.

21. Steyvers M, Griffiths T. Probabilistic Topic Models. In: Handbook of Latent Semantic Analysis. Hove, East Sussex: Psychology Press; 2007:439-480.

22. Valdez J. Bias in Public Health Researchthical Implications and Objective Assessment Tools. OAKTrust. 2018. URL: https://oaktrust.library.tamu.edu/handle/1969.1/174067 [accessed 2020-01-02]

23. Funk B, Sadeh-Sharvit S, Fitzsimmons-Craft EE, Trockel MT, Monterubio GE, Goel NJ, et al. A Framework for Applying Natural Language Processing in Digital Health Interventions. J Med Internet Res 2020 Feb 19;22(2):e13855 [FREE Full text] [doi: 10.2196/13855] [Medline: 32130118]

24. Chen H, Boore J. Translation and back-translation in qualitative nursing research: methodological review. J Clin Nurs 2010 Jan;19(1-2):234-239. [doi: 10.1111/j.1365-2702.2009.02896.x] [Medline: 19886874]

25. Dinakar K, Chaney A, Lieberman H, Blei D. Real-time Topic Models for Crisis Counseling. URL: https://affect.media.mit.edu/pdfs/Realtime-Topic-Modeling-Crisis-Counseling.pdf [accessed 2021-11-30]

26. Gethers M, Poshyvanyk D. Using Relational Topic Models to capture coupling among classes in object-oriented software systems. 2010 Presented at: 2010 IEEE International Conference on Software Maintenance; September 12-18, 2010; Timisoara. [doi: 10.1109/icsm.2010.5609687]

27. Potnis D, Tahamtan I. Hashtags for gatekeeping of information on social media. J Assoc Inf Sci Technol 2021 Mar 11;72(10):1234-1246. [doi: 10.1002/asi.24467]

28. Watts C, Hefler M, Freeman B. 'We have a rich heritage and, we believe, a bright future': how transnational tobacco companies are using Twitter to oppose policy and shape their public identity. Tob Control 2019 Mar;28(2):227-232. [doi: 10.1136/tobaccocontrol-2017-054188] [Medline: 29666168]

29. Dai H, Hao J. Online popularity of JUUL and Puff Bars in the USA: 2019-2020. Tob Control 2020 Oct 13. [doi: 10.1136/tobaccocontrol-2020-055727] [Medline: 33051277]

30. Wang W, Chen L, Thirunarayan K, Sheth A. Cursing in English on twitter. 2014 Presented at: CSCW'14: Computer Supported Cooperative Work; February 15-19, 2014; Baltimore, MD p. 415-425. [doi: 10.1145/2531602.2531734]

31. Hong YH, Soh CH, Khan N, Abdullah MMB, Teh BH. Effectiveness of Anti-Smoking Advertising: The Roles of Message and Media. IJBM 2013 Sep 22;8(19). [doi: 10.5539/ijbm.v8n19p55]

32. Social Media Fact Sheet. Pew Research Center. URL: https://www.pewresearch.org/internet/fact-sheet/social-media/ [accessed 2020-01-10]

33. Gore RJ, Diallo S, Padilla J. You Are What You Tweet: Connecting the Geographic Variation in America's Obesity Rate to Twitter Content. PLoS One 2015;10(9):e0133505 [FREE Full text] [doi: 10.1371/journal.pone.0133505] [Medline: 26332588]

34. Padilla JJ, Kavak H, Lynch CJ, Gore RJ, Diallo SY. Temporal and spatiotemporal investigation of tourist attraction visit sentiment on Twitter. PLoS One 2018;13(6):e0198857 [FREE Full text] [doi: 10.1371/journal.pone.0198857] [Medline: 29902270]

35.  Dai Z, Yan C, Wang Z, Wang J, Xia M, Li K, et al. Discriminative analysis of early Alzheimer's disease using multi-modal imaging and multi-level characterization with multi-classifier (M3). Neuroimage 2012 Mar 01;59(3):2187-2195. [doi: 10.1016/j.neuroimage.2011.10.003] [Medline: 22008370]

36.  Klein EG, Czaplicki L, Berman M, Emery S, Schillo B. Visual Attention to the Use of #ad versus #sponsored on e-Cigarette Influencer Posts on Social Media: A Randomized Experiment. J Health Commun 2020 Dec 01;25(12):925-930. [doi: 10.1080/10810730.2020.1849464] [Medline: 33238805]

## Abbreviations

**API:** application programming interface
**FDA:** Food and Drug Administration
**LDA:** Latent Dirichlet Allocation
**NLP:** natural language processing

XSL·FO
**RenderX**

Original Paper

# Infodemic Signal Detection During the COVID-19 Pandemic: Development of a Methodology for Identifying Potential Information Voids in Online Conversations

Tina D Purnat[1*], MSc; Paolo Vacca[2*], BSt; Christine Czerniak[3*], PhD; Sarah Ball[2*], BA; Stefano Burzo[4*], MA; Tim Zecchin[2*], BA; Amy Wright[2*], BA; Supriya Bezbaruah[5*], PhD; Faizza Tanggol[6], BA; Ève Dubé[7*], PhD; Fabienne Labbé[7*], PhD; Maude Dionne[7*], MSc; Jaya Lamichhane[3*], MA, MBA; Avichal Mahajan[3*], PhD; Sylvie Briand[3*], MPH, MD, PhD; Tim Nguyen[3*], MSc

[1]Digital Health and Innovation, Science Division, World Health Organization, Geneva, Switzerland

[2]Media Measurement Ltd, London, United Kingdom

[3]Emergency Preparedness, World Health Organization, Geneva, Switzerland

[4]Department of Political Science, University of British Columbia, Vancouver, BC, Canada

[5]Health Emergencies Programme, World Health Organization Regional Office for South East Asia, New Delhi, India

[6]World Health Organization Country Office Malaysia, Brunei Darussalam and Singapore, Putrajaya, Malaysia

[7]Institut national de santé publique du Québec, Montreal, QC, Canada

[*]these authors contributed equally

**Corresponding Author:**
Christine Czerniak, PhD
Emergency Preparedness
World Health Organization
20 Avenue Appia
Geneva, 1211
Switzerland
Phone: 41 (0)227912111
Email: czerniakc@who.int

## Abstract

**Background:** The COVID-19 pandemic has been accompanied by an *infodemic*: excess information, including false or misleading information, in digital and physical environments during an acute public health event. This infodemic is leading to confusion and risk-taking behaviors that can be harmful to health, as well as to mistrust in health authorities and public health responses. The World Health Organization (WHO) is working to develop tools to provide an evidence-based response to the infodemic, enabling prioritization of health response activities.

**Objective:** In this work, we aimed to develop a practical, structured approach to identify narratives in public online conversations on social media platforms where concerns or confusion exist or where narratives are gaining traction, thus providing actionable data to help the WHO prioritize its response efforts to address the COVID-19 infodemic.

**Methods:** We developed a taxonomy to filter global public conversations in English and French related to COVID-19 on social media into 5 categories with 35 subcategories. The taxonomy and its implementation were validated for retrieval precision and recall, and they were reviewed and adapted as language about the pandemic in online conversations changed over time. The aggregated data for each subcategory were analyzed on a weekly basis by volume, velocity, and presence of questions to detect signals of information voids with potential for confusion or where mis- or disinformation may thrive. A human analyst reviewed and identified potential information voids and sources of confusion, and quantitative data were used to provide insights on emerging narratives, influencers, and public reactions to COVID-19–related topics.

**Results:** A COVID-19 public health social listening taxonomy was developed, validated, and applied to filter relevant content for more focused analysis. A weekly analysis of public online conversations since March 23, 2020, enabled quantification of shifting interests in public health–related topics concerning the pandemic, and the analysis demonstrated recurring voids of verified health information. This approach therefore focuses on the detection of infodemic signals to generate actionable insights to rapidly inform decision-making for a more targeted and adaptive response, including risk communication.

XSL•FO

RenderX

**Conclusions:** This approach has been successfully applied to identify and analyze infodemic signals, particularly information voids, to inform the COVID-19 pandemic response. More broadly, the results have demonstrated the importance of ongoing monitoring and analysis of public online conversations, as information voids frequently recur and narratives shift over time. The approach is being piloted in individual countries and WHO regions to generate localized insights and actions; meanwhile, a pilot of an artificial intelligence–based social listening platform is using this taxonomy to aggregate and compare online conversations across 20 countries. Beyond the COVID-19 pandemic, the taxonomy and methodology may be adapted for fast deployment in future public health events, and they could form the basis of a routine social listening program for health preparedness and response planning.

## Introduction

### Background

Since the beginning of the COVID-19 pandemic, digital communication and social networking have supported the rapid growth of real-time information sharing about the virus that causes COVID-19 (SARS-CoV-2) and the disease in the public domain and across borders. The breadth of conversation, diversity of sources, and polarity of opinions have sometimes resulted in excessive information, including false or misleading information, in digital and physical environments during an acute public health event; this can lead to confusion and risk-taking behaviors that can harm health, trust in health authorities, and the public health response [1]. The excess of information can amplify and protract outbreaks, and it can reduce the effectiveness of pandemic response efforts and interventions.

To address this challenge, the World Health Organization (WHO) Information Network for Epidemics (EPI-WIN), in collaboration with digital research partners, developed a methodology for weekly analysis of digital social media data to identify, categorize, and understand the key concerns expressed in online conversations [2]. The application of this methodology provided the WHO with week-on-week analysis for the prioritization of actions to address online information voids and sources of confusion using verified health information as part of ongoing emergency response planning. When there is a lack of quality information about topics of concern for online users, these topics can be quickly filled with conjecture, low-quality health information, and viral misleading content [3,4], thus potentially causing harm to communities. This approach therefore focuses on the detection of infodemic signals—identifying or predicting rising areas of concern and information voids in the online information ecosystem on a weekly basis to generate actionable insights to rapidly inform decision-making for a more effective response, including adapting risk communication [5].

### Infodemic Management During a Health Emergency

Previous research has explored the use of data produced and consumed on the web to inform public health officials, agencies, and policy—a concept known as *infodemiology* [6]. Initially, the concept of infodemiology aimed to identify the gap between expert knowledge and public practice [7], and it has since evolved to detect and analyze health information on the web through publicly shared search queries, blogs, websites, and social media posts.

The design of interventions for infodemic response must account for an ecosystem where information flow online can cause public health harm offline. Metrics and frameworks related to digital information flows and online behavior are most useful to practitioners when they can be coupled with other online and offline sources of public health data that inform public health decision-making. The WHO has therefore expanded the concept of infodemiology into a multidisciplinary scientific field that amalgamates cross-disciplinary and mixed methods approaches designed to inform the health emergency response [8].

Health emergencies give rise to information overload, which has been shown to influence people's risk perceptions and protective actions during health emergencies [9]. Overload of information of variable quality, timeliness, and relevance is strongly associated with people's experience of information anxiety, which in turn can give rise to information avoidance. Recent examples, from HIV to Ebola virus to Zika virus to polio, have demonstrated the high cost to public health and health systems when misinformation sows distrust, exacerbated by ineffective public health communication and community engagement [3,10]. A lack of active community collaboration in the health response early on deepened distrust, especially as these epidemics unfolded. Currently, most emergency and outbreak recommendations emphasize the value of listening to communities, involving them early in the response, and communicating clearly with them in a timely manner [11,12].

Health authorities therefore not only face the challenge of providing relevant, high-quality health information but also must provide it at the right time, in the right format, and with collaborative engagement of communities [13]. Social listening can help overcome barriers to acceptance of high-quality health information and enactment of healthy behaviors by enabling better understanding of community questions, confusion, information seeking, or intensified attention for given topics. Critical information voids can be identified and characterized in both the online and offline information ecosystems. Our research focuses on the identification and characterization of points of confusion, harmful narratives, and key questions that can reveal information voids in the online social media space

during a health emergency, thereby adding analytical methods to the field of infodemiology that are practical and can directly inform the public health response during a health emergency.

## Analytical Approaches and Metrics To Date

The rise of social media platforms has generated a readily available source of real-time data related to what people express and share in online communities. The 2009 H1N1 influenza pandemic was the first pandemic to occur in the era of social media and was one of the earliest outbreaks informed by analysis of online conversations and information-seeking behaviors. The previous pandemic offers a case study that evidences how online social listening has been used to follow rapidly evolving public sentiment, track actual disease activity, and monitor the emergence of misinformation [14-16]. Although social media platforms have been used to quantify public concerns and sentiment and to monitor real-time pandemic data, they have also been identified as a medium that can enable the spread of low-quality information. For example, within health emergencies, false information has been shown to be posted twice as frequently as evidence-based information, although it is retweeted less frequently [17]. Provision of targeted, relevant, timely, understandable, and resonant health information can therefore benefit from upstream infodemic management activities of public health authorities, including more robust social listening programs.

The onset of the COVID-19 pandemic has exacerbated concerns about misinformation. Throughout the pandemic, there has been a demand for information; at first, this demand was for information about the origin of the virus, and now it is focused mainly on the response to the virus, particularly vaccination and wider public health and social measures. Similar to information voids [5], COVID-19 misinformation trackers have defined the concept of *data deficits* in the online space when there are "high levels of demand for information about a topic, but credible information is in low supply" [3]. The issue with a lack of quality information is that the conversation space can be much more readily filled by misinformation, which may be faster to create and share, more emotive (resonant), and better promoted by content promotion algorithms than factual health information.

Despite the influx of studies as to how information is being spread and shared in the era of COVID-19 and how information is influencing people's health practices, major gaps in knowledge remain as to how best to monitor, understand, and respond to it [8]. Among many possible solutions, social listening, content pretesting, and other computational social science methods have been identified as ways to detect and analyze information voids and viral misinformation narratives [13]. Misinformation research has focused on social media platforms with easier access to data, such as Twitter and YouTube [18,19]; however, misinformation is prevalent across the digital ecosystem (as well as offline). Culture and access to the internet can also affect the nature of misinformation and how it spreads [20]. Beyond identifying what misinformation looks like, studies have also attempted to identify how it emerges [21]*,* aligning with the concept of information voids. Although social listening has tended to focus

on spotting myths and rumors, as well as content items with high engagement and reshare rates, the methodology introduced in this paper expands the scope of social listening and positions it as a core practice of emergency response. This includes prioritizing detection of information voids for more proactive infodemic management before these gaps in understanding are filled with more speculation, misleading information, and counterproductive narratives.

Detecting viral misinformation narratives and information voids in real-time data is crucial to a rapid, comprehensive response by authorities for effective delivery of health information to populations during a health emergency, although this does not ensure that people will necessarily act in accordance with that health information. Previous research has evaluated the correction of misinformation and the role of individuals versus organizations in using real-time data [22-24]. The pathway from receiving information, to intent, to action is understudied and a priority area for future research [8]. Evaluation of intervention impact is challenging [4], but evaluation of interventions must be integrated as part of adaptive infodemic management, including social listening.

Interventions need to address the different aspects of the information ecosystem that influence the spread and health impact of an infodemic. For platforms, content moderation policies, modification of content promotion algorithms, and designing for friction can discourage sharing of misinformation and unverified information [25], while supporting literacies such as health, media, information, digital, and data literacies can promote resilience [26]. The literature highlights the value of a multipronged approach for addressing infodemics at various levels in the digital information ecosystem. However, although public health authorities can influence and interact with the other participants in this space, there is a need to suggest immediate and practical tools that public health authorities can deploy within their mandate in a health emergency context in support of their health operations and communication activities [8].

## A Need for Practical Tools for Health Authorities

Research is ongoing to assist policy makers in understanding public concerns and sentiment around the pandemic as well as in tracking information outbreaks and the emergence of misinformation. However, there is little to no empirical evidence on how this research can be used to develop practical tools for an outbreak response by public health authorities. More collaborations between researchers and public health practitioners are needed to fill this gap. As a contribution to the infodemic response toolbox, the taxonomy and methodology in this study offer a practical, structured approach for identifying information voids and narratives of concern that warrant attention and action. This approach has already provided actionable data to help the WHO focus its efforts for the COVID-19 pandemic response.
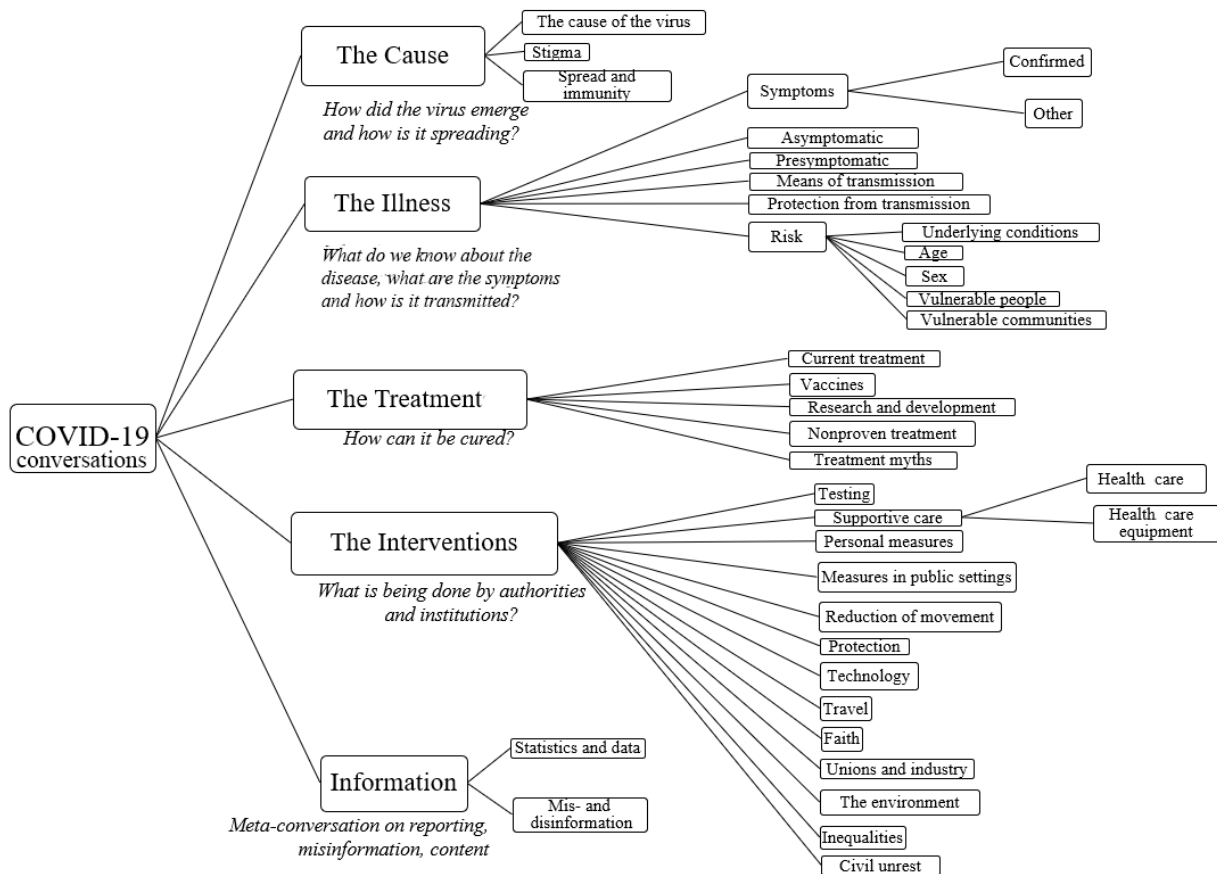
XSL•FO

**RenderX**

## Methods

### Development of a Public Health Taxonomy for Social Listening

A social listening taxonomy for COVID-19 conversations was developed specifically for this analysis. It was designed to filter digital content referring to COVID-19 (and synonyms) for items of relevance in a public health context and to classify that content into categories. The taxonomy consisted of 35 keyword-based searches (one set of searches for each of two languages, namely English and French) which were grouped into 5 overarching topic categories representing thematic areas in which people were engaging, writing, or searching for information.

The 5 top-level categories and corresponding 35 subcategories of this social listening taxonomy for COVID-19 conversations

were defined based on established epidemic management and public health practices during an outbreak of infectious disease [27] (Figure 1). The first 4 categories refer to the focus of epidemic management activities during the pandemic: (1) the cause of the disease—what do we know about the virus, and how is it spreading? (2) the illness—what are the symptoms, and how is it transmitted? (3) the treatment—how can it be cured? and (4) the interventions—what is being done by authorities and institutions? In addition, a fifth category was included to examine public perceptions on circulating information (ie, metaconversations about evidence and statistics, mis- and disinformation, successful and harmful content, or key influencers who have been actively amplifying information on COVID-19). This category was designed because misinformation, rumors, and polarization of factual versus misleading narratives are common challenges in epidemic management.

**Figure 1.** Structure of the social listening taxonomy for COVID-19 conversations.



Each of these 5 categories were segmented into subcategory levels that are familiar to the epidemiologist's investigation and management of the outbreak, resulting in a total of 35 taxonomy subcategory levels (Figure 1). For example, the taxonomy category about the illness was further defined by subcategories to identify conversations, questions, or confusion about the symptoms of the illness, how it transmits, and what populations may be affected by it (across demographics, vulnerable populations, and people with underlying conditions). By defining the social listening taxonomy across the investigation areas of epidemic management [27], the resulting infodemic insights can be more quickly evaluated by public health professionals

and turned into actionable recommendations to inform the epidemic response.

Each of the 35 taxonomy subcategories encompassed a list of topics that captured different aspects of that segment of the online conversation on COVID-19. Keywords for the 35 subcategory searches were generated based on expert knowledge from the WHO EPI-WIN team and translated into Boolean search strings to identify topic-related language for review of relevant social media posts and news content. The keywords of the taxonomy are available on request at the contact address listed in the *Acknowledgments* section.

In addition to the taxonomy subcategory levels, the keyword-based Boolean search string was created to also identify posts containing a question; this enabled analysis of categories for which people were seeking information and, therefore, potential information voids. The question search string was designed to be paired with each of the 35 taxonomy Boolean search strings to identify posts referring to the topic and containing question words, verb-subject inversions, and auxiliaries.

Finally, the sum of the total volume of the social media conversation (on all topics) was estimated by monitoring the number of posts mentioning at least one of the most commonly used words in English (eg, *the*, *and*, *or*, *I*) and French (eg, *le*, *la*, *ou*, *et*). The data were collected via a Boolean search string comprised of these most commonly used words.

## Data Sources and Data Collection

The analysis was based on the weekly aggregation of publicly available social media data in English and French using Meltwater Explore. Institutional Review Board review was not sought, as the analysis used large bodies of text written by humans on the internet and on some social media platforms. The analysis and resulting reports focused on the identification of conversation narratives and thematic questions instead of on individual statements and users.

The Meltwater social listening platform was configured to collect verbatim mentions of keywords associated with the 35 predefined taxonomy category Boolean searches from 9 open data sources and fora (Twitter, blog entries, Facebook, Reddit posts and comments, other unspecified message boards or fora, comments under news articles and blog entries, Instagram posts, product reviews, and YouTube video titles and comments). A total of 87.02% of the resulting analysis data set was sourced from Twitter. Blogs (5.34%) and, specifically, the Reddit platform (4.34%) were the next most prominent sources in the data set. These were followed by message boards (2.14%), comments under news articles (0.89%), online review websites (0.13%), Instagram (0.12%), and Facebook (0.03%).

For each of the 35 taxonomy subcategories, the global daily total volume of posts, and the volume of posts posing a question, were recorded on a weekly basis. Tracking changes in volume from week to week also enabled determination of the velocity for a given subcategory.

## Testing and Validation of the Taxonomy

The methodology used to test and validate the retrieval and classification in this study used both retrieval precision and retrieval recall, which are related to how much retrieved data is relevant and how much relevant data is retrieved, respectively [28,29]. These validation metrics are useful for assessing the performance of machine learning models in information retrieval and have been used for metrics on content retrieved and classified via Boolean searches for news media content [28] and Twitter data [29].

To test whether the taxonomy categories captured the intended information (retrieval precision), a random sample of content captured by each of the 35 Boolean searches was human-coded for relevance (10,500 posts in total) by a single reviewer, with a second reviewer validating the coding. The post was coded as either relevant to the search subcategory (1) or not relevant (0).

The aim of the coding was to determine the proportion of relevant (R; also, "true positive" [TP]) results as a percentage of the retrieved sample. The coders judged whether a post was relevant according to the intended definition of the specific subcategory search for which the post was returned. For example, if a post had been returned for the "The Illness – Confirmed Symptoms" search, the coder would check if the post referred to a confirmed symptom of COVID-19 (TP) or whether the matched keywords were mentioned in a different context (irrelevant [I]; also, false positive [FP]). For instance, if a post had been returned by the Boolean search for COVID-19 vaccines, did the post refer to COVID-19 vaccinations? If yes, the post was coded as a TP. If the post in question mentioned COVID-19, but the part of the post mentioning vaccines was about the influenza vaccine, the post would be coded as an FP.

The initial retrieval precision testing showed an average result of 82% for the 35 taxonomy subcategory searches. The retrieval precision rate was calculated as precision = [TP ÷ (TP + FP)] × 100%.

A total of 7 searches returned content below the target minimum retrieval precision rate of 70%, with a range of 42% to 100% (Table 1). To reduce the rate of false positives, the keywords for the 7 searches that performed below the target minimum rate were subsequently reviewed and updated to exclude keywords and phrases returning irrelevant content. On retesting, the average retrieval precision rate for the 35 searches was 87%, with a range of 72% to 100%. The full results of the retrieval precision testing and subsequent retesting can be seen in Table 1.

To spot-check the coding for reliability, we deployed a second reviewer to analyze 10% of the posts (30 per taxonomy category search, 1500 in total). We calculated the Cohen kappa to determine intercoder reliability, which was found to be high ($\kappa$=0.81, observed agreement [$p_o$]=0.95, expected agreement [$p_e$]=0.76).

A further test was performed to assess retrieval recall: whether content of relevance to the research aims failed to be retrieved by the taxonomy searches. To test this, a random sample of 1000 items of content, mentioning COVID-19 (and synonyms) but excluding the taxonomy category keywords (the "not retrieved" sample in Table 2), was human-coded for relevance from a public health perspective. Posts in this sample were determined by the coder to be relevant (R) to the aims of the public health research (false negative [FN]), or irrelevant (I) to the research aims (true negative [TN]). Coding was performed by the same reviewer and was binary; content was irrelevant (I, and therefore also TN) or was relevant (R, and therefore also FN) and deemed to have been missed in taxonomy category searches.

**Table 1.** Results of retrieval precision testing and retesting with a sample size of 300 posts analyzed per subcategory.

| Subcategory | Posts retrieved by the taxonomy category search human-coded as true positives and retrieval precision rate, n (%) |
| --- | --- |
| The Cause – The Cause | 217 (72.3) |
| The Cause – Further Spread – Stigma | 260 (86.7) |
| The Cause – Further Spread – Immunity | 245 (81.7) |
| The Illness – Confirmed Symptoms | 189[a] (63)/239[b] (79.7) |
| The Illness – Other Discussed Symptoms | 141[a] (47)/218[b] (72.7) |
| The Illness – Asymptomatic | 300 (100) |
| The Illness – Presymptomatic | 300 (100) |
| The Illness – Means of Transmission | 295 (98.3) |
| The Illness – Protection From Transmission | 299 (99.7) |
| The Illness – Underlying Conditions | 238 (79.3) |
| The Illness – Demographics – Sex | 215 (71.7) |
| The Illness – Demographics – Age | 215 (71.7) |
| The Illness – Vulnerable People | 287 (95.7) |
| The Illness – Vulnerable Communities | 269 (89.7) |
| Treatment – Vaccines | 300 (100) |
| Treatment – Current Treatment | 144[a] (48)/224[b] (74.7%) |
| Treatment – Research & Development | 290 (96.7) |
| Treatment – Nonproven Treatment (Nutrition) | 245 (81.7) |
| Treatment – Myths | 126[a] (42)/221[b] (73.7) |
| Interventions – Measures in Public Settings | 243 (81) |
| Interventions – Testing | 280 (93.3) |
| Interventions – Supportive Care – Equipment | 204[a] (68)/257[b] (85.7%) |
| Interventions – Supportive Care – Health Care | 289 (96.3) |
| Interventions – Personal Measures | 298 (99.3) |
| Interventions – Reduction of Movement | 256 (85.3) |
| Interventions – Protection | 276 (92) |
| Interventions – Technology | 278 (92.7) |
| Interventions – Travel | 250 (83.3) |
| Interventions – Faith | 201[a] (67)/269[b] (89.7) |
| Interventions – Unions and Industry | 223 (74.3) |
| Interventions – The Environment | 183[a] (61)/290[b] (96.7) |
| Interventions – Inequalities | 280 (93.3) |
| Interventions – Civil Unrest | 280 (93.3) |
| Information – Misinformation | 273 (91) |
| Information – Statistics | 244 (81.3) |

[a]Indicates a taxonomy subcategory search that performed below minimum requirements and was subsequently updated and retested to yield better performance.

[b]Number and percentage of posts in the sample coded as true positives in the retesting of the taxonomy subcategory search following the update.

**Table 2.** Results of human coding of retrieved and unretrieved samples for calculation of retrieval recall and F-scores.

| Sample | Coded relevant | | Coded irrelevant | | Total coded sample, n |
|---|---|---|---|---|---|
| | Samples, n | Description | Samples, n | Description | |
| Retrieved[a] | 875 | True positive | 125 | False positive | 1000a |
| Not retrieved | 304 | False negative | 696 | True negative | 1000 |
| Total | 1179 | N/A[b] | 821 | N/A | 2000 |

[a]The "retrieved" sample size was downweighted to equal the "not retrieved" sample size.

[b]N/A: not applicable.

The results of the coding of the "not retrieved" sample indicated the proportion of TN results as a proportion of the sample; 70% of content was judged not to be relevant to the research aims, and therefore it was deemed correct that this content was not retrieved by our taxonomy. From the data, we also calculated the retrieval recall rate as recall = [TP ÷ (TP + FN)] × 100%.

The overall retrieval recall rate was 74%. This coding process enabled identification of areas where the existing Boolean string could be expanded to include more relevant keywords to retrieve more relevant content, or where the taxonomy could be expanded to include new and emerging issues. From the content that was not retrieved but was judged to be of potential relevance to the research aims (false negatives, FNs), 3 topics were identified that will be added to the taxonomy in a pending update: mutations/variants of the COVID-19 virus; "long covid" (long-term symptoms of COVID-19); and the impact of the pandemic on mental health and well-being.

To validate the coding of the sample of "not retrieved" content for reliability, we deployed a second reviewer to analyze 10% of the posts (100 posts). We calculated the Cohen kappa to determine intercoder reliability, which was found to be high ($\kappa=0.86$, $p_o=0.93$, $p_e=0.50$).

From the results of the coding of the retrieved and unretrieved samples, we calculated an F1 score and an F0.5 score with the following formulas: $F_1 = [(2 \times \text{precision} \times \text{recall}) \div (\text{precision} + \text{recall})]$, and $F_{0.5} = [(1.25 \times \text{precision} \times \text{recall}) \div (0.25 \times \text{precision} + \text{recall})]$.

The F1 score (harmonic mean of precision and recall) for the searches was 0.80, and the F0.5 score was 0.84. F1 and F0.5 scores range from 0 to 1, with 1 representing perfect performance. A higher F1 or F0.5 score is considered reasonable, with a score closer to 1 indicating stronger performance of a retrieval and classification approach. The inclusion of the F0.5 measure reflects the greater importance of retrieval precision in this study: given the vast number of potentially relevant pieces of content, it is more important to the aims of this project to correctly classify the retrieved posts than to collect every possibly relevant post. Therefore, we consider it a positive result to achieve a higher F0.5 score than F1 score. This is because in this study, it is more important that the results are not impacted by a high number of false positives and that the true positives are classified into the correct subcategory. The retrieval recall testing is also helpful because it enables identification of new or changing pandemic issues,

such as new terminology being used that can be added to the taxonomy category search language over time.

## Quantitative Data Analysis

Potential information voids were identified based on 3 parameters within the weekly data set: the volume (ie, how many social media items referred to topic X?), the velocity (ie, the rate of increase of the number of social media items that have engaged with topic X over the course of the past week), and the presence of questions about the topic. The volume was the sum of the online items that mentioned COVID-19 together with a keyword related to each tracked topic. Velocity was determined as the percentage increase of the volume of content items aggregated under each topic from week to week, where velocity = [(current week's total number of mentions – previous week's total number of mentions) ÷ (previous week's total number of mentions)] × 100%.

Starting in late March 2020, weekly global analysis reports were produced that supplied the EPI-WIN team with early warnings of points of concern expressed in public comments by online users [2,4]. By May 4, 2021, the data sample consisted of a sum of 1.02 billion unique social media posts. This was a subset of the larger pool of 1.3 billion total public social media posts in English and French mentioning COVID-19 gathered by the data aggregator. The sample of 1.02 billion posts consisted of approximately 3% of the pool of all public social media posts written in English and French that had been gathered by the data aggregator since March 2020. The data set of total public social media posts gathered by the aggregator was verified through the automated search of mention of the most common words in English and French (eg, *the*, *le*, *and*, *et*).

Each week, social media conversations were segmented based on levels of velocity and quantitatively examined for public engagement (eg, likes, shares, poll votes, reactions), hashtags, and most-used keywords and phrases. From this weekly quantitative data, up to 10 topics with high velocity and/or a large proportion of social media posts expressing a question, and/or with high levels of engagement, were identified as potential priority information voids or sources of confusion or concern.

The identified issues on social media were then further evaluated using engagement data and Google search trends to determine whether a significant number of online users had also been looking for information on these topics to help determine whether the information void was more widespread.

## Qualitative Analysis

Each week, we used the quantitative analysis to identify up to 10 topics reflecting potential information voids and areas of concern. These topics were then examined in more detail via qualitative analysis to understand the context and identify where action may need to be taken in line with a sequential explanatory design approach [30]. The qualitative analysis involved ad hoc human-led review of the key narratives, influencers, and public reactions as reflected in the content.

This analysis prioritized the flagging of widespread confusion or frequently asked questions, the rapid amplification of misinformation, or ad hoc aspects of the conversation that were particularly relevant to public health, such as vaccine questioning ahead of and during a vaccination campaign.

## Reporting

The quantitative data were compiled in a web-based dashboard accessible to the emergency responders in the EPI-WIN team, and insights were discussed with EPI-WIN emergency responders on a weekly basis. The dashboard was updated weekly to allow investigation of short- and long-term trends in volumes, changes in velocity, and the volume of questions for each topic.

Weekly written reports outlined quantitative and qualitative findings about the 5 to 10 topics of concern, included visualizations from the dashboard, and summarized recommendations for action when needed [31].

## Results

Quantitative analysis of the volume changes indicated that the narratives and questions in the online conversations shifted as the pandemic evolved over the course of 2020 and into 2021 (Table 3). Based on the average weekly rises of the topics within each of the 5 taxonomy categories in the yearly quarters between March 23, 2020, and March 31, 2021, it was observed that the second quarter (Q2) and third quarter (Q3) of 2020 were characterized by a steady increase in conversations about "the interventions." Although discussion of "the illness" decreased in 2020, it surged again in the first quarter (Q1) of 2021. In the fourth quarter (Q4) of 2020, "the treatment" had the highest velocity in digital conversations, while the metaconversation on COVID-19 information experienced the greatest velocity in Q1 of 2021.

**Table 3.** Most discussed topics by month and results of the pivoted data set by month, sorted by volume of social media mentions.

| Year and month | Most discussed topic | Volume (millions of social media mentions) |
| --- | --- | --- |
| **2020** | | |
| March | Interventions – Testing | 37 |
| April | Interventions – Testing | 18 |
| May | Interventions – Measures in Public Settings | 12 |
| June | Interventions – Testing | 11 |
| July | Interventions – Testing | 14 |
| August | Interventions – Testing | 9 |
| September | Interventions – Testing | 8 |
| October | Interventions – Testing | 17 |
| November | Interventions – Testing | 8 |
| December | Treatment – Vaccines | 15 |
| **2021** | | |
| January | Treatment – Vaccines | 15 |
| February | Treatment – Vaccines | 12 |
| March | Treatment – Vaccines | 15 |
| April | Treatment – Vaccines | 15 |

At the same time, topics re-emerged periodically in terms of popularity. All 35 categories of topics that were tracked resumed a higher velocity throughout the reporting period for an average of 18 weeks combined (Table 4). The 2 topics that attracted increasing interest most frequently were "myths" and "risk based on age demographics" (rising for 26 and 24 weeks, respectively) followed by "the cause" of the virus and "reduction of movement" (both 23 weeks) "vaccines" and "stigma" (both 22 weeks), and "other discussed symptoms" (21 weeks). Digital conversations on "the cause" of the epidemic, "misinformation" as a phenomenon, and "immunity" had the longest continuous periods of surge in volume of social media posts discussing these topics in the context of the COVID-19 pandemic; the conversation on "the cause" increased in both the first and second half of the analysis period for 7 continuous weeks during the first half of the reporting period, while the metaconversation about misinformation increased for 6 consecutive weeks. Conversations about "immunity" increased for 5 consecutive weeks in June-July 2020 and in November-December 2020.

**Table 4.** Frequency of weekly velocity growth (number of weeks in which a topic experienced positive velocity) and average weekly increase rate (or decrease, when a negative value is returned) by topic.

| Topic | Number of weeks in which a topic experienced positive velocity (increase of social media mentions since previous week) | Average weekly increase in number of social media mentions (%) |
| --- | --- | --- |
| Treatment – Myths | 26 | 31 |
| The Illness – Demographics – Age | 24 | 8 |
| Interventions – Reduction of Movement | 23 | 14 |
| The Cause – The Cause | 23 | 7 |
| Vaccines | 22 | 7 |
| The Cause – Further Spread: Stigma | 22 | 10 |
| Interventions – Faith | 21 | 3 |
| The Illness – Other Discussed Symptoms | 21 | 52 |
| Treatment – Current Treatment | 20 | 3 |
| Interventions – Travel | 20 | 3 |
| Interventions – The Environment | 20 | 7 |
| The Illness – Confirmed Symptoms | 20 | 1 |
| The Illness – Asymptomatic transmission | 19 | 6 |
| The Illness – Means of Transmission | 19 | 10 |
| Interventions – Measures in Public Settings | 18 | 4 |
| The Cause – Further Spread: Immunity | 18 | 9 |
| Treatment – Nonproven Treatment (Nutrition) | 18 | 4 |
| Information – Statistics and Data | 18 | 4 |
| Interventions – Technology | 18 | 5 |
| Information – Misinformation | 17 | 5 |
| The Illness – Vulnerable Communities | 17 | 4 |
| Interventions – Testing | 17 | -1 |
| Interventions – Supportive Care – Health Caren | 17 | 3 |
| Interventions – Protection | 17 | -3 |
| The Illness – Presymptomaticn | 16 | 40 |
| The Illness – Underlying Conditionsn | 16 | 10 |
| Interventions – Supportive Care – Equipment | 16 | 0 |
| The Illness – Protection From Transmission | 16 | 0 |
| Interventions – Personal Measures | 15 | -3 |
| The Illness – Demographics – Sex | 15 | 8 |
| Treatment – Research & Development | 15 | 8 |
| Interventions – Unions and Industry | 15 | -1 |
| Interventions – Inequalitiesn | 14 | -1 |
| The Illness – Vulnerable Peoplen | 14 | 1 |
| Interventions – Civil Unrest | 14 | 32 |

Analysis of the peaks in discussion of 2 of the leading recurring topics, "risk related to age demographics" and "the cause," provided insight into how narratives around these topics were fueled by real-life events. The conversation on "risk related to age demographics" increased in velocity 24 times throughout the period studied. A total of 3 million public social media posts

engaged with the topic: 64% of these posts were focused on children, whereas 30% focused on older people. The volume of conversation on children and COVID-19 risk increased above the yearly average for 133 days. Speculation about the severity of COVID-19 infection in children was raised consistently throughout the evaluation period, and it represented fertile

XSL·FO

RenderX

ground for confusion and potential misinformation. Major triggers included news reports of child deaths (560,000 public posts discussed children and mortality), reports of symptoms observed in children in particular (300,000 public posts discussed children and symptoms), the debate over school reopenings, particularly with regard to transmission (818,000 public posts) and, most recently, COVID-19 immunization (656,000 public posts). In relation to this topic, doubts resurfaced repeatedly about the threat of COVID-19 to children; however, there was a diversity of narrative foci for these doubts, linked to changing events during the pandemic.

By contrast, public discussion on the possible origins of the pandemic ("the cause") had a persistent narrative throughout the evaluation period. "The cause" of the epidemic was a focus of 3.26 million public social media posts throughout the period monitored. The size of the conversation was most prominent at the beginning of the pandemic and diminished as of June 2020, but with periodically recurring peaks in the number of posts. Conspiracy theories suggesting the artificial origin of the virus as a bioweapon were persistent in online discussion, and prominent influencers operating in the conspiracy theory space were often linked to resurgent peaks in public online discussion. The phrase "biological weapon" was mentioned 326,000 times in the public social media space (cf. 141 million mentions of COVID-19 vaccines in the same period). The rate of mentions decreased by 65% from Q2 to Q3 2020 (as it decreased to 34,000 mentions globally), but it surged to 110,000 in Q4 as the theory regained prominence in the public discourse, in part driven by the release of a preprint paper claiming that the virus was an "unrestricted bioweapon" [31,32]. In Q4 2020, 16% of posts referring to the virus as a bioweapon referenced the authors of that paper. Although the nature of the narrative around COVID-19 as a "bioweapon" was relatively constant, our findings indicate that existing conspiracy theories can be fueled with new details in debates about science [32], underscoring a need to improve science literacy and communication.

## Discussion

### Principal Findings

The insights obtained in this study have afforded public health experts the opportunity for a more rapid and targeted assessment of a subsample of narratives across the English and French languages using public digital sources. These insights can be combined with others to better understand whether and how people are understanding public health and social measures and putting them into practice to protect themselves and their communities.

The application of this taxonomy to successive weekly online social listening analysis has resulted in a better understanding of the evolution and dynamics of high-velocity conversations about COVID-19 worldwide in English and French during the pandemic. The taxonomy also provides a quantifiable approach to support more adaptive and targeted planning and prioritization of health response activities. For example, monitoring and characterizing re-emerging topics can guide re-evaluation and updating of risk communication and community engagement initiatives to improve understandability and resonance, or

highlight where adjustments in technical guidance, public health policy, and social measures may be needed. In addition, the fact that narratives discussed online often overlap across different categories reveals the breadth of this taxonomy, and this overlap enables emerging narratives and potential information voids to be picked up through velocity alerts raised in different elements of the taxonomy.

The testing process described in this article forms the basis of the taxonomy review and maintenance process. Updates to the taxonomy are also informed by observations from the weekly analysis and reporting of the data, and public health expert knowledge via WHO, the wider news agenda, and epidemic management context of the pandemic. The taxonomy has been updated twice since its creation in March 2020, with a third update forthcoming in 2021. The aim of the taxonomy updates is to ensure that important new and emerging topics are captured as the pandemic evolves (as in the examples of variants/mutations and "long covid" above) and that the taxonomy includes the latest language and terms being used by the public [29]. For example, as the pandemic progressed, members of the public increasingly dropped the use of formal terms, such as referring to the virus as "Covid" rather than "COVID-19"; therefore, the taxonomy keywords were expanded to reflect this change. When the taxonomy was updated and validated, the database was also updated back to the start date of the research to ensure consistency in the analysis data set and to allow for analysis of long-term trends.

There is added value in using a common social listening taxonomy for integration of insights from a variety of data sources and research methods in online and offline communities. This can provide a more systematic way to integrate analysis of different data sources and facilitate complementarity of digital social listening data with other data such as knowledge, attitudes, and practices research to help uncover drivers of online discussion, and to support social listening in vulnerable or more marginalized communities, including those with limited access to online platforms.

A challenge of this analysis approach is the need for human analysts to continuously monitor and evolve the taxonomy in line with the developing narratives and emerging topics as well as the changing language used in discussions of the COVID-19 pandemic. Ideally, taxonomies would be tested, reviewed, and updated frequently, particularly when a new stage of the pandemic begins (eg, when the vaccine rollout started), as such events in the pandemic timeline can generate new topics of discussion and new terminology (eg, "Covid passports"). However, the benefit of more frequent updates is balanced by the need for comparability of data across time as well as by the fact that this analytical method needs to be rapidly reproducible, including in more resource-constrained environments, to have real, practical use week-on-week during the pandemic to inform the immediate needs of the health authority response, including risk communication and community engagement in any country context.

To help identify actionable insights, the weekly analysis was focused less on exact counts of mentions and more on relative changes, narratives, and topic signals to evaluate and

contextualize infodemic signals. When rapidly identifying up to 10 information voids in large weekly data sets, absolute precision was less important than the early detection of an actionable signal to help trigger a timely response. For example, if there was a sudden rise in online narratives expressing concern over a treatment, coupled with other information available from the emergency response, the exact number of mentions was less important than signal detection, analysis, and recommendations for possible action. Despite this, more research is needed to refine and streamline the process for rapidly updating and publishing such taxonomies, especially in protracted epidemics, where shifts in concerns and conversations are bound to occur.

A key takeaway from the analysis that can be applied during the current pandemic is the frequent recurrence of topics of concern and its implications for communication. Public health authorities, governments, and nongovernmental organizations must be prepared to communicate repeatedly on the same issue, adapting frames, approaches, and content as public perceptions of issues and topics shift. Our analysis shows that areas of concern wax and wane, with confusion disappearing and re-emerging as new information comes to light or new events occur. Monitoring the changing narratives on a weekly basis and over time using a taxonomy, such as the one used in this study, can enable health authorities to assess longer-term trends and to be more nimble in adapting approaches to respond effectively to topics of concern and to counter misinformation. Further research can help to adapt these digital social listening approaches to provide metrics for evaluation of infodemic management interventions.

The taxonomy has been adapted, translated, and applied in a number of country-level studies in Mali, the Philippines, and Malaysia [33-35]. Applying the approach at the country level included the localization of keywords and their validation. Once this work was completed, the taxonomy and methodological approach proved to be a useful tool for generating insights into narratives in public discourse and potential information voids at the national level. Furthermore, the research framework is now being applied in Canada by the National Institute of Public Health of Québec as an input into the public health response and risk communication in that province, showing that the taxonomy is also applicable at the subnational level [36]; Institut national de santé publique du Québec [forthcoming].

A pilot project by WHO EPI-WIN and research partners, Early Artificial intelligence–supported Response with Social Listening (EARS) [37], also built on the taxonomy from this research and applied it to an automated classification of content and analysis of publicly shared opinions and concerns in 20 countries. The EARS project is enabling both country-level analysis and cross-country comparisons of themes in online conversations, although obtaining in-depth contextual insights still requires human-led analysis of potential information voids and sources of confusion. Therefore, more investment in analytical capacities in social listening at the country level is needed to provide more contextual analysis, interpretation of infodemic insights, and formulation of recommendations for action, as well as to build capacities for using social listening for health response evaluation and adaptation.

There is an opportunity to apply the taxonomy and methodology described in this paper to detect information voids during future, as yet unknown, pandemics and other public health crises. The 5 top-level categories and some of the 35 sub-categories are relevant to social listening in any outbreak but would need to be adapted to the type of pathogen. If, for example, the HIV/AIDS epidemic had started in the digital, connected world of 2020 rather than in the 1980s, the online social listening taxonomy structure would have needed some adjustment to filter and segment public discourse related to the epidemic and identify information voids. For example, a "Demographics – Men Who Have Sex With Men" topic could be added under the category "the illness" to better hear questions and concerns from this particular demographic group. This approach could also include adjustments to subcategories under "the intervention" to remove irrelevant subcategories of "Reduction of Movement" and "Unions and Industry." After such a taxonomy review and adjustment, the keywords used to capture content related to each category and subcategory would also need to be systematically reviewed to ensure they were appropriate to the narratives in relation to specific illness in question. For example, terms relating to injected drug use, sex between men, sex between a man and woman, and mother-to-child transmission could be added under "The Illness – Modes of Transmission". Having a taxonomy structure and methodology already in place as a starting point would enable faster deployment of digital social listening activities in a future outbreak.

## Limitations

Interpretation of the analysis must account for the limitations of the data sources included in the content aggregator. During health emergencies, health authorities require surge support in social listening, response, and evaluation functions. Analysis services from a central analytics unit or from commercial or academic institutions need to be set up quickly to use a systematic approach to detect and understand people's changing concerns, questions, and possible areas of confusion shared publicly online. The overhead in management of data from open sources can be high, and in settings where the social listening analytics capacity is not yet in place for routine analysis, content aggregators can be used to rapidly set up an analysis workflow. The media content aggregation platform used for this study offers firehose access to Twitter, ensuring a complete set of data for analysis, subject to privacy limitations. Other sources in the platform are either sampled from or limited to public posts only [38]. This is a limitation that applies to most analytics of this type, as Facebook and other social media platforms set limitations on the data they make available due to their privacy policies and commercial interests. As a result, there is an overrepresentation of Twitter content in this analysis [39,40]. The use of private data aggregators may lead to the use of unconventional, uncontrolled samples whose breadth and comprehensiveness are constrained by practical and legal limitations. Other methods would be required to characterize conversations in hidden online communities, closed groups, and closed messaging apps, and thorough consideration of the ethics of social listening would be warranted in such contexts.

This research is global and is limited to two major languages (English and French). As a result, only major online narrative

themes and information voids were identified, and the resulting interpretations may not be representative of trends and patterns that could be observed in digital communities for other languages. Moreover, in a global weekly analysis, smaller or more localized conversations may go undetected. One of the aims of this work is to apply and advance the methods to develop taxonomies that can be rapidly applied to any linguistic context for different geographies and public health events.

It has also been observed that the global English-language data set is prone to overrepresent the voice of social media in geographic regions or communities that are more digitally active than others. A key challenge in this study was the digital amplification of discourse pertaining to US politics, the elections, and the digital prominence of US civil society thereof [41]. In such situations, exclusion keywords may be used to exclude major events or large-scale media coverage from analysis so that they do not mask citizens' publicly shared narratives that are more relevant for public health authorities. This can also be addressed when presenting the analysis results. For example, the weekly reports presented analysis of the narratives from the United States and the United Kingdom separately from the analysis of data from other countries where English was the language of online conversation. This helped to uncover previously undetected narratives outside the United States and the United Kingdom. Future research is needed to assess how results may vary in different linguistic communities and to evaluate the effects of geographies that may be superinfluencers of global discourse.

Another limitation of this research is the start date of the project, March 23, 2020, which is several weeks after COVID-19 was declared a public health emergency of international concern; however, data prior to this date (back to January 2020) have been retrieved and stored for future analysis, ensuring that it is possible to analyze a longer timeline.

Adaptation and application of the taxonomy structure in future outbreaks must also take into account validation of information retrieval and recall. The test scores referenced in the taxonomy testing and validation section should be taken as estimates of the accuracy of the retrieval process by the taxonomy category searches, and function most effectively as a tool for identifying areas for improvement. A key limitation of the test results is that human coders can make errors [29]. The human coders involved in the testing and validation were highly experienced in coding and highly familiar with the topic in question, which can help minimize the incidence of coding errors. Future applications of this validation approach could also deploy more coders in an effort to remove potential bias introduced by reliance on a small number of coders.

## Conclusions

This research focuses on the identification of potential information voids and sources of confusion in online social conversations to provide actionable insights for risk communication and community engagement and other health response activities. While it can provide insight into the opinions expressed online, integration with other analyses, including from listening to offline communities is needed. Applying this methodology globally has provided the added and needed insight, inspiring new ways of thinking and use of information in support of risk communication during health emergencies. Much of the value of the taxonomy we developed is in the capacity to rapidly deploy and provide ongoing insights about information voids during an outbreak, which then allows a health authority to take evidence-informed action and course-correct risk communication during an epidemic. The application of the taxonomy and methodology for social listening at regional, country, and subnational levels in the COVID-19 pandemic—which is already being tested—offers possibilities for more actionable insights that must increasingly support a localized response. Moreover, this method offers an approach for monitoring of concerns, questions, and information voids in future outbreaks, enabling a faster response by the health authorities in affected countries during the next acute health event.

## Data Availability

The listing of keywords and search terms per taxonomy subcategory is available upon request by contacting enquiry@mediameasurement.com.

## Conflicts of Interest

S Ball, PV, AW, and TZ are employed by a media monitoring company that provides a wide range of services to clients in media monitoring and listening, including the WHO. The work described in this paper was part of the contractual service to the WHO. The other authors have no conflicts to declare.

## References

1. 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]

2. Coronavirus disease 2019 (COVID-19) Situation Report 100: 29 April 2020. World Health Organization. 2020 Apr 29. URL: https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200429-sitrep-100-covid-19.pdf?sfvrsn=bbfbf3d1_6 [accessed 2021-07-16]

3. Shane T, Noel P. Data deficits: why we need to monitor the demand and supply of information in real time. First Draft News. 2020 Sep 28. URL: https://firstdraftnews.org/long-form-article/data-deficits/ [accessed 2021-07-16]

4. Coronavirus disease 2019 (COVID-19) Situation Report 128: 27 May 2020. World Health Organization. 2020 May 27. URL: https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200527-covid-19-sitrep-128.pdf?sfvrsn=11720c0a_2 [accessed 2021-07-16]

5. Zecchin T. WHO ad-hoc online consultation on managing the COVID-19 infodemic. World Health Organization. URL: https://www.who.int/teams/risk-communication/infodemic-management/who-ad-hoc-online-consultation-on-managing-the-covid-19-infodemic [accessed 2021-07-16]

6. Eysenbach G. Infodemiology and infoveillance: framework for an emerging set of public health informatics methods to analyze search, communication and publication behavior on the Internet. J Med Internet Res 2009 Mar 27;11(1):e11 [FREE Full text] [doi: 10.2196/jmir.1157] [Medline: 19329408]

7. Eysenbach G. Infodemiology: the epidemiology of (mis)information. The American Journal of Medicine 2002 Dec;113(9):763-765. [doi: 10.1016/S0002-9343(02)01473-0]

8. WHO public health research agenda for managing infodemics. World Health Organization. 2021 Feb 03. URL: https://www.who.int/publications/i/item/9789240019508 [accessed 2021-07-16]

9. Soroya SH, Farooq A, Mahmood K, Isoaho J, Zara S. From information seeking to information avoidance: understanding the health information behavior during a global health crisis. Inf Process Manag 2021 Mar;58(2):102440 [FREE Full text] [doi: 10.1016/j.ipm.2020.102440] [Medline: 33281273]

10. Vicol D, Tannous N, Belesiotis P, Tchakerian N. Health misinformation in Africa, Latin America and the UK: Impacts and possible solutions. Full Fact. 2020 May. URL: https://fullfact.org/media/uploads/en-tackling-health-misinfo.pdf [accessed 2021-07-16]

11. Risk communication and community engagement readiness and response to coronavirus disease (COVID-19): interim guidance, 19 March 2020. World Health Organization. 2020 Mar 19. URL: https://www.who.int/publications-detail/risk-communication-and-community-engagement-readiness-and-initial-response-for-novel-coronaviruses-(-ncov) [accessed 2021-07-16]

12. WHO outbreak communication guidelines. World Health Organization. 2020 Jan 15. URL: https://www.who.int/publications/i/item/who-outbreak-communication-guidelines [accessed 2021-07-16]

13. Purnat T, Wilhelm E. Building systems for respond to infodemics and build resilience to misinformation. LinkedIn. 2020 Dec 02. URL: https://www.linkedin.com/pulse/building-systems-respond-infodemics-build-resilience-tina-d-purnat/ [accessed 2021-07-16]

14. Ahmed W, Bath PA, Sbaffi L, Demartini G. Novel insights into views towards H1N1 during the 2009 Pandemic: a thematic analysis of Twitter data. Health Info Libr J 2019 Mar 20;36(1):60-72. [doi: 10.1111/hir.12247] [Medline: 30663232]

15. Chew C, Eysenbach G. Pandemics in the age of Twitter: content analysis of Tweets during the 2009 H1N1 outbreak. PLoS One 2010 Nov 29;5(11):e14118 [FREE Full text] [doi: 10.1371/journal.pone.0014118] [Medline: 21124761]

16. Signorini A, Segre AM, Polgreen PM. The use of Twitter to track levels of disease activity and public concern in the U.S. during the influenza A H1N1 pandemic. PLoS One 2011 May 04;6(5):e19467 [FREE Full text] [doi: 10.1371/journal.pone.0019467] [Medline: 21573238]

17. Pulido CM, Villarejo-Carballido B, Redondo-Sama G, Gómez A. COVID-19 infodemic: more retweets for science-based information on coronavirus than for false information. Int Sociol 2020 Apr 15;35(4):377-392. [doi: 10.1177/0268580920914755]

18. Medford R, Saleh S, Sumarsono A, Perl T, Lehmann C. An "infodemic": leveraging high-volume Twitter data to understand early public sentiment for the coronavirus disease 2019 outbreak. Open Forum Infect Dis 2020 Jul;7(7):ofaa258 [FREE Full text] [doi: 10.1093/ofid/ofaa258] [Medline: 33117854]

19. Tsao S, Chen H, Tisseverasinghe T, Yang Y, Li L, Butt ZA. What social media told us in the time of COVID-19: a scoping review. Lancet Digit Health 2021 Mar;3(3):e175-e194. [doi: 10.1016/s2589-7500(20)30315-0]

20. Nsoesie EO, Cesare N, Müller M, Ozonoff A. COVID-19 misinformation spread in eight countries: exponential growth modeling study. J Med Internet Res 2020 Dec 15;22(12):e24425 [FREE Full text] [doi: 10.2196/24425] [Medline: 33264102]

21. Pascual-Ferrá P, Alperstein N, Barnett DJ. Social network analysis of COVID-19 public discourse on Twitter: implications for risk communication. Disaster Med Public Health Prep 2020 Sep 10:1-9 [FREE Full text] [doi: 10.1017/dmp.2020.347] [Medline: 32907685]

22. Vraga EK, Bode L. Using expert sources to correct health misinformation in social media. Science Communication 2017 Sep 14;39(5):621-645. [doi: 10.1177/1075547017731776]

23. Fung IC, Fu K, Chan C, Chan BSB, Cheung C, Abraham T, et al. Social media's initial reaction to information and misinformation on Ebola, August 2014: facts and rumors. Public Health Rep 2016 May;131(3):461-473 [FREE Full text] [doi: 10.1177/003335491613100312] [Medline: 27252566]

24. Bode L, Vraga EK. In related news, that was wrong: the correction of misinformation through related stories functionality in social media. J Commun 2015 Jun 23;65(4):619-638. [doi: 10.1111/jcom.12166]

25. Simpson E, Conner A. Fighting coronavirus misinformation and disinformation: preventive product recommendations for social media platforms. Center for American Progress. 2020 Aug 18. URL: https://www.americanprogress.org/issues/technology-policy/reports/2020/08/18/488714/fighting-coronavirus-misinformation-disinformation/ [accessed 2021-07-16]

26. Carmi E, Yates S, Lockley E, Pawluczuk A. Data citizenship: rethinking data literacy in the age of disinformation, misinformation, and malinformation. Internet Policy Rev 2020;9(2):1-22. [doi: 10.14763/2020.2.1481]

27. Managing epidemics: key facts about deadly diseases. World Health Organization. 2018. URL: https://www.who.int/emergencies/diseases/managing-epidemics-interactive.pdf [accessed 2021-07-16]

28. Stryker JE, Wray RJ, Hornik RC, Yanovitzky I. Validation of Database Search Terms for Content Analysis: The Case of Cancer News Coverage. Journal Mass Commun Q 2016 Jun 25;83(2):413-430. [doi: 10.1177/107769900608300212]

29. Kim Y, Huang J, Emery S. Garbage in, Garbage Out: Data Collection, Quality Assessment and Reporting Standards for Social Media Data Use in Health Research, Infodemiology and Digital Disease Detection. J Med Internet Res 2016 Feb 26;18(2):e41 [FREE Full text] [doi: 10.2196/jmir.4738] [Medline: 26920122]

30. Snelson CL. Qualitative and mixed methods social media research: a review of the literature. Int J Qual Methods 2016 Mar 01;15(1):160940691562457. [doi: 10.1177/1609406915624574]

31. Media Measurement. COVID-19 Infodemic Digital Intelligence Reports, 16-22 June 2021. Google Docs. 2021 Jun. URL: https://drive.google.com/drive/folders/1G_72yWkmFg6q7kWmukGyFd6LRY5XeyPp?usp=sharing [accessed 2021-06-29]

32. Yan L, Kang S, Guan J, Hu S. SARS-CoV-2 is an unrestricted bioweapon: a truth revealed through uncovering a large-scale, organized scientific fraud. Zenodo Preprint posted on October 8, 2020. [doi: 10.5281/zenodo.4073131]

33. Media Measurement. COVID-19 Localized Infodemic Digital Intelligence: Philippines, 15 February – 11 April 2021. Google Docs. 2021. URL: https://drive.google.com/file/d/1At76iihKfz7HaxG8KVsDttRrPb80iwi7/view?usp=sharing [accessed 2021-06-29]

34. Media Measurement. COVID-19 Localized Infodemic Digital Intelligence: Brunei Darussalam, Malaysia, Singapore, 1 January – 31 December 2020. Google Docs. 2021. URL: https://drive.google.com/file/d/1dC3NhjiCxFN2FpLpiPCkYGOUKns7lFLd/view?usp=sharing [accessed 2021-06-29]

35. Media Measurement. Analyses des discours sur la COVID-19 dans les médias sociaux au Mali 18/01-18/04 2021. Google Docs. URL: https://drive.google.com/file/d/1fXOkQl0Q-lzlIya3Ln4t2s1a9np5bgne/view?usp=sharing [accessed 2021-06-29]

36. Chouinard É. Vaccins : les échanges sur les réseaux sociaux assombrissent le portrait de la réalité. Radio-Canada Québec. URL: https://ici.radio-canada.ca/nouvelle/1791421/internet-vaccins-reseaux-sociaux-etude-intelligence-artificielle-eve-dube [accessed 2021-06-25]

37. Early AI-supported Response with Social Listening. World Health Organization. URL: https://apps.who.int/ears [accessed 2021-07-16]

38. What are the data sources in Explore? Meltwater. URL: https://help.meltwater.com/en/articles/4064549-what-are-the-data-sources-in-explore [accessed 2021-07-16]

39. Gerts D, Shelley CD, Parikh N, Pitts T, Watson Ross C, Fairchild G, et al. "Thought I'd share first" and other conspiracy theory tweets from the COVID-19 infodemic: exploratory study. JMIR Public Health Surveill 2021 Apr 14;7(4):e26527 [FREE Full text] [doi: 10.2196/26527] [Medline: 33764882]

40. Hayes JL, Britt BC, Evans W, Rush SW, Towery NA, Adamson AC. Can Social Media Listening Platforms' Artificial Intelligence Be Trusted? Examining the Accuracy of Crimson Hexagon's (Now Brandwatch Consumer Research's) AI-Driven Analyses. Journal of Advertising 2020 Sep 17;50(1):81-91. [doi: 10.1080/00913367.2020.1809576]

41. Samoilenko SA, Miroshnichenko A. Profiting from the "Trump Bump": the effects of selling negativity in the media. In: Handbook of Research on Deception, Fake News, and Misinformation Online. Hershey, PA: IGI Global; 2018.

## Abbreviations

**EARS:** Early Artificial intelligence–supported Response with Social Listening
**EPI-WIN:** World Health Organization Information Network for Epidemics
**FN:** false negative
**FP:** false positive
**pe:** expected agreement
**po:** observed agreement
**Q1:** first quarter
**Q2:** second quarter
**Q3:** third quarter
**Q4:** fourth quarter
**TN:** true negative
**TP:** true positive
**WHO:** World Health Organization

XSL•FO
**RenderX**

Original Paper

# COVID-19 Information Sources and Health Behaviors During Pregnancy: Results From a Prenatal App-Embedded Survey

James Bohnhoff[1], MD; Alexander Davis[2], PhD; Wändi Bruine de Bruin[3,4], PhD; Tamar Krishnamurti[5], PhD

[1]Division of General Pediatrics, University of Pittsburgh School of Medicine, Pittsburgh, PA, United States

[2]Department of Engineering and Public Policy, Carnegie Mellon University, Pittsburgh, PA, United States

[3]Sol Price School of Public Policy, University of Southern California, Los Angeles, CA, United States

[4]Schaeffer Center for Health Policy and Economics, University of Southern California, Los Angeles, CA, United States

[5]Division of General Internal Medicine, University of Pittsburgh School of Medicine, Pittsburgh, PA, United States

Corresponding Author:
Tamar Krishnamurti, PhD
Division of General Internal Medicine
University of Pittsburgh School of Medicine
200 Meyran Ave, Suite 200
Pittsburgh, PA, 15213
United States
Phone: 1 412 692 4855
Email: tamark@pitt.edu

## Abstract

**Background:**   Pregnancy is a time of heightened COVID-19 risk. Pregnant individuals' choice of specific protective health behaviors during pregnancy may be affected by information sources.

**Objective:**   This study examined the association between COVID-19 information sources and engagement in protective health behaviors among a pregnant population in a large academic medical system.

**Methods:**   Pregnant patients completed an app-based questionnaire about their sources of COVID-19 information and engagement in protective health behaviors. The voluntary questionnaire was made available to patients using a pregnancy app as part of their routine prenatal care between April 21 and November 27, 2020.

**Results:**   In total, 637 pregnant responders routinely accessed a median of 5 sources for COVID-19 information. The most cited source (79%) was the Centers for Disease Control and Prevention (CDC). Self-reporting evidence-based protective actions was relatively common, although 14% self-reported potentially harmful behaviors to avoid COVID-19 infection. The CDC and other sources were positively associated with engaging in protective behaviors while others (eg, US president Donald Trump) were negatively associated with protective behaviors. Participation in protective behaviors was not associated with refraining from potentially harmful behaviors (*P*=.93). Moreover, participation in protective behaviors decreased (*P*=.03) and participation in potentially harmful actions increased (*P*=.001) over the course of the pandemic.

**Conclusions:**   Pregnant patients were highly engaged in COVID-19–related information-seeking and health behaviors. Clear, targeted, and regular communication from commonly accessed health organizations about which actions may be harmful, in addition to which actions offer protection, may offer needed support to the pregnant population.

## Introduction

Pregnant people are at higher risk for severe COVID-19 illness and adverse pregnancy outcomes such as hypertensive disorders, preterm birth, and cesarean delivery [1,2]. However, the risk of vertical transmission of COVID-19 is still being studied [3], and data on efficacy and safety of COVID-19 vaccination for pregnant women lag behind those for other populations [4];

XSL•FO
**RenderX**

furthermore, recommendations on appropriate health action differ by information source, including conflicting advice by professional health organizations [5,6]. Thus, pregnant people are faced with heightened risk and less certain information when seeking knowledge of appropriate COVID-19–related health choices. Even if highly motivated to engage in positive health behaviors, pregnant people have not always known what actions would offer them appropriate protection, with many doubting the benefits of protective behaviors such as vaccination [7,8].

In the broader US population, adoption of protective behaviors continues to be uneven, despite a growing scientific consensus on effective protective behaviors to decrease the transmission and contraction of COVID-19 [9,10]. Individuals' information sources may be an important determinant of health beliefs, behaviors, and the acceptance of health guidance [11]. Use of news sources such as the Centers for Disease Control and Prevention (CDC) has been associated with COVID-19 knowledge [12] and protective action such as social distancing [13]. People have also sought COVID-19 information from other sources including social media [12], where evidence-based guidelines were often drowned out by misinformation [14] and "echo chambers" [15]. The sources that pregnant people access may, therefore, inform their willingness to implement protective behaviors.

Here, we used data collected through a pregnancy health tracking app with the aim of examining the relationship between the sources from which pregnant people seek COVID-19–related health information and their engagement in protective health behaviors. Specifically, we examined whether the information source chosen for learning about COVID-19 was associated with (1) engagement in evidence-based protective health behaviors, such as hand-washing and social distancing, and (2) potentially harmful behaviors that have been perpetuated through misinformation, such as personal use of UV radiation to treat or prevent infection [16,17] while accounting for demographic and clinical covariates. We also examined (3) whether participation in evidence-based behaviors was associated with refraining from harmful actions. We hypothesized that pregnant individuals' COVID-19 information sources are associated with their effective and potentially harmful behaviors, and that higher levels of effective health behaviors would be associated with lower levels of harmful behaviors.

## Methods

### Data Collection Tool

Providers at the University of Pittsburgh Medical Center health system prescribed the MyHealthyPregnancy (MHP) app (iOS version 1.4.7, Android version 1.8) to pregnant patients at their first prenatal appointment as part of routine prenatal care. All content was developed in conjunction with a clinical education team employed by the health care system. MHP applies machine learning algorithms to patient-entered data to model an individual patient's likelihood of adverse pregnancy events. The app offers relevant resources (eg, local health services) or actions (eg, prompts to call their provider), depending on the information that is entered into the app, as well as notifying their provider if critical health risks are documented. From April

2020, MHP added questions about COVID-19 symptoms (COVID-19 screening tool), responding to symptom reports with care-seeking guidance, and a separate COVID-19 behaviors questionnaire that included questions about COVID-19 information sources and engagement in specific protective behaviors. App users were then also offered some additional education about appropriate protective behaviors. Surveys were checked against the Checklist for Reporting the Results of E-Surveys (CHERRIES), focusing on items relevant to an app-based survey [18].

The internal protocol for prescribing MHP was to send a weblink to the patient's phone. App users electronically consented to share identifiable data with their health care provider and anonymized aggregate data for research. Participants did not receive any financial compensation for app use.

During the patient's first use of MHP (onboarding), they were prompted with 26 multiple-choice questions, over 4 screens of questioning, which included questions on demographics and pregnancy history. During the study period (April 21 to November 27, 2020), participants were invited via an SMS text message and in-app notification to voluntarily complete the COVID-19 screening tool (4 questions) and COVID-19 behaviors questionnaire (8 questions). The app's "Learning Center" was then updated for all app users, regardless of use of the screening tool or participation in the COVID-19 behaviors questionnaire. The COVID-19 screening tool remained available for use at any time.

### COVID-19–Related Information Sources and Protective Actions

As part of the COVID-19 behaviors questionnaire, participants indicated where they received their coronavirus-related information from a list of choices composed of government entities, media sources, the internet and social media, and personal contacts, with the option to list additional sources through free text. Participants were also asked to select actions that they had taken in the last month to keep themselves safe from COVID-19. These actions are enumerated in the Results section. The research team reviewed the first behaviors questionnaire completed by each participant. We categorized actions on their evidence base and potential health risk. Three actions, including (1) avoiding public spaces, gatherings, or crowds, (2) washing hands with soap or using hand sanitizer several times per day, and (3) wearing a face mask, were categorized as "most effective" in accordance with CDC recommendations [19]. Other actions, such as "cancel[ing] or postpon[ing] air travel for work" were categorized as "protective actions" on the basis that they were of known benefit but overlapped with the "most effective" actions or were not applicable to every individual. Actions related to scheduling or canceling medical appointments were categorized as indeterminate. Actions identified by the CDC and the World Health Organization (WHO) as commonly reported misconceptions were categorized as "Other unnecessary or ineffective" for preventing COVID-19 (eg, "Stockpiled food or water") or "Ineffective and potentially harmful" (eg, "Used antibiotics") [20,21]. For each participant, we recorded the number of "most effective" actions selected as well as the

XSL•FO

RenderX

selection of any "potentially harmful" actions, focusing on these 2 categories as the most likely to be of interest to organizations hoping to reduce COVID-19 spread and prevent harm. This categorization system was developed during analysis in late 2020 but attempted to describe recommendations which had been relatively consistent throughout the pandemic. In particular, wearing a face mask, though initially discouraged by the CDC, was recommended as a voluntary, protective health measure beginning early April 2020 [22].

### Other Health Information

Respondents were designated as having a high-risk pregnancy history if they reported any of the following at baseline: use of in vitro fertilization or ovulation-inducing medications, prior pregnancy loss, prior premature birth (<37 weeks) or newborn with an extended hospital stay, prior premature rupture of membranes, or diagnosis of autoimmune disease, hypertension, chronic kidney disease, or diabetes. Respondents were designated as having COVID-19–relevant symptoms during the time of survey response if they reported current fever, cough, or shortness of breath in the COVID-19 screening tool. They were also asked, "Are you experiencing financial or other personal difficulties as a result of this pandemic?"

### Statistical Power

Most of our analyses are comparisons of proportions or odds ratios of respondents reporting protective or harmful actions depending on their reported information sources. With 637 respondents, we detected a statistical difference of 10 percentage points (0.6 vs 0.5) with 80% power and a Cronbach $\alpha$ of .05 using a 2-sample test of proportions.

### Analysis

All analyses were performed using STATA (version 15.1; StataCorp, LLC). Missing data were imputed on the basis of median and mode responses. To test the association between information sources and health behaviors, we performed 2 regression analyses. First, we used linear regression analysis to assess the association between the use of individual sources and the number of "most effective" protective actions engaged in, also factoring in the model demographics (age, race, education,

number of children, and COVID-19–related distress), health characteristics (high-risk pregnancy history and COVID-19 symptoms), and survey date, which we included to measure population-level changes over the course of the pandemic. Second, we performed logistic regression analysis to test the association between the use of individual information sources and engagement in any "potentially harmful" actions. Finally, we tested the association between engagement in any "potentially harmful" factors and number of "most effective" actions undertaken using logistic regression analysis.

### Ethics Approval and Consent to Participate

The health care system's quality improvement review board approved the research. Informed consent was obtained from all study participants.

## Results

### Results Overview

In total, 637 women (22% of the 2906 active app users during the study period) completed the app-based COVID-19 survey, at a median gestational age of 15 weeks (IQR 10-24 weeks). Table 1 shows respondent demographics. The demographic characteristics of survey respondents were similar to those of all active app users. Respondents reported receiving information about COVID-19 from a median of 5 sources (IQR 3-7 sources). From the least to the most used source, 49 (8%) participants indicated receiving information from MSNBC, ranging to 505 (79%) participants who received information from the CDC. The most frequently cited free-text source was Dr Anthony Fauci (6 participants, 1%).

Table 2 shows the rates of respondent-reported COVID-19 protective actions, ranging from 2% (n=11, for each of "Used antibiotics" and "Used an ultraviolet disinfection lamp") to 98% (n=626, for "Washed your hands with soap or used hand sanitizer several times per day"). Regarding the actions categorized as most effective, only 1% (n=7) of those surveyed reported practicing none of these actions and 80% (n=512) had practiced all three. In total, 89 (14%) individuals reported at least one misguided/potentially harmful action.

**Table 1.** Demographic characteristics of the surveyed women (N=637).

| Characteristic | Value |
| --- | --- |
| Age (years), mean (SD) | 30.4 (7.2) |
| **Race, n (%)** | |
| White | 517 (81) |
| Black | 63 (10) |
| Other | 57 (9) |
| **Children, n (%)** | |
| 0 | 367 (57) |
| 1 | 167 (26) |
| ≥2 | 33 (16) |
| Income[a] (US $), mean (SD) | 66,656 (33,356) |
| **Education[b], n (%)** | |
| High school or less | 164 (26) |
| 2 or 4 years of college | 263 (41) |
| Postgraduate | 203 (32) |
| Prefer not to answer | 7 (1) |

[a]Income was collected as a categorical variable but treated as a continuous variable in analysis.

[b]In the United States, "High School" is a general education intended to be universal and to continue till the age of 18 years, followed by college and postgraduate training for some individuals. For reference, according to the most recently released educational data from the US Census bureau, 39% of the population aged ≥18 years had completed high school or had a lower education, 49% had completed some college but no postgraduate degree, and 12% had completed a postgraduate degree [23].

**Table 2.** Actions taken by study participants to decrease COVID-19 infection risk.

| Self-reported protective actions | Participants, n (%) |
|---|---|
| **Most effective actions** | |
| Washed hands with soap or used hand sanitizer several times per day | 626 (98) |
| Wore a face mask | 614 (97) |
| Avoided public spaces, gatherings, or crowds | 525 (82) |
| **Other effective actions** | |
| Avoided contact with people who could be at high risk | 503 (79) |
| Avoided eating at restaurants | 479 (75) |
| Canceled or postponed personal or social activities | 460 (72) |
| Worked or studied at home | 345 (54) |
| Ordered meals or groceries to be delivered | 340 (53) |
| Avoided your place of worship | 285 (45) |
| Canceled or postponed work or school activities | 243 (38) |
| Canceled or postponed air travel for pleasure | 201 (32) |
| Canceled or postponed air travel for work | 109 (17) |
| **Indeterminate effectiveness actions** | |
| Visited a doctor | 270 (42) |
| Canceled a doctor's appointment | 97 (15) |
| **Other unnecessary or ineffective actions** | |
| Wiped down items from the grocery store | 294 (46) |
| Wiped down packages with disinfectant | 262 (41) |
| Stockpiled food or water | 217 (34) |
| Took a hot bath | 139 (22) |
| Ate garlic | 84 (13) |
| **Ineffective and potentially harmful actions** | |
| Used a hand dryer instead of hand washing to kill the virus with heat | 32 (5) |
| Rinsed nose with saline | 22 (4) |
| Sprayed self with alcohol or chlorine | 26 (4) |
| Used other medicines or supplements not prescribed by a doctor | 17 (3) |
| Used an ultraviolet disinfection lamp | 11 (2) |
| Used antibiotics | 11 (2) |

## Information Source and Most Effective Actions

In regression analysis, those who were more likely to seek information from the CDC ($P$=.002), the WHO ($P$=.01), local health departments ($P$=.006), health care workers ($P$=.03), and public media ($P$=.04) practiced more of the 3 most effective protective actions (Table 3). Those who were less likely to obtain information from the US president (Donald Trump) or Vice-President (Mike Pence) at the time ($P$=.02) also practiced more of the most effective actions. The number of most effective actions engaged in was also positively associated with older age ($P$=.006) and negatively associated with later date of surveying ($P$=.003), higher number of children ($P$=.02), and the presence of COVID-19 symptoms ($P$=.04).

**Table 3.** COVID-19 information sources and most effective and potentially harmful actions.

| News source | Respondents citing this source or trait, n (%) | Regression coefficient for the most effective actions (95% CI) | P value | Log odds ratio for potentially harmful actions (95% CI) | P value |
|---|---|---|---|---|---|
| Centers for Disease Control and Prevention | 505 (79) | *0.18 (0.07 to 0.29)* [a] | *.002* | *−0.82 (−1.52 to −0.11)* | *.02* |
| Local Department of Heath | 425 (67) | *0.12 (0.04 to 0.21)* | *.006* | −0.17 (−0.76 to 0.41) | .57 |
| World Health Organization | 313 (49) | *0.11 (0.02 to 0.20)* | *.01* | 0.29 (−0.31 to 0.90) | .34 |
| US Department of Health | 208 (33) | 0.01 (−0.08 to 0.10) | .80 | *0.74 (0.14 to 1.34)* | *.02* |
| President Donald Trump or Vice-President Mike Pence | 75 (12) | *−0.16 (−0.29 to −0.03)* | *.02* | 0.04 (−0.79 to 0.87) | .92 |
| Health care workers | 405 (64) | *0.09 (0.01 to 0.17)* | *.03* | *0.59 (0.03 to 1.16)* | *.04* |
| Friends and family | 206 (32) | 0.03 (−0.06 to 0.13) | .51 | *1.04 (0.44 to 1.64)* | *.001* |
| Internet or social media | 137 (22) | 0.08 (−0.02 to 0.19) | .13 | −0.17 (−0.86 to 0.51) | .62 |
| Coworkers | 123 (19) | 0.01 (−0.09 to 0.12) | .82 | 0.06 (−0.62 to 0.74) | .86 |
| Local news | 192 (30) | 0 (−0.1 to 0.09) | .92 | 0.57 (−0.03 to 1.16) | .06 |
| Public media | 161 (25) | *0.10 (0.00 to 0.19)* | *.04* | −0.63 (−1.34 to 0.03) | .06 |
| National newspapers | 120 (19) | 0.04 (−0.07 to 0.15) | .45 | −0.70 (−1.48 to 0.09) | .08 |
| CNN | 112 (18) | 0.07 (−0.05 to 0.19) | .28 | −0.2 (−1.03 to 0.62) | .63 |
| NBC news | 72 (11) | −0.04 (−0.2 to 0.11) | .60 | −1.39 (−2.73 to 0.04) | .04 |
| Fox News | 65 (10) | 0.01 (−0.14 to 0.16) | .87 | 0.76 (−0.13 to 1.65) | .09 |
| ABC news | 61 (10) | −0.05 (−0.23 to 0.13) | .57 | −0.13 (−1.31 to 1.05) | .83 |
| MSNBC | 49 (8) | 0.05 (−0.13 to 0.23) | .57 | 0.26 (−0.90 to 1.43) | .66 |
| CBS news | 52 (8) | −0.07 (−0.28 to 0.14) | .51 | *1.48 (0.21 to 2.76)* | *.02* |
| **Other covariates** | | | | | |
|    Age (per 10 years) | N/A[b] | *0.08 (0.02 to 0.14)* | *.006* | 0.24 (−0.17 to 0.66) | .25 |
|    **Race** | | | | | |
|       White | N/A | Reference | | Reference | |
|       Black | N/A | 0.12 (−0.03 to 0.26) | .11 | 0.25 (−0.55 to 1.06) | .54 |
|       Other | N/A | 0.11 (−0.02 to 0.24) | .09 | *1.15 (0.40 to 1.90)* | *.003* |
|    Income (per US $10,000) | N/A | −0.00 (−0.02 to 0.01) | .75 | −0.04 (−0.14 to 0.06) | .40 |
|    **Education** | | | | | |
|       High school or less | N/A | −0.07 (−0.20 to 0.07) | .33 | 0.16 (−0.69 to 1.01) | .72 |
|       Collegiate | N/A | 0.02 (−0.08 to 0.11) | .73 | 0.39 (−0.27 to 1.04) | .25 |
|       Postgraduate | N/A | Reference | | Reference | |
|       Prefer not to answer | N/A | −0.25 (−0.62 to 0.12) | .18 | −0.30 (−2.82 to 2.23) | .82 |
|    Number of children | N/A | *−0.05 (−0.09 to −0.01)* | *.02* | −0.01 (−0.27 to 0.28) | .97 |
|    COVID-19–related distress | 193 (30) | 0.01 (−0.08 to 0.09) | .87 | −0.23 (−0.80 to 0.34) | .44 |
|    COVID-19 symptoms | 14 (2) | *−0.26 (−0.52 to −0.01)* | *.04* | −0.10 (−1.90 to 1.71) | .92 |
|    High risk pregnancy history | 276 (43) | 0.01 (−0.07 to 0.09) | .78 | −0.44 (−1.00 to 0.12) | .13 |
|    Date of survey completion | N/A | *−0.026 (−0.04 to −0.01)* | *.003* | *0.15 (0.03 to 0.26)* | *.01* |

[a]Italicized values are statistically significant. Regression coefficients for the most effective actions were generated through linear regression predicting each additional "most effective" action. Log-transformed odds ratios for harmful actions were generated through logistic regression analysis predicting *any* "misguided and potentially harmful" action. For "date of survey completion," an increase in the regressor of 1 corresponds to a 30-day (1-month) change.

[b]N/A: not applicable.

## Information Source and Potentially Harmful Actions

Citing the following information sources was positively associated with engagement in any potentially harmful actions: the US Department of Health (*P*=.02), health care workers (*P*=.04), friends and family (*P*=.001), and CBS news (*P*=.02) (Table 3). Citing the CDC (*P*=.02) was negatively associated with engaging in harmful actions. Potentially harmful actions were positively associated with a later date (*P*=.01) and being of race other than White or Black. (*P*=.003).

## Most Effective Actions and Potentially Harmful Actions

On logistic regression analysis, the number of most effective actions engaged in was not associated with participation in any potentially harmful actions (*P*=.93) (Table 3).

Regression analyses showing associations between information sources and both "other effective" and "other unnecessary or ineffective" actions are shown in Multimedia Appendix 1.

# Discussion

## Principal Findings

In this local sample, pregnant people surveyed in the first 10 months of the COVID-19 pandemic in the United States reported multiple, varied COVID-19 information sources. These information sources were associated with individuals' actions in several cases. Most significantly, we found that using the CDC as an information source was associated with most effective actions and negatively associated with harmful actions. The associations we found may have resulted from traits in the individuals we studied. For example, reporting the US president Donald Trump and Vice-President Mike Pence as an information source was associated with engaging in less protective actions, resonating with prior evidence that individuals' political affiliations often influence their information source [24] and are associated with multiple COVID-19 protective actions [25-28]. Alternatively, sources may have been actively providing different information [14] or may have communicated similar information but with different levels of clarity or different degrees of targeting information specifically to pregnant audiences [29]. While trust in the CDC has repeatedly been shown to be associated with protective actions during the COVID-19 pandemic [13,30], other associations, including those of individual networks, are more difficult to explain. It is notable that professional pregnancy and maternal health organizations were not listed as sources in open-ended responses since these may be the most reliable sources of targeted up-to-date scientific information for this population.

Overall, we found that participation in the most effective protective actions was relatively high, with more than 90% of our pregnant sample reporting mask-wearing and frequent hand-washing. However, a nontrivial minority reported participating in at least one misguided or potentially harmful action, and we did not detect an association between participating in effective actions and abstaining from harmful actions. Trusted public health sources may need to directly address which actions are not helpful, particularly for populations for whom there may be additional uncertainty

around the risk that is posed to them. Indeed, it is possible that an excess of fear and uncertainty, rather than lack of information, drives engagement in behaviors without an evidence base. This echoes a prior research finding that COVID-19 conspiracy theorists showed increased rates of protective behaviors, both those that were and those that were not recommended by governmental bodies [31].

In regression analysis controlling for information source, respondents' reports of participation in effective actions were lower, and reports of participation in potentially harmful actions was higher, over the time course of the pandemic. This finding suggests that pregnant people, although often more likely than nonpregnant people to participate in evidence-based health behaviors [32,33], may have experienced flagging motivation to adhere to guidelines—that is, "pandemic fatigue" [34,35]—over time. Alternatively, given our findings suggesting the importance of information sources, it is also possible that as pregnant people encountered more sources of information and disinformation, they had decreasing clarity over which of their actions were evidence-based, leading to a perverse use of harmful actions in pursuit of greater protection. Our survey was launched early in the pandemic, when the pregnant population may have had to rely largely on their own "mental models" of safe behavior when engaging in proactive actions. It is possible that these intuitions about safe behavior, earlier during the pandemic, were clearer to pregnant people than the conflicting or unclear recommendations they may have received from other channels over the course of the survey time period.

To our knowledge, this is the first study to report that a decline in effective actions over the course of a pandemic may be associated with a rise in spurious or dangerous actions [12]. This finding requires confirmation, ideally through longitudinal surveys that can track individual rather than population-level changes in participation in recommended and spurious, potentially harmful actions over time during a public health emergency. While we would expect time-related behavioral changes attributable or related to pandemic fatigue to be replicated in other populations, at other times, and in other public health crises, this might not be true for changes related to unclear recommendations from information sources.

The data analyzed here, which were collected through a health app integrated into routine care, also demonstrates the potential role that health apps may play in alerting clinicians to health behaviors of patients at the individual or population level. We have previously reported how the MyHealthyPregnancy app collects user-reported risk information, such as violence toward intimate partners or drug adherence, directly providing resources and alerting clinicians when critical risks are identified [36,37]. Such tools could also serve as a platform to deliver responsive information campaigns to counter health misconceptions or misinformation.

## Limitations

This study has several limitations. These cross-sectional data can demonstrate associations between information sources and behaviors but cannot prove causality, and can determine changes within a population but not within individuals, as might be achieved through repeat sampling [34]. We focus here on a

XSL•FO
**RenderX**

select population of pregnant people who engaged with a health-tracking app, which may limit the generalizability of our findings to other populations. In addition, the continuously changing dynamics of the COVID-19 pandemic were both informative and limiting. We were able to comment on changes in actions as the pandemic progressed. However, as the pandemic continues to evolve, information sources may shift owing to elections and changing media landscapes, and pregnant people will face new decisions around health behaviors as well.

## Conclusions

Pregnant people are now faced with the need to make decisions regarding COVID-19 vaccination and booster vaccination [38,39] and are adjusting their health behaviors as those around them are vaccinated. Pregnant people may further adjust their health behaviors in response to SARS-CoV-2 variants and other developments. As they continue to face additional contexts with uncertainty, dissemination of and adherence to health guidance will continue to be an important determinant of health at the population level. We found that our respondents accessed health information from several sources and that health behaviors may shift over time from effective to potentially harmful behaviors, regardless of information source. As pregnancy-relevant data continues to be gathered across agencies and institutions, it is critical that it be made widely publicly available and disseminated in accordance with best practices in health communication [14]. Most strikingly, perhaps, our findings show that even those individuals motivated to engage in "best-practice" behaviors were not necessarily less likely to also practice ineffective or harmful behaviors. Moving forward, key health organizations that are routinely viewed as sources of reliable health information, such as the CDC, should make a concerted effort to offer structured guidance on what *should* and *should not* be practiced in terms of protective behaviors for this population specifically.

Multimedia Appendix 1
COVID-19 information sources and other effective and other unnecessary or ineffective actions.
[DOCX File , 20 KB - infodemiology_v1i1e31774_app1.docx ]

## References

1. Stafford IA, Parchem JG, Sibai BM. The coronavirus disease 2019 vaccine in pregnancy: risks, benefits, and recommendations. Am J Obstet Gynecol 2021 May;224(5):484-495 [FREE Full text] [doi: 10.1016/j.ajog.2021.01.022] [Medline: 33529575]

2. Zambrano LD, Ellington S, Strid P, Galang RR, Oduyebo T, Tong VT, CDC COVID-19 Response Pregnancy and Infant Linked Outcomes Team. Update: Characteristics of Symptomatic Women of Reproductive Age with Laboratory-Confirmed SARS-CoV-2 Infection by Pregnancy Status - United States, January 22-October 3, 2020. MMWR Morb Mortal Wkly Rep 2020 Nov 06;69(44):1641-1647 [FREE Full text] [doi: 10.15585/mmwr.mm6944e3] [Medline: 33151921]

3. Halici-Ozturk F, Ocal FD, Aydin S, Tanacan A, Ayhan SG, Altinboga O, et al. Investigating the risk of maternal-fetal transmission of SARS-CoV-2 in early pregnancy. Placenta 2021 Mar;106:25-29 [FREE Full text] [doi: 10.1016/j.placenta.2021.02.006] [Medline: 33610934]

4. Gray KJ, Bordt EA, Atyeo C, Deriso E, Akinwunmi B, Young N, et al. COVID-19 vaccine response in pregnant and lactating women: a cohort study. medRxiv Mar 08 Preprint posted online March 8, 2021 [FREE Full text] [doi: 10.1101/2021.03.07.21253094] [Medline: 33758889]

5. ACOG and SMFM Joint Statement on WHO Recommendations Regarding COVID-19 Vaccines and Pregnant Individuals. American College of Obstetricians and Gynecologists. URL: https://www.acog.org/news/news-releases/2021/01/acog-and-smfm-joint-statement-on-who-recommendations-regarding-covid-19-vaccines-and-pregnant-individuals [accessed 2021-04-20]

6. Giles ML, Gunatilaka A, Palmer K, Sharma K, Roach V. Alignment of national COVID-19 vaccine recommendations for pregnant and lactating women. Bull World Health Organ 2021 Oct 01;99(10):739-746 [FREE Full text] [doi: 10.2471/BLT.21.286644] [Medline: 34621092]

7. Battarbee AN, Stockwell MS, Varner M, Newes-Adeyi G, Daugherty M, Gyamfi-Bannerman C, et al. Attitudes Toward COVID-19 Illness and COVID-19 Vaccination among Pregnant Women: A Cross-Sectional Multicenter Study during August-December 2020. Am J Perinatol 2021 Oct 01. [doi: 10.1055/s-0041-1735878] [Medline: 34598291]

XSL•FO

RenderX

8.    Hirshberg JS, Huysman BC, Oakes MC, Cater EB, Odibo AO, Raghuraman N, et al. Offering onsite COVID-19 vaccination to high-risk obstetrical patients: initial findings. Am J Obstet Gynecol MFM 2021 Nov;3(6):100478 [FREE Full text] [doi: 10.1016/j.ajogmf.2021.100478] [Medline: 34481996]

9.    Matrajt L, Leung T. Evaluating the Effectiveness of Social Distancing Interventions to Delay or Flatten the Epidemic Curve of Coronavirus Disease. Emerg Infect Dis 2020 Aug;26(8):1740-1748 [FREE Full text] [doi: 10.3201/eid2608.201093] [Medline: 32343222]

10.   Social Distancing: Keep a Safe Distance to Slow the Spread. Centers for Disease Control and Prevention. URL: https://stacks.cdc.gov/view/cdc/90522 [accessed 2021-04-20]

11.   Underwood NL, Gargano LM, Jacobs S, Seib K, Morfaw C, Murray D, et al. Influence of Sources of Information and Parental Attitudes on Human Papillomavirus Vaccine Uptake among Adolescents. J Pediatr Adolesc Gynecol 2016 Dec;29(6):617-622. [doi: 10.1016/j.jpag.2016.05.003] [Medline: 27216710]

12.   Ali SH, Foreman J, Tozan Y, Capasso A, Jones AM, DiClemente RJ. Trends and Predictors of COVID-19 Information Sources and Their Relationship With Knowledge and Beliefs Related to the Pandemic: Nationwide Cross-Sectional Study. JMIR Public Health Surveill 2020 Oct 08;6(4):e21071 [FREE Full text] [doi: 10.2196/21071] [Medline: 32936775]

13.   Fridman I, Lucas N, Henke D, Zigler CK. Association Between Public Knowledge About COVID-19, Trust in Information Sources, and Adherence to Social Distancing: Cross-Sectional Survey. JMIR Public Health Surveill 2020 Sep 15;6(3):e22060 [FREE Full text] [doi: 10.2196/22060] [Medline: 32930670]

14.   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]

15.   Zollo F, Bessi A, Del Vicario M, Scala A, Caldarelli G, Shekhtman L, et al. Debunking in a world of tribes. PLoS One 2017;12(7):e0181821 [FREE Full text] [doi: 10.1371/journal.pone.0181821] [Medline: 28742163]

16.   Coronavirus: Outcry after Trump suggests injecting disinfectant as treatment. BBC News. 2020 Apr 24. URL: https://www.bbc.com/news/world-us-canada-52407177 [accessed 2021-04-20]

17.   Nierenberg A. Please Do Not Eat Disinfectant. The New York Times. URL: https://www.nytimes.com/article/coronavirus-disinfectant-inject-ingest.html [accessed 2021-04-20]

18.   Eysenbach G. Improving the quality of Web surveys: the Checklist for Reporting Results of Internet E-Surveys (CHERRIES). J Med Internet Res 2004 Sep 29;6(3):e34 [FREE Full text] [doi: 10.2196/jmir.6.3.e34] [Medline: 15471760]

19.   How to Protect Yourself & Others. Centers for Disease Control and Prevention. URL: https://www.cdc.gov/coronavirus/2019-ncov/prevent-getting-sick/prevention.html [accessed 2021-03-22]

20.   Coronavirus disease (COVID-19) advice for the public: Mythbusters. World Health Organization. URL: https://www.who.int/emergencies/diseases/novel-coronavirus-2019/advice-for-public/myth-busters [accessed 2021-05-05]

21.   Centers for Disease Control and Prevention. URL: https://www.cdc.gov/coronavirus/2019-ncov/daily-life-coping/essential-goods-services.html [accessed 2021-05-05]

22.   CDC Now Recommends Americans Consider Wearing Cloth Face Coverings In Public. NPR. 2020 Apr 03. URL: https://tinyurl.com/t46yh6x7 [accessed 2021-10-25]

23.   Educational Attainment in the United States: 2019. United States Census Bureau. 2020 Mar 30. URL: https://www.census.gov/content/census/en/data/tables/2019/demo/educational-attainment/cps-detailed-tables.html [accessed 2021-08-08]

24.   Bruine de Bruin W, Saw H, Goldman DP. Political polarization in US residents' COVID-19 risk perceptions, policy preferences, and protective behaviors. J Risk Uncertain 2020 Nov 18:1-18 [FREE Full text] [doi: 10.1007/s11166-020-09336-3] [Medline: 33223612]

25.   Educational Attainment in the United States: 2019. Pew Research Center. 2020 Jun 25. URL: https://www.pewresearch.org/politics/2020/06/25/republicans-democrats-move-even-further-apart-in-coronavirus-concerns/ [accessed 2020-08-31]

26.   Leventhal AM, Dai H, Barrington-Trimis JL, McConnell R, Unger JB, Sussman S, et al. Association of Political Party Affiliation With Physical Distancing Among Young Adults During the COVID-19 Pandemic. JAMA Intern Med 2021 Mar 01;181(3):399-403. [doi: 10.1001/jamainternmed.2020.6898] [Medline: 33315091]

27.   Travis J, Harris S, Fadel T, Webb G. Identifying the determinants of COVID-19 preventative behaviors and vaccine intentions among South Carolina residents. PLoS One 2021;16(8):e0256178 [FREE Full text] [doi: 10.1371/journal.pone.0256178] [Medline: 34432817]

28.   Naeim A, Baxter-King R, Wenger N, Stanton AL, Sepucha K, Vavreck L. Effects of Age, Gender, Health Status, and Political Party on COVID-19-Related Concerns and Prevention Behaviors: Results of a Large, Longitudinal Cross-sectional Survey. JMIR Public Health Surveill 2021 Apr 28;7(4):e24277 [FREE Full text] [doi: 10.2196/24277] [Medline: 33908887]

29.   Krishnamurti T, Bruine de Bruin W. Developing Health Risk Communications: Four Lessons Learned. In: Raue M, Lermer E, Streicher B, editors. Psychological Perspectives on Risk and Risk Analysis. Cham: Springer; 2018:299-309.

30.   Sinicrope PS, Maciejko LA, Fox JM, Steffens MT, Decker PA, Wheeler P, et al. Factors associated with willingness to wear a mask to prevent the spread of COVID-19 in a Midwestern Community. Prev Med Rep 2021 Dec;24:101543 [FREE Full text] [doi: 10.1016/j.pmedr.2021.101543] [Medline: 34493965]

31. Juanchich M, Sirota M, Jolles D, Whiley LA. Are COVID-19 conspiracies a threat to public health? Psychological characteristics and health protective behaviours of believers. Eur J Soc Psychol 2021 Jun 16 [FREE Full text] [doi: 10.1002/ejsp.2796] [Medline: 34518709]

32. Vaz MJR, Barros SMO, Palacios R, Senise JF, Lunardi L, Amed AM, et al. HIV-infected pregnant women have greater adherence with antiretroviral drugs than non-pregnant women. Int J STD AIDS 2007 Jan;18(1):28-32. [doi: 10.1258/095646207779949808] [Medline: 17326859]

33. Ickovics JR, Wilson TE, Royce RA, Minkoff HL, Fernandez MI, Fox-Tierney R, Perinatal Guidelines Evaluation Group. Prenatal and postpartum zidovudine adherence among pregnant women with HIV: results of a MEMS substudy from the Perinatal Guidelines Evaluation Project. J Acquir Immune Defic Syndr 2002 Jul 01;30(3):311-315. [doi: 10.1097/00126334-200207010-00007] [Medline: 12131568]

34. MacIntyre CR, Nguyen P, Chughtai AA, Trent M, Gerber B, Steinhofel K, et al. Mask use, risk-mitigation behaviours and pandemic fatigue during the COVID-19 pandemic in five cities in Australia, the UK and USA: A cross-sectional survey. Int J Infect Dis 2021 May;106:199-207 [FREE Full text] [doi: 10.1016/j.ijid.2021.03.056] [Medline: 33771668]

35. Gidengil CA, Parker AM, Zikmund-Fisher BJ. Trends in risk perceptions and vaccination intentions: a longitudinal study of the first year of the H1N1 pandemic. Am J Public Health 2012 Apr;102(4):672-679 [FREE Full text] [doi: 10.2105/AJPH.2011.300407] [Medline: 22397349]

36. Krishnamurti T, Davis AL, Quinn B, Castillo AF, Martin KL, Simhan HN. Mobile Remote Monitoring of Intimate Partner Violence Among Pregnant Patients During the COVID-19 Shelter-In-Place Order: Quality Improvement Pilot Study. J Med Internet Res 2021 Feb 19;23(2):e22790 [FREE Full text] [doi: 10.2196/22790] [Medline: 33605898]

37. Krishnamurti T, Davis AL, Rodriguez S, Hayani L, Bernard M, Simhan HN. Use of a Smartphone App to Explore Potential Underuse of Prophylactic Aspirin for Preeclampsia. JAMA Netw Open 2021 Oct 01;4(10):e2130804 [FREE Full text] [doi: 10.1001/jamanetworkopen.2021.30804] [Medline: 34714341]

38. Male V. Are COVID-19 vaccines safe in pregnancy? Nat Rev Immunol 2021 Apr;21(4):200-201 [FREE Full text] [doi: 10.1038/s41577-021-00525-y] [Medline: 33658707]

39. Pasricha T. Should You Get a Covid Booster if You Are Pregnant? The New York Times. URL: https://www.nytimes.com/2021/10/19/well/family/pregnancy-covid-booster.html? [accessed 2021-10-25]

## Abbreviations

**CDC:** Centers for Disease Control and Prevention
**MHP:** MyHealthyPregnancy
**WHO:** World Health Organization

XSL·FO
**RenderX**

Original Paper

# Examining the Public's Most Frequently Asked Questions Regarding COVID-19 Vaccines Using Search Engine Analytics in the United States: Observational Study

Nicholas B Sajjadi[1*], BSc; Samuel Shepard[1*], BSc; Ryan Ottwell[2,3*], DO; Kelly Murray[4*], PharmD; Justin Chronister[5*], DO; Micah Hartwell[1,6*], PhD; Matt Vassar[1,6*], PhD

[1]Office of Medical Student Research, College of Osteopathic Medicine, Oklahoma State University Center for Health Sciences, Tulsa, OK, United States

[2]Department of Internal Medicine, University of Oklahoma School of Community Medicine, Tulsa, OK, United States

[3]Department of Dermatology, St. Joseph Mercy Hospital, Ann Arbor, MI, United States

[4]Department of Emergency Medicine, College of Osteopathic Medicine, Oklahoma State University Center for Health Sciences, Tulsa, OK, United States

[5]Department of Internal Medicine, College of Osteopathic Medicine, Oklahoma State University Center for Health Sciences, Tulsa, OK, United States

[6]Department of Psychiatry and Behavioral Sciences, College of Osteopathic Medicine, Oklahoma State University Center for Health Sciences, Tulsa, OK, United States

[*]all authors contributed equally

Corresponding Author:
Nicholas B Sajjadi, BSc
Office of Medical Student Research
College of Osteopathic Medicine
Oklahoma State University Center for Health Sciences
1111 W 17th Street
Tulsa, OK, 74104
United States
Phone: 1 9185821972
Fax: 1 9185821972
Email: nicholas.sajjadi@okstate.edu

## Abstract

**Background:** The emergency authorization of COVID-19 vaccines has offered the first means of long-term protection against COVID-19–related illness since the pandemic began. It is important for health care professionals to understand commonly held COVID-19 vaccine concerns and to be equipped with quality information that can be used to assist in medical decision-making.

**Objective:** Using Google's RankBrain machine learning algorithm, we sought to characterize the content of the most frequently asked questions (FAQs) about COVID-19 vaccines evidenced by internet searches. Secondarily, we sought to examine the information transparency and quality of sources used by Google to answer FAQs on COVID-19 vaccines.

**Methods:** We searched COVID-19 vaccine terms on Google and used the "People also ask" box to obtain FAQs generated by Google's machine learning algorithms. FAQs are assigned an "answer" source by Google. We extracted FAQs and answer sources related to COVID-19 vaccines. We used the Rothwell Classification of Questions to categorize questions on the basis of content. We classified answer sources as either academic, commercial, government, media outlet, or medical practice. We used the Journal of the American Medical Association's (JAMA's) benchmark criteria to assess information transparency and Brief DISCERN to assess information quality for answer sources. FAQ and answer source type frequencies were calculated. Chi-square tests were used to determine associations between information transparency by source type. One-way analysis of variance was used to assess differences in mean Brief DISCERN scores by source type.

**Results:** Our search yielded 28 unique FAQs about COVID-19 vaccines. Most COVID-19 vaccine–related FAQs were seeking factual information (22/28, 78.6%), specifically about safety and efficacy (9/22, 40.9%). The most common source type was media outlets (12/28, 42.9%), followed by government sources (11/28, 39.3%). Nineteen sources met 3 or more JAMA benchmark criteria with government sources as the majority (10/19, 52.6%). JAMA benchmark criteria performance did not significantly

differ among source types ($\chi^2_4$=7.40; *P*=.12). One-way analysis of variance revealed a significant difference in mean Brief DISCERN scores by source type ($F_{4,23}$=10.27; *P*<.001).

**Conclusions:** The most frequently asked COVID-19 vaccine–related questions pertained to vaccine safety and efficacy. We found that government sources provided the most transparent and highest-quality web-based COVID-19 vaccine–related information. Recognizing common questions and concerns about COVID-19 vaccines may assist in improving vaccination efforts.

## Introduction

As of August 01, 2021, COVID-19 has affected over 198 million people and has been responsible for over 4.2 million deaths worldwide [1,2]. In response to the pandemic, the US Food and Drug Administration issued emergency use authorizations for 2 COVID-19 vaccines in late 2020, 1 manufactured by Pfizer-BioNTech and the second by Moderna [3,4]. Overcoming logistical barriers will be crucial for enabling successful vaccine campaigns. Additionally, addressing the public's perception of COVID-19 vaccines and the quality of available information is vital for promoting positive public reception and reducing vaccine hesitancy. Vaccine hesitancy, which refers to reluctance or refusal to receive vaccines, is complex and is determined by numerous factors such as trust in vaccine safety and efficacy, perceived risk of receiving or refusing a vaccine, and accessibility to and affordability of vaccines [5]. Hesitancy toward COVID-19 vaccines may hinder successful vaccination efforts.

The pace of vaccine development, misinformation, and overall growth in vaccine hesitancy are factors potentially contributing to COVID-19 vaccine refusal [5,6]>. Identifying factors associated with COVID-19 vaccine refusal may assist in developing strategies to reduce vaccine hesitancy. To identify demographic factors associated with COVID-19 vaccine acceptance, Lazarus et al [7] surveyed individuals in 19 countries and reported that individuals who reported a high degree of trust in the government were more likely to report vaccine acceptance than those with low trust. In the United States, a survey study by the US Census Bureau showed that 49% of respondents were reluctant to receive a COVID-19 vaccine. Of those reluctant to receive COVID-19 vaccines, the most common reason for reluctance was concern for side effects. The second most common reason was planning to wait and see if the vaccines were safe [8]. A US survey conducted early in the pandemic sought to predict COVID-19 vaccine acceptance in the United States and found that several vulnerable populations reported low willingness [9]. The growing prevalence of vaccine hesitancy highlights the importance of clinician preparedness to address patients' concerns as access to COVID-19 vaccines grows. Health care professionals should serve as reliable sources of vaccine information, instilling confidence in patients and potentially enhancing vaccine acceptance [10], especially for COVID-19 vaccines [11].

Apart from consulting health care professionals, individuals frequently use the internet when seeking health care information; some use the internet as their primary source for health information [12]. In the United States, 61% of adults have searched the internet for medical information [13]. Searching the internet for medical information simultaneously presents benefits and challenges regarding patient-provider interactions [14]. The increasingly common practice of using the internet to obtain health care information makes it possible to study commonly held medical concerns by examining searching patterns and behaviors. Previous studies have documented the prevalence of COVID-19 vaccine hesitancy in the United States [8,15] and globally [7], but none of these studies explored the content of COVID-19 vaccine concerns evidenced by internet searching. Moreover, the quality of COVID-19 vaccine information resulting from internet searching has yet to be investigated. Thus, the primary objective of this study was to use Google's RankBrain machine learning algorithm to characterize the content of the most frequently asked questions (FAQs) about COVID-19 vaccines in the United States. Secondarily, we sought to grade the transparency and quality of suggested information regarding COVID-19 vaccines. We aim to equip health care professionals and researchers with information about the common concerns regarding COVID-19 vaccines, possibly supporting more successful vaccination efforts. We hypothesize that most COVID-19 vaccine–related FAQs in the United States will pertain to safety and efficacy, as survey studies have indicated these concerns as the most important driver of COVID-19 vaccine hesitancy in the United States.

## Methods

### Background

We used Google to perform our search as it is the most frequently used search engine globally as of 2015 [16]. Moreover, Google's search engine uses a powerful machine learning system called RankBrain [17] alongside the natural language processing technology known as Bidirectional Encoder Representations from Transformers [18] to detect patterns from large volumes of search queries. Google assesses the intent of a search query using rigorous language processing algorithms to sort through billions of indexed webpages and to suggest the ones most relevant to the search [19]. The resulting patterns and data are used to formulate lists of FAQs related to the original search contents. FAQs are found in boxes labeled

"People also ask" or "Common questions." Google assigns each FAQ a link to information that "answers" the question [20]. Google uses its webmaster guidelines to remove low-quality spam websites from search results and prioritize high-quality sources using a system called PageRank [19]. Taken together, these FAQs represent millions of common inquiries regarding medical information. Linked answers to each FAQ reveal which information sources individuals are likely to encounter when searching Google for medical information. Our methodology was adapted from a study by Shen et al [21], who used Google FAQs to reliably reveal common concerns about orthopedic procedures and to assess the transparency of the suggested information.

## Systematic Search

On January 23, 2021, using a newly installed web browser to minimize personalized advertisement algorithms, we separately searched Google [22] for the following three terms: "covid 19 vaccine," "pfizer covid vaccine," and "moderna covid vaccine." We selected these terms to capture the most likely general inquiries concerning the only 2 COVID-19 vaccines available at the time of our search. For each inquiry, we refreshed the list of FAQs found in the "Common questions" or "People also ask" box generated by Google. By expanding the tab on a FAQ, additional FAQs appear. We repeated this process until reaching a minimum of 150 FAQs for each search, as studies using similar methodology have recommended using 50-150 sources [21]. We used the high end of the recommended number of sources (150) for two reasons: to increase the likelihood of encountering an FAQ that would be pertinent to the current study and to reflect the precedent set in the literature. Since

query results are tailored to the user's location, search history, and search settings, we used clean browsers to minimize any influence of history and settings while allowing results to reflect queries from the United States [19].

## Data Extraction

Of the resultant FAQs, we extracted only those directly pertaining to or mentioning COVID-19 vaccines along with their answer links. In a masked duplicated fashion, investigators NS and SS extracted these data using a Google Form on January 23, 2021. FAQ data extraction was completed on January 23, 2021. After extraction, any duplicate FAQs from the individual searches were removed, followed by the removal of any duplicate FAQs among the 3 searches. After the screening and reduction process, our searches resulted in a compilation of unique FAQs regarding COVID-19 vaccines.

## Question Classification and Answer Source Type

Applying methodology adapted from previous studies [16,21], we first used the Rothwell Classification of Questions [23] to categorize FAQs under three broad categories: fact, policy, and value. Fact questions were further subclassified into four groups: safety and efficacy, vaccine administration schedule, cost, and technical details. Policy questions were subclassified into two groups: indications and complications. Value questions were subclassified into two groups: evaluation of credibility and appraisal of risk or benefit. Next, we categorized answer sources as either commercial, academic, medical practice, government, or media outlet according to previously established classification schemes [21,24]. Table 1 shows the *Question Classification* and *Answer Source Type* definitions. For each answer source, we extracted the country of origin.

**Table 1.** The Rothwell Classification of Questions, Question Classification by Topic, and Answer Source Type.

| Rothwell classification | Description |
| --- | --- |
| Fact | Asks objective, factual information regarding COVID-19 vaccines (ie, "How long does it take the vaccine to work?") |
| Policy | Asks information on a specific course of action under given circumstances related to COVID-19 (ie, should people on immunosuppressants get the vaccine?) |
| Value | Asks to conceptually evaluate COVID-19 vaccines (ie, "Will the COVID-19 vaccine work better than masks?") |
| **Question subclassification by topic** | |
| **Fact** | |
| Safety and efficacy | Questions about vaccine safety including side effects and how well the vaccine works |
| Vaccine administration schedule | Specific questions about the vaccine schedule, number of shots, and vaccine distribution |
| Cost | Cost of the vaccine, whether it is free, or who is paying for it |
| Technical Details | Mechanism by which the vaccine works, including specific questions about immunologic responses |
| **Policy** | |
| Indications | Who should or should not receive a COVID-19 vaccine |
| Complications | Questions about specific complications after being vaccinated |
| **Value** | |
| Evaluation of Credibility | Seeking authoritative approval from a trustworthy source; seeking ethos |
| Appraisal of Risk or Benefit | Necessity of preventive measures after vaccination (ie, "Is getting vaccinated worth it?") |
| **Answer source type** | |
| Commercial | Organization that publishes medical information that is not otherwise associated with an academic institution, government agency, health care system, or nonmedical news outlet such as WebMD and Healthline |
| Academic | Institution with clear academic affiliations, as evidenced by information on the website that did not better meet criteria for another classification or website ending in ".*edu*," such as Mayo Clinic and Harvard University |
| Medical practice | Affiliation with a health care system or individual health care professional who did not explicitly state a commercial, academic, or government affiliation, such as private practice and a hospital system |
| Government | Websites hosted by government organizations or sources from websites ending in ".gov," such as the Centers for Disease Control and the US Food and Drug Administration |
| Media outlet | Nonmedical organizations or social media pages claiming to publish news-related stories for the purpose of information-sharing in the form of interviews, blog posts, or articles, such as the National Public Radio, Wall Street Journal, and USA Today |

## Information Transparency and Quality

The Journal of the American Medical Association's (JAMA's) benchmark criteria [25] was then used to assess information transparency for each answer source. JAMA benchmark criteria have been used to effectively screen web-based information for fundamental aspects of information transparency [21,26-28].

JAMA benchmark criteria were also used to characterize web-based misinformation regarding COVID-19 in early 2020 [29]. Sources meeting 3 more criteria are considered to have high transparency, while sources meeting less than 3 criteria have poor transparency. Table 2 lists the JAMA benchmark criteria definitions.

**Table 2.** Journal of the American Medical Association's benchmark criteria.

| Criteria | Description |
| --- | --- |
| Authorship | Clearly identifiable author and contributors with affiliations and relevant credentials present. |
| Attribution | References and sources clearly listed with any copyright information disclosed. |
| Currency | Clearly identifiable posting date of any content as well as the date of any revisions. |
| Disclosure | Website ownership clearly disclosed along with any sponsorship, advertising, underwriting, and financial support. |

The information quality was assessed using the Brief DISCERN information quality assessment tool. DISCERN is a series of questions originally developed by Charnock et al [30] as a means for patients and providers to quickly and reliably ascertain the

quality of written health care information regarding medical treatments. The DISCERN quality assessment tool has been used to assess the quality of internet sources in a variety of medical fields [31-33]. Khazaal et al [34] developed an abbreviated 6-item version (Brief DISCERN) with comparable reliability and validity, which preserves the advantages of the original tool while affording a potentially more user-friendly format. Thus, we used the Brief DISCERN quality assessment

tool, which has been previously used [35,36]. Sources are scored from 1 to 5 based on the criteria listed in Table 3.

Authors NS and SS applied the JAMA benchmark criteria and the Brief DISCERN tool in a masked duplicate fashion, and author MH resolved any discrepancies. This protocol was submitted to the institutional review board of Oklahoma State University Center for Health Sciences and was determined to be non–Human Subjects Research.

**Table 3.** Brief DISCERN questions and scoring.

| Questions | Score | | |
|---|---|---|---|
| | Low (1) "No" | Moderate (2-4) "Partially" | High (5) "Yes" |
| Is it clear what sources of information were used to compile the publication (other than the author or producer)? | No sources of evidence for the information are mentioned | The sources are clear to some extent and are referenced in the text *or* in a bibliography | The sources are very clear and are referenced in text *and* in a bibliography |
| Is it clear when the information used or reported in the publication was produced? | No dates have been given | Only the date of the publication itself is clear, or dates for some of but not all acknowledged sources are given | Dates for all acknowledged sources are clear |
| Does it describe how each treatment works? | None of the descriptions about treatments include details of how it works | Descriptions of some but not all treatments are given *or* the details provided are unclear or incomplete | The description of treatment includes details of how it works |
| Does it describe the benefits of each treatment? | No benefits are described | A benefit is described for some but not all treatments | A benefit is described for each treatment |
| Does it describe the risk of each treatment? | No risks are described for any of the treatments. | A risk is described for some but not all treatments. | A risk is described for each treatment. |
| Does it describe how the treatment choices affect overall quality of life? | There is no reference to overall quality of life in relation to treatment choices. | The publication includes a reference to overall quality of life in relation to treatment choices, but the information is unclear or incomplete. | The publication includes a clear reference to overall quality of life in relation to any of the treatment choices mentioned. |

## Analyses

Frequencies and percentages were reported for each FAQ's classification. Chi-square tests were used to determine associations between JAMA benchmark criteria by source type. One-way analysis of variance was used to determine whether the mean Brief DISCERN score differed by source type. Post hoc comparisons, performed using *t* tests with Bonferroni correction, were used to identify mean differences between source type categories. Interrater agreement for each assessment was determined using intraclass correlation coefficients.

## *Results*

A total of 467 FAQs were generated from all 3 searches: 161 from "covid 19 vaccine," 155 from "moderna covid vaccine," and 151 from "pfizer covid vaccine." Of these, "covid 19 vaccine" yielded 5 vaccine-related FAQs, "moderna covid vaccine" yielded 22, and "pfizer covid vaccine" yielded 14. After removing duplicates, our searches yielded a total of 28 unique FAQs regarding COVID-19 vaccines (Table 4).

**Table 4.** List of the 28 unique frequently asked questions regarding COVID-19 vaccines.

| Frequently asked questions | Rothwell classification | Subclassification | Answer source | JAMA benchmark criteria (≥3) | Brief DISCERN score |
|---|---|---|---|---|---|
| Are both Covid vaccines 2 doses? | Fact | Vaccine administration schedule | Commercial | No | 15 |
| Are you immune to Covid after vaccine? | Fact | Safety and efficacy | Media outlet | No | 21 |
| Can I get COVID-19 right after being vaccinated? | Fact | Technical details | Government | Yes | 29 |
| Can the COVID-19 vaccine make you sick? | Fact | Safety and efficacy | Government | Yes | 29 |
| Can you still get Covid after first vaccine? | Fact | Technical details | Media outlet | No | 18 |
| Can you test positive for Covid after vaccine? | Fact | Technical details | Media outlet | No | 9 |
| Do COVID-19 vaccines require more than one shot? | Fact | Vaccine administration schedule | Government | Yes | 29 |
| Do you have to wait 90 days after Covid to get the vaccine? | Fact | Vaccine administration schedule | Media outlet | Yes | 15 |
| Do you have to wear mask after Covid vaccine? | Value | Risk/benefit appraisal | Media outlet | Yes | 28 |
| Does Covid vaccine Stop Spread? | Value | Risk/benefit appraisal | Media outlet | Yes | 22 |
| Has the Pfizer-BioNTech COVID-19 vaccine been authorized by the FDA? | Fact | Safety and efficacy | Government | Yes | 30 |
| How does the COVID-19 mRNA vaccine work? | Fact | Technical details | Government | No | 25 |
| How effective is the Pfizer COVID-19 vaccine? | Fact | Safety and efficacy | Media outlet | No | 17 |
| How long do you have to wait between Covid vaccines? | Fact | Vaccine administration schedule | Media outlet | No | 16 |
| How many shots of Moderna COVID-19 vaccine should I get? | Fact | Vaccine administration schedule | Government | Yes | 29 |
| Is it safe to take the COVID-19 vaccine? | Fact | Safety and efficacy | Government | Yes | 29 |
| Is the Moderna vaccine for COVID-19 approved by the FDA? | Fact | Safety and efficacy | Academic | Yes | 30 |
| Should you get the Covid vaccine if you were previously infected with Covid? | Policy | Indications | Media outlet | Yes | 15 |
| What are some common side effects of the COVID-19 vaccine? | Fact | Safety and efficacy | Government | Yes | 29 |

## Question Classification

Using the Rothwell classification system, the majority of FAQs were seeking factual information (22/28;78.6%). Among these factual questions, the most common topic was safety and efficacy (9/22, 40.9%) followed by technical details (6/22, 27.3%), vaccine administration schedule (6/22, 27.3%), and cost (1/22, 4.5%) (Table 4).

## Answer Sources

The most common answer source type overall was media outlets (12/28, 42.9%), followed by government sources (11/28, 39.3%), commercial sources (3/28, 10.7%), academic sources (1/28, 3.55%), and medical practice (1/28, 3.55%). FAQs classified as technical details were most frequently answered by a media outlet (4/6, 66.7%). Of FAQs classified as fact, most were answered by government sources (11/22, 50%). Government sources also most commonly answered FAQs related to safety

and efficacy (5/9, 55.6%), cost (1/1, 100%), and vaccine administration schedule (3/6, 50%) (Table 4). In total, 26 of 28 (92.8%) answer sources were from the United States, 1 was from the United Kingdom (3.6%), and 1 was from Australia (3.6%).

## Information Transparency

In total, 19 sources met 3 or more JAMA benchmark criteria, of which government sources were the majority (10/19, 52.6%), followed by media outlets (7/19, 36.8%), commercial sources (1/19, 5.3%), and academic sources (1/19, 5.3%). Among sources meeting less than 3 criteria, media outlets were the most common (5/9, 55.6%), followed by commercial sources (2/9, 22.2%), medical practice (1/9, 11.1%), and government sources (1/9, 11.1%). Approximately 92.7% (11/12) of government sources met 3 or more JAMA benchmark criteria, whereas 58.3% (7/12) of media outlets met 3 or more criteria. The overall JAMA Benchmark Criteria performance did not significantly

XSL•FO

RenderX

differ among source types ($\chi^2_4$=7.40; $P$=.12); however, we found significant associations between individual source's performance on meeting JAMA benchmark criteria for authorship and the source type ($\chi^2_4$=*18.03*, *P*<.001), with 11/28 (39.3%) media outlet sources meeting authorship criteria compared to 10/28

(35.7%) government sources not meeting the authorship criteria. We also found a similar but negative relationship with JAMA benchmark criteria's disclosure criteria and source type ($\chi^2_4$=15.36; $P$=.004) with 10/28 (35.7%) government sources meeting these criteria compared to 9/28 (32.1%) media outlets not meeting these criteria (Tables 5 and 6).

**Table 5.** Journal of the American Medical Association's benchmark criteria and by source type.

| Sources meeting 3 or more JAMA benchmark criteria | Source type, n (%) | | | | | Total | Chi-square (*df*) | *P* value |
|---|---|---|---|---|---|---|---|---|
| | Academic | Commercial | Government | Medical practice | Media outlet | | | |
| **Journal of the American Medical Association's benchmark criteria** | | | | | | | 7.40 (*4*) | .12 |
| 3+ | 1 (3.6) | 1 (3.6) | 10 (35.7) | 0 (0.0) | 7 (25.0) | 19 (67.9) | | |
| <3 | 0 (0.0) | 2 (7.1) | 1 (3.6) | 1 (3.6) | 5 (17.9) | 9 (32.1) | | |
| **Authorship** | | | | | | | 18.03 (*4*) | .001 |
| No | 1 (3.6) | 2 (7.1) | 10 (35.7) | 0 (0.0) | 1 (3.6) | 14 (50.0) | | |
| Yes | 0 (0.0) | 1 (3.6) | 1 (3.6) | 1 (3.6) | 11 (39.3) | 14 (50.0) | | |
| **Attribution** | | | | | | | 7.21 (*4*) | .13 |
| No | 0 (0.0) | 2 (7.1) | 1 (3.6) | 1 (3.6) | 4 (14.3) | 8 (28.9) | | |
| Yes | 1 (3.6) | 1 (3.6) | 10 (35.7) | 0 (0.0) | 8 (28.9) | 20 (71.4) | | |
| **Currency** | | | | | | | 1.60 (*4*) | .81 |
| No | 0 (0.0) | 0 (0.0) | 1 (3.6) | 0 (0.0) | 0 (0.0) | 1 (3.6) | | |
| Yes | 1 (3.6) | 3 (10.7) | 10 (35.7) | 1 (3.6) | 12 (42.9) | 27 (96.4) | | |
| **Disclosure** | | | | | | | 15.36 (*4*) | .004 |
| No | 0 (0.0) | 3 (10.7) | 1 (3.6) | 1 (3.6) | 9 (32.1) | 14 (50.0) | | |
| Yes | 1 (3.6) | 0 (0.0) | 10 (35.7) | 0 (0.0) | 3 (10.7) | 14 (50.0) | | |

**Table 6.** Brief DISCERN scores by source type.

| | Source type | | | | | Average (SD) | *F* value (*df*) | *P* value |
|---|---|---|---|---|---|---|---|---|
| | Academic | Commercial | Government | Medical practice | Media outlet | | | |
| Brief DISCERN score, mean (SD) | 30.0 (0.0) | 17 (2.6) | 28.6 (1.4) | 18.0 (0.0) | 19.6 (5.6) | 23.2 (6.2) | 10.27 (4, 23) | <.001 |

## Information Quality

ANOVA revealed significant differences in mean Brief DISCERN scores by source type ($F_{4,23}$=10.27; $P$<.001), suggesting important differences in quality among the different source types. Post hoc analysis with Bonferroni correction revealed significant differences in Brief DISCERN scores between government and commercial sources ($P$=.002) and between government sources and media outlets ($P$<.001). Mean (SD) values of Brief DISCERN scores by source are provided in Table 6. Interrater agreement for our analyses was high (interclass correlation=0.96; 95% CI 0.95-0.97).

## Discussion

### Principal Findings

Using Google and its search analytics, we were able to identify the most frequently asked questions regarding COVID-19

vaccines in the United States. Google generated these FAQs by using millions of search queries nationwide. Additionally, we evaluated the assigned "answer" source for each FAQ, assessing each source's information transparency and quality. To our knowledge, this study is the first of its kind to evaluate the public's most frequently asked questions concerning the COVID-19 vaccines in the United States using Google search analytics. Our study is also the first of its kind to identify common answer sources used to address COVID-19 vaccine–related concerns and to assess their transparency and quality. In the following discourse, we discuss the importance of knowing COVID-19 FAQs in the context of the current COVID-19 vaccination campaigns while also providing recommendations for improving the public's confidence and willingness to be vaccinated.

## FAQs

The most popular COVID-19 vaccine–related questions sought factual information regarding safety and efficacy, indicating greater public concern regarding these topics. Consistent with our findings, survey studies found that safety and efficacy were among the most common COVID-19 vaccine concerns reported by the public and health care workers [37-40]. Additionally, studies have identified safety concerns as being one of the most common reasons for COVID-19 vaccine hesitancy [8,38-42]. In the United States, surveys indicate that 10% to 20% of adults and an estimated 8% of health care workers will refuse COVID-19 vaccines [8,37,39,43]. While the willingness to receive the COVID-19 vaccines has increased, the alarmingly high percentage of adults refusing vaccination creates a significant barrier to protecting our most vulnerable populations [43-45]. The potential cost of vaccine hesitancy and refusal in the United States is not exclusive to the COVID-19 pandemic. For example, an outbreak of measles virus, a pathogen for which vaccines effectively control outbreaks, occurred in Clark County, Washington, in 2019 [46]. Of 71 individuals involved, 61 (86%) were unvaccinated and 52 (73%) were children [46,47]. Moreover, vaccination rates in Clark County have been 10%-14% below the national average (88%) since 2013. The measles outbreak in 2019 was estimated to cost US $3.3 million to $3.5 million in labor, direct medical costs, and productivity losses [48]. It is likely that the cost of the Clark County measles outbreak could have been mitigated or reduced with adequate vaccination [47]. Thus, to prevent similar, but likely far worse, outcomes with COVID-19, effectively educating the public on the safety of COVID-19 vaccines is paramount for enhancing COVID-19 vaccine acceptance [49].

## Answer Sources

Overall, COVID-19 vaccine FAQs were most often answered by media outlets, followed by government sources. FAQs about safety and efficacy were answered more often by government sources, while media outlets frequently answered FAQs about technical details. The answer sources linked to each FAQ are found in "People also ask" or "Common concerns" boxes and are direct answers generated by Google [50]. These direct answers are supplied from Google's "trusted entities" database and are based on relational topics and machine learning [50]. While "trusted entities" seems rather vague, it appears that Google considers direct answers to be "trusted" based on clarity, completeness, and the lack of excessive promotional jargon. With the public's trust and willingness to accept the vaccine being a key element in a successful vaccination campaign [44,51-53], it may be more appropriate for direct answers addressing COVID-19 vaccine FAQs to be based on scientific integrity, objectivity, and transparency.

## Transparency and Quality of the Answer Source

The FAQs with direct answers from government sources were more likely to meet 3 or more JAMA benchmark criteria, indicating that government answers were more transparent. Additionally, government and academic sources were found to be of significantly higher quality. While media outlets are unquestionably an important source of health information to the public, these findings suggest that government sources may

be better for addressing the public's COVID-19 vaccine concerns. Although media outlets had moderate transparency and quality, there are notable reasons to use more reliable and objective sources. Generally, COVID-19 misinformation is rampant and the public opinion can be easily manipulated [29,45]. Indeed, media outlets are a frequent source of COVID-19 misinformation, and false claims are amplified by widespread news coverage [29,54]. For example, news stories early in the pandemic touting hydroxychloroquine as a "cure" perpetuated this misinformation in the absence of evidence [55]. More recently and more specifically related to the COVID-19 vaccines, rumors that COVID-19 vaccines cause infertility in women have circulated on social media [56]. Lastly, the politicization and polarization of news coverage surrounding the COVID-19 pandemic heavily influenced the public's attitude to COVID-19 response policies [55,57-60]. Taken together, trouble with media outlets as trustworthy sources further supports the use of unbiased answer sources such as government agencies.

## Recommendations

Above all, we recommend that individuals consider health care professionals as the primary source of information regarding COVID-19 vaccines. However, in cases where access to a health care professional is limited, web-based sources unquestionably present opportunities to quickly provide high-quality and accurate information regarding COVID-19 vaccines. We agree with Mills and Sivelä [61] that a successful COVID-19 vaccination campaign depends on gaining the public's trust in health care systems and government agencies, such as the Centers for Disease Control and Prevention and the World Health Organization, while also minimizing vaccine misinformation. Additionally, government sources must strive to translate scientifically dense literature into easily understandable information that answers widespread concerns. Therefore, the dissemination of this study's findings may promote the public's trust in these institutions as we have shown that government and academic sources provided the most transparent and highest-quality information addressing COVID-19 vaccine–related concerns.

Google recently demonstrated their willingness to support these COVID-19 vaccination campaigns by collaborating with Ohio State University to combat COVID-19 misinformation [62]. This partnership aims to ensure that people receive accurate information about COVID-19 vaccines to increase the public's confidence and willingness to be vaccinated. Thus, in alignment with Google's current intentions, we recommend that all COVID-19 vaccine FAQs be linked to government and academic answer sources; this would provide people with transparent and quality vaccine information. At a minimum, FAQs on safety and efficacy should be answered by government sources, as safety and efficacy concerns are among the primary drivers of COVID-19 vaccine hesitancy [39-42].

## Strength and Limitations

Our study's primary strength is the incorporation of Google FAQs as a novel source of insight regarding millions of individual inquiries about COVID-19 vaccines, which is an application of methodology adapted from the published literature

[21,26-28,34,35] and improved upon herein. Using FAQs generated by Google to explore the content of concerns regarding COVID-19 vaccine safety and efficacy may prevent common limitations of survey studies such as low response rates, reporting biases, and selection bias. Additionally, Google's large data set is continuously analyzed in real time and may offer improved and more specific targets when approaching the public's medical concerns. All classifications and assessments were performed in a masked duplicate fashion in accordance with standards set by the Cochrane Review and experts in the meta-research field [63,64] with high interrater reliability between investigators.

Our study is not without limitations though, such as those due to the dynamic nature of Google's search outputs. As searching for COVID-19 vaccine–related information continues, new and updated FAQs will be generated, limiting the generalizability of our study to the time when our search was performed. Additionally, the transparency and quality assessments we used do not check for information accuracy, as this would require source-by-source comparison to generally accepted truths regarding COVID-19 vaccines, rendering our assessments as gauges of information transparency and not of information accuracy. Lastly, the categorizing of FAQs and answer sources was limited owing to their subjectivity. Although the categories were developed in line with previous reports and had high interobserver reliability, there is still potential for overlap between categories.

## Conclusions

The expedient development and approval of COVID-19 vaccines is the culmination of the world's greatest scientific achievements; however, without positive public reception and adequate counseling and education, COVID-19 vaccination efforts may be hindered. Using Google allowed us to obtain a list of FAQs based on millions of searches for content related to COVID-19 vaccines, which reflected widespread and common concerns. We found that the most common COVID-19 vaccine–related questions pertained to vaccine safety and efficacy, which is supported by the findings of survey studies. We found that government and academic sources provided the most transparent and highest-quality web-based information for answering the public's most frequently asked questions about COVID-19 vaccines. Recognizing common concerns about COVID-19 vaccines may better assist health care professionals, researchers, and government agencies in improving vaccination efforts. Ensuring a successful vaccination campaign requires the public's trust, which may be enhanced through the availability of high-quality and transparent COVID-19 vaccine information, such as that provided by government sources.

## Conflicts of Interest

## References

1. Dong E, Du H, Gardner L. An interactive web-based dashboard to track COVID-19 in real time. Lancet Infect Dis 2020 May;20(5):533-534 [FREE Full text] [doi: 10.1016/S1473-3099(20)30120-1] [Medline: 32087114]
2. COVID-19 Dashboard. Johns Hopkins Coronavirus Resource Center. URL: https://coronavirus.jhu.edu/map.html [accessed 2021-08-01]
3. Pfizer-BioNTech COVID-19 Vaccine. US Food and Drug Aminidtration. URL: https://www.fda.gov/emergency-preparedness-and-response/coronavirus-disease-2019-covid-19/pfizer-biontech-covid-19-vaccine [accessed 2021-01-21]
4. Moderna COVID-19 Vaccine. US Food and Drug Administration. URL: https://www.fda.gov/emergency-preparedness-and-response/coronavirus-disease-2019-covid-19/moderna-covid-19-vaccine [accessed 2021-01-21]
5. MacDonald NE, SAGE Working Group on Vaccine Hesitancy. Vaccine hesitancy: Definition, scope and determinants. Vaccine 2015 Aug 14;33(34):4161-4164 [FREE Full text] [doi: 10.1016/j.vaccine.2015.04.036] [Medline: 25896383]
6. Kennedy J. Vaccine Hesitancy: A Growing Concern. Paediatr Drugs 2020 Apr;22(2):105-111. [doi: 10.1007/s40272-020-00385-4] [Medline: 32072472]
7. Lazarus JV, Ratzan SC, Palayew A, Gostin LO, Larson HJ, Rabin K, et al. Author Correction: A global survey of potential acceptance of a COVID-19 vaccine. Nat Med 2021 Feb;27(2):354 [FREE Full text] [doi: 10.1038/s41591-020-01226-0] [Medline: 33432176]
8. File T, Mohanty A. Around Half of Unvaccinated Americans Indicate They Will "Definitely" Get COVID-19 Vaccine. United States Census Bureau. 2021 Jan 27. URL: https://www.census.gov/library/stories/2021/01/around-half-of-unvaccinated-americans-indicate-they-will-definitely-get-covid-19-vaccine.html [accessed 2021-02-25]
9. Kelly BJ, Southwell BG, McCormack LA, Bann CM, MacDonald PDM, Frasier AM, et al. Predictors of willingness to get a COVID-19 vaccine in the U.S. BMC Infect Dis 2021 Apr 12;21(1):338 [FREE Full text] [doi: 10.1186/s12879-021-06023-9] [Medline: 33845781]
10. Shen SC, Dubey V. Addressing vaccine hesitancy: Clinical guidance for primary care physicians working with parents. Can Fam Physician 2019 Mar;65(3):175-181 [FREE Full text] [Medline: 30867173]
11. Bailey SR. Physicians provide key voice in building vaccine confidence. American Medical Association. 2021 Feb 25. URL: https://tinyurl.com/3jzer8s7 [accessed 2021-01-27]

12.  Gualtieri L. The doctor as the second opinion and the internet as the first. In: CHI '09 Extended Abstracts on Human Factors in Computing Systems. 2009 Presented at: CHI '09: CHI Conference on Human Factors in Computing Systems; April 4-9, 2009; Boston, MA p. 2489-2498 URL: https://dl.acm.org/doi/abs/10.1145/1520340.1520352 [doi: 10.1145/1520340.1520352]

13.  Cohen RA, Adams PF. Use of the internet for health information: United States, 2009. NCHS Data Brief 2011 Jul(66):1-8 [FREE Full text] [Medline: 22142942]

14.  Tan SS, Goonawardene N. Internet Health Information Seeking and the Patient-Physician Relationship: A Systematic Review. J Med Internet Res 2017 Jan 19;19(1):e9 [FREE Full text] [doi: 10.2196/jmir.5729] [Medline: 28104579]

15.  Largent EA, Persad G, Sangenito S, Glickman A, Boyle C, Emanuel EJ. US Public Attitudes Toward COVID-19 Vaccine Mandates. JAMA Netw Open 2020 Dec 01;3(12):e2033324 [FREE Full text] [doi: 10.1001/jamanetworkopen.2020.33324] [Medline: 33337490]

16.  Kanthawala S, Vermeesch A, Given B, Huh J. Answers to Health Questions: Internet Search Results Versus Online Health Community Responses. J Med Internet Res 2016 Apr 28;18(4):e95 [FREE Full text] [doi: 10.2196/jmir.5369] [Medline: 27125622]

17.  Schachinger K. A Complete Guide to the Google RankBrain Algorithm. Search Engine J 2017 [FREE Full text]

18.  Devlin J, Chang M, Lee K, Toutanova K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. arXiv. Preprint posted online May 24, 2019 [FREE Full text]

19.  How Search algorithms work. Google Search. URL: https://www.google.com/search/howsearchworks/algorithms/ [accessed 2021-04-21]

20.  Nayak P. Understanding searches better than ever before. The Keyword. 2019. URL: https://blog.google/products/search/search-language-understanding-bert/ [accessed 2021-07-27]

21.  Shen TS, Driscoll DA, Islam W, Bovonratwet P, Haas SB, Su EP. Modern Internet Search Analytics and Total Joint Arthroplasty: What Are Patients Asking and Reading Online? J Arthroplasty 2021 Apr;36(4):1224-1231 [FREE Full text] [doi: 10.1016/j.arth.2020.10.024] [Medline: 33162279]

22.  Google. URL: http://google.com [accessed 2021-03-12]

23.  Rothwell J. In Mixed Company: Communicating in Small Groups. Boston, MA: Wadsworth Publishing; 2012.

24.  Starman JS, Gettys FK, Capo JA, Fleischli JE, Norton HJ, Karunakar MA. Quality and content of Internet-based information for ten common orthopaedic sports medicine diagnoses. J Bone Joint Surg Am 2010 Jul 07;92(7):1612-1618. [doi: 10.2106/JBJS.I.00821] [Medline: 20595567]

25.  Silberg WM, Lundberg GD, Musacchio RA. Assessing, controlling, and assuring the quality of medical information on the Internet: Caveant lector et viewor--Let the reader and viewer beware. JAMA 1997 Apr 16;277(15):1244-1245. [Medline: 9103351]

26.  Cassidy JT, Baker JF. Orthopaedic Patient Information on the World Wide Web: An Essential Review. J Bone Joint Surg Am 2016 Feb 17;98(4):325-338. [doi: 10.2106/JBJS.N.01189] [Medline: 26888683]

27.  Kartal A, Kebudi A. Evaluation of the Reliability, Utility, and Quality of Information Used in Total Extraperitoneal Procedure for Inguinal Hernia Repair Videos Shared on WebSurg. Cureus 2019 Sep 04;11(9):e5566 [FREE Full text] [doi: 10.7759/cureus.5566] [Medline: 31695985]

28.  Corcelles R, Daigle CR, Talamas HR, Brethauer SA, Schauer PR. Assessment of the quality of Internet information on sleeve gastrectomy. Surg Obes Relat Dis 2015;11(3):539-544. [doi: 10.1016/j.soard.2014.08.014] [Medline: 25604832]

29.  Cuan-Baltazar JY, Muñoz-Perez MJ, Robledo-Vega C, Pérez-Zepeda MF, Soto-Vega E. Misinformation of COVID-19 on the Internet: Infodemiology Study. JMIR Public Health Surveill 2020 Apr 09;6(2):e18444 [FREE Full text] [doi: 10.2196/18444] [Medline: 32250960]

30.  Charnock D, Shepperd S, Needham G, Gann R. DISCERN: an instrument for judging the quality of written consumer health information on treatment choices. J Epidemiol Community Health 1999 Feb;53(2):105-111 [FREE Full text] [doi: 10.1136/jech.53.2.105] [Medline: 10396471]

31.  Fan KS, Ghani SA, Machairas N, Lenti L, Fan KH, Richardson D, et al. COVID-19 prevention and treatment information on the internet: a systematic analysis and quality assessment. BMJ Open 2020 Sep 10;10(9):e040487 [FREE Full text] [doi: 10.1136/bmjopen-2020-040487] [Medline: 32912996]

32.  Haragan AF, Zuwiala CA, Himes KP. Online Information About Periviable Birth: Quality Assessment. JMIR Pediatr Parent 2019 Jun 07;2(1):e12524 [FREE Full text] [doi: 10.2196/12524] [Medline: 31518325]

33.  Azer SA, AlOlayan TI, AlGhamdi MA, AlSanea MA. Inflammatory bowel disease: An evaluation of health information on the internet. World J Gastroenterol 2017 Mar 07;23(9):1676-1696 [FREE Full text] [doi: 10.3748/wjg.v23.i9.1676] [Medline: 28321169]

34.  Khazaal Y, Chatton A, Cochand S, Coquard O, Fernandez S, Khan R, et al. Brief DISCERN, six questions for the evaluation of evidence-based content of health-related websites. Patient Educ Couns 2009 Oct;77(1):33-37. [doi: 10.1016/j.pec.2009.02.016] [Medline: 19372023]

35.  Banasiak NC, Meadows-Oliver M. Evaluating asthma websites using the Brief DISCERN instrument. J Asthma Allergy 2017;10:191-196 [FREE Full text] [doi: 10.2147/JAA.S133536] [Medline: 28670135]

36. Zheluk A, Maddock J. Plausibility of Using a Checklist With YouTube to Facilitate the Discovery of Acute Low Back Pain Self-Management Content: Exploratory Study. JMIR Form Res 2020 Nov 20;4(11):e23366 [FREE Full text] [doi: 10.2196/23366] [Medline: 33216003]

37. Shekhar R, Sheikh AB, Upadhyay S, Singh M, Kottewar S, Mir H, et al. COVID-19 Vaccine Acceptance among Health Care Workers in the United States. Vaccines (Basel) 2021 Feb 03;9(2):119 [FREE Full text] [doi: 10.3390/vaccines9020119] [Medline: 33546165]

38. Akarsu B, Canbay Özdemir D, Ayhan Baser D, Aksoy H, Fidancı İ, Cankurtaran M. While studies on COVID-19 vaccine is ongoing, the public's thoughts and attitudes to the future COVID-19 vaccine. Int J Clin Pract 2021 Apr;75(4):e13891 [FREE Full text] [doi: 10.1111/ijcp.13891] [Medline: 33278857]

39. Fisher KA, Bloomstone SJ, Walder J, Crawford S, Fouayzi H, Mazor KM. Attitudes Toward a Potential SARS-CoV-2 Vaccine : A Survey of U.S. Adults. Ann Intern Med 2020 Dec 15;173(12):964-973 [FREE Full text] [doi: 10.7326/M20-3569] [Medline: 32886525]

40. Verger P, Scronias D, Dauby N, Adedzi KA, Gobert C, Bergeat M, et al. Attitudes of healthcare workers towards COVID-19 vaccination: a survey in France and French-speaking parts of Belgium and Canada, 2020. Euro Surveill 2021 Jan;26(3):2002047 [FREE Full text] [doi: 10.2807/1560-7917.ES.2021.26.3.2002047] [Medline: 33478623]

41. Wang K, Wong EL, Ho K, Cheung AW, Yau PS, Dong D, et al. Change of Willingness to Accept COVID-19 Vaccine and Reasons of Vaccine Hesitancy of Working People at Different Waves of Local Epidemic in Hong Kong, China: Repeated Cross-Sectional Surveys. Vaccines (Basel) 2021 Jan 18;9(1):62 [FREE Full text] [doi: 10.3390/vaccines9010062] [Medline: 33477725]

42. Su Z, McDonnell D, Cheshmehzangi A, Li X, Maestro D, Šegalo S, et al. With Great Hopes Come Great Expectations: A Viewpoint on Access and Adoption Issues Associated with COVID-19 Vaccines. JMIR Public Health Surveill 2021 Feb 01 [FREE Full text] [doi: 10.2196/26111] [Medline: 33560997]

43. Funk C, Tyson A. Intent to Get a COVID-19 Vaccine Rises to 60% as Confidence in Research and Development Process Increases. Pew Research Center. 2020 Dec 03. URL: https://www.pewresearch.org/science/2020/12/03/intent-to-get-a-covid-19-vaccine-rises-to-60-as-confidence-in-research-and-development-process-increases/ [accessed 2021-02-13]

44. Schaffer DeRoo S, Pudalov NJ, Fu LY. Planning for a COVID-19 Vaccination Program. JAMA 2020 Jun 23;323(24):2458-2459. [doi: 10.1001/jama.2020.8711] [Medline: 32421155]

45. Nahum A, Drekonja D, Alpern J. The Erosion of Public Trust and SARS-CoV-2 Vaccines- More Action Is Needed. Open Forum Infect Dis 2021 Feb;8(2):ofaa657 [FREE Full text] [doi: 10.1093/ofid/ofaa657] [Medline: 34141815]

46. Carlson A, Riethman M, Gastañaduy P, Lee A, Leung J, Holshue M, et al. Notes from the Field: Community Outbreak of Measles - Clark County, Washington, 2018-2019. MMWR Morb Mortal Wkly Rep 2019 May 17;68(19):446-447 [FREE Full text] [doi: 10.15585/mmwr.mm6819a5] [Medline: 31095534]

47. Porter A, Goldfarb J. Measles: A dangerous vaccine-preventable disease returns. Cleve Clin J Med 2019 Jun;86(6):393-398 [FREE Full text] [doi: 10.3949/ccjm.86a.19065] [Medline: 31204978]

48. Pike J, Melnick A, Gastañaduy PA, Kay M, Harbison J, Leidner AJ, et al. Societal Costs of a Measles Outbreak. Pediatrics 2021 Apr;147(4):e2020027037. [doi: 10.1542/peds.2020-027037] [Medline: 33712549]

49. Answering Patients' Questions About COVID-19 Vaccine and Vaccination. Centers for Disease Control and Prevention. URL: https://www.cdc.gov/vaccines/covid-19/hcp/answering-questions.html [accessed 2021-03-09]

50. Hill J. People Also Ask Boxes and Related Questions. Hill Web Creations. 2020. URL: https://www.hillwebcreations.com/people-also-ask-related-questions/ [accessed 2021-02-14]

51. COCONEL Group. A future vaccination campaign against COVID-19 at risk of vaccine hesitancy and politicisation. Lancet Infect Dis 2020 Jul;20(7):769-770 [FREE Full text] [doi: 10.1016/S1473-3099(20)30426-6] [Medline: 32445713]

52. Schwartz JL. Evaluating and Deploying Covid-19 Vaccines - The Importance of Transparency, Scientific Integrity, and Public Trust. N Engl J Med 2020 Oct 29;383(18):1703-1705. [doi: 10.1056/NEJMp2026393] [Medline: 32966716]

53. Trogen B, Oshinsky D, Caplan A. Adverse Consequences of Rushing a SARS-CoV-2 Vaccine: Implications for Public Trust. JAMA 2020 Jun 23;323(24):2460-2461. [doi: 10.1001/jama.2020.8917] [Medline: 32453392]

54. Evanega S, Lynas M, Adams J, Smolenyak K. Coronavirus misinformation: quantifying sources and themes in the COVID-19 'infodemic'. JMIR Preprints 2020 [FREE Full text] [doi: 10.2196/preprints.25143]

55. Motta M, Stecula D, Farhart C. How Right-Leaning Media Coverage of COVID-19 Facilitated the Spread of Misinformation in the Early Stages of the Pandemic in the U.S. Can J Pol Sci 2020 May 01;53(2):335-342 [FREE Full text] [doi: 10.1017/S0008423920000396]

56. Rodriguez A. No, the COVID-19 vaccine doesn't cause infertility in women. USA Today. 2020. URL: https://www.usatoday.com/story/news/health/2020/12/10/covid-vaccine-debunking-claims-causes-infertility-sterilization/6497018002/ [accessed 2021-02-04]

57. Political polarisation impedes the public policy response to COVID-19. VoxEU & CEPR. 2020. URL: https://voxeu.org/article/political-polarisation-impedes-public-policy-response-covid-19 [accessed 2021-02-11]

58. Hart PS, Chinn S, Soroka S. Politicization and Polarization in COVID-19 News Coverage. Science Communication 2020 Aug 25;42(5):679-697 [FREE Full text] [doi: 10.1177/1075547020950735]

59.  Greiner B, Ottwell R, Vassar M, Hartwell M. Public Interest in Preventive Measures of Coronavirus Disease 2019 Associated With Timely Issuance of Statewide Stay-at-Home Orders. Disaster Med Public Health Prep 2020 Dec;14(6):765-768 [FREE Full text] [doi: 10.1017/dmp.2020.189] [Medline: 32498752]

60.  Hartwell M, Greiner B, Kilburn Z, Ottwell R. Association of Public Interest in Preventive Measures and Increased COVID-19 Cases After the Expiration of Stay-at-Home Orders: A Cross-Sectional Study. Disaster Med Public Health Prep 2020 Sep 10:1-5 [FREE Full text] [doi: 10.1017/dmp.2020.333] [Medline: 32907675]

61.  Mills MC, Sivelä J. Should spreading anti-vaccine misinformation be criminalised? BMJ 2021 Feb 17;372:n272. [doi: 10.1136/bmj.n272] [Medline: 33597153]

62.  Foresman B. Google, Ohio State join to fight COVID-19 vaccine misinformation. EDSCOOP. 2021. URL: https://edscoop.com/google-ohio-state-join-to-fight-covid-19-vaccine-misinformation/ [accessed 2021-02-20]

63.  Higgins J, Thomas J, Chandler J. Cochrane Handbook for Systematic Reviews of Interventions. Hoboken, NJ: John Wiley & Sons; 2019. URL: https://play.google.com/store/books/details?id=cTqyDwAAQBAJ [accessed 2020-05-03]

64.  Mbuagbaw L, Lawson DO, Puljak L, Allison DB, Thabane L. A tutorial on methodological studies: the what, when, how and why. BMC Med Res Methodol 2020 Sep 07;20(1):226 [FREE Full text] [doi: 10.1186/s12874-020-01107-7] [Medline: 32894052]

## Abbreviations

**FAQ:** frequently asked question
**JAMA:** Journal of the American Medical Association

Original Paper

# The Impact of COVID-19 on Conspiracy Hypotheses and Risk Perception in Italy: Infodemiological Survey Study Using Google Trends

Alessandro Rovetta[1], SRSCI

Mensana srls, Brescia, Italy

**Corresponding Author:**
Alessandro Rovetta, SRSCI
Mensana srls
Via Malta 12
Brescia, 25124
Italy
Phone: 39 3927112808
Email: rovetta.mresearch@gmail.com

## Abstract

**Background:** COVID-19 has caused the worst international crisis since World War II. Italy was one of the countries most affected by both the pandemic and the related infodemic. The success of anti–COVID-19 strategies and future public health policies in Italy cannot separate itself from the containment of fake news and the divulgation of correct information.

**Objective:** The aim of this paper was to analyze the impact of COVID-19 on web interest in conspiracy hypotheses and risk perception of Italian web users.

**Methods:** Google Trends was used to monitor users' web interest in specific topics, such as conspiracy hypotheses, vaccine side effects, and pollution and climate change. The keywords adopted to represent these topics were mined from Bufale.net—an Italian website specializing in detecting online hoaxes—and Google Trends suggestions (ie, related topics and related queries). Relative search volumes (RSVs) of the time-lapse periods of 2016-2020 (pre–COVID-19) and 2020-2021 (post–COVID-19) were compared through percentage difference ($\Delta_\%$) and the Welch $t$ test ($t$). When data series were not stationary, other ad hoc criteria were used. The trend slopes were assessed through Sen slope (SS). The significance thresholds have been indicatively set at $P=.05$ and $t=1.9$.

**Results:** The COVID-19 pandemic drastically increased Italian netizens' interest in conspiracies ($\Delta_\% \in [60, 288]$, $t \in [6, 12]$). Web interest in conspiracy-related queries across Italian regions increased and became more homogeneous compared to the pre–COVID-19 period (average RSV=80±2.8, $t_{min}=1.8$, $\Delta_{min\%}=+12.4$, $min\Delta_{SD\%}=-25.8$). In addition, a growing trend in web interest in the infodemic YouTube channel ByoBlu has been highlighted. Web interest in hoaxes has increased more than interest in antihoax services ($t_1=11.3$ vs $t_2=4.5$; $\Delta_{1\%}=+157.6$ vs $\Delta_{2\%}=+84.7$). Equivalently, web interest in vaccine side effects exceeded interest in pollution and climate change ($SS_{vaccines}=0.22$, $P<.001$ vs $SS_{pollution}=0.05$, $P<.001$; $\Delta_\%=+296.4$). To date, a significant amount of fake news related to COVID-19 vaccines, unproven remedies, and origin has continued to circulate. In particular, the creation of SARS-CoV-2 in a Chinese laboratory constituted about 0.04% of the entire web interest in the pandemic.

**Conclusions:** COVID-19 has given a significant boost to web interest in conspiracy hypotheses and has made it more uniform across regions in Italy. The pandemic accelerated an already-growing trend in users' interest toward some fake news sources, including the 500,000-subscriber YouTube channel ByoBlu, which was removed from the platform by YouTube for disinformation in March 2021. The risk perception related to COVID-19 vaccines has been so distorted that vaccine side effect–related queries outweighed those relating to pollution and climate change, which are much more urgent issues. Moreover, a large amount of fake news has circulated about COVID-19 vaccines, remedies, and origin. Based on these findings, it is recommended that the Italian authorities implement more effective infoveillance systems, and that communication by the mass media be less sensationalistic and more consistent with the available scientific evidence. In this context, Google Trends can be used to monitor users' response to specific infodemiological countermeasures. Further research is needed to understand the psychological mechanisms that regulate risk perception.

XSL•FO
**RenderX**

## KEYWORDS

COVID-19; fake news; Google Trends; infodemiology; Italy; risk perception

## Introduction

COVID-19 was responsible for one of the most dramatic global crises after World War II. As of April 24, 2021, the official global toll was 144 million cases and 3.1 million deaths [1]. Such a pandemic has also triggered a vast infodemic, capable of seriously damaging the economic and health systems of many countries as well as enabling the spread of the novel coronavirus itself [2]. Specifically, an infodemic is defined as an excessive amount of unfiltered information concerning a problem, such that the solution is made more difficult [3]. However, it is not the first time that the world has been forced to face a vast infodemic; for example, during the HCoV-EMC/2012 (human coronavirus–Erasmus Medical Center/2012) epidemic generated by a previous coronavirus, some flawed denominations, such as "Middle East respiratory syndrome" and "swine flu," have caused unintentional adverse social and economic impacts by stigmatizing industries and communities [4]. In addition, the adoption of improper names has also led to medical and nursing errors concerning drug administration [5]. To deal with this growing problem, fomented by increasingly rapid mass media such as that provided by the internet, Dr Gunther Eysenbach has devised a scientific branch called "infodemiology," which encompasses all the techniques for monitoring and analyzing information [3]. In general, the infodemiological approach is based on the collection of information circulating in a network—not necessarily online—with the following possible purposes: (1) investigate the mental and physical health of a group or community, (2) identify the dangers and extent of disinformation or misinformation regarding a specific topic, and (3) carry out assessments in the field of public health (eg, using web searches to obtain information about symptoms and spread of a disease).

As reported by the World Health Organization (WHO), the COVID-19 infodemic can intensify or lengthen outbreaks. For this reason, a huge infodemiological effort has been made to study the information circulating on the web and contain the spread of fake news [6]. In this context, 132 nations worldwide signed a document to guarantee their commitment to the battle against disinformation and misinformation [7]. On the operational level, infodemic management takes place through four key steps: (1) listening to community concerns and questions, (2) promoting understanding of risk and health expert advice, (3) building resilience to misinformation, and (4) engaging and empowering communities to take positive action [2]. This paper focuses on points 2 and 3 as concerns Italy, one of the nations most affected by COVID-19 [1]. The objective is to analyze and quantify the impact of COVID-19 on Italian netizens' risk perception and new and pre-existing conspiracy hypotheses through Google Trends, an infoveillance tool provided by Google that returns users' web interest in specific topics in the form of normalized values called relative search volumes (RSVs) [8]. In this regard, pre-existing conspiracies are defined as those conspiracies that existed even before COVID-19 and are not directly related to it. The denomination

"conspiracy hypotheses" aims to underline the absence of the scientific background necessary to call them theories. Google Trends has been exploited extensively in the scientific community to conduct infodemiological, medical, psychological, economic, and even epidemiological studies [9-14]. Indeed, although the media can influence users' web searches [15], Google Trends provides valuable details on the dynamics of users' online interests, including the influence of the media on collective thinking [16,17].

As of April 2021, the success of the vaccination campaign has been crucial in the fight against COVID-19 [18,19]. Conspiracy hypotheses, inadequate risk perception, and unjustified fears have already undermined nonpharmacological containment measures and can reduce the effectiveness of pharmacological ones [20]. Furthermore, these factors can compromise the management of equally serious problems such as pollution and climate change, which are often linked to COVID-19 incidence and mortality [21-23]. Therefore, such a scenario requires careful surveillance of the online information flow as well as thoughtful communication. As we will show in this research, Google Trends can help achieve this goal.

## Methods

### Data Collection

#### Overview

For each topic, appropriate keywords were selected according to the methods explained in the following subsections. Each keyword was searched on Google Trends under the category "all categories." The time-lapse period was set to 5 years (April 21, 2016, to April 21, 2021). Only the most relevant queries were included in the results (ie, ). The selection of the queries with the highest RSVs was conducted by consulting related topics and related queries provided by Google Trends. In this way, it was possible to select the most relevant queries, including those containing typos. All keywords were collected for at least 7 consecutive days in order to highlight potential anomalies and significant variations [24].

#### Pre-existing Fake News

Pre-existing fake news and disinformation channels were mined from the specialized antihoax website Bufale.net [25]. The selection of the keywords to search on Google Trends took place through the following steps: (1) consultation of the blacklisted infodemic sources drawn up by the authors of the Bufale.net website, (2) search of all the aforementioned infodemic sources on Google Trends, and (3) selection of infodemic sources with . By doing so, four keywords that represent the main conspiracy-related web interests on Google Trends were identified: "cospirazione + nuovo ordine mondiale + complotto" (conspiracy + new world order + plot), "byoblu" (a 500,000-subscriber YouTube channel removed in March 2021 for disinformation), "Maurizio Blondet" (an Italian journalist who supports conspiracy hypotheses), and "luogocomune" (a

Facebook page sharing conspiracy hypotheses). All of these keywords have been independently searched on the web to verify the actual presence of hoaxes and fake news, understood as information that conflicts with current scientific literature. The details of this examination are reported in Multimedia Appendix 1.

### Risk Perception

RSVs of the query "fake news + bufale + notizie false" (fake news + hoaxes + false news) and the previous queries (ie, pre-existing fake news) were compared. By doing so, it was possible to observe the impact of the pandemic on web interest in antihoax services and the hoaxes themselves. The same procedure was carried out for the queries "vaccini effetti collaterali + vaccino effetti collaterali" (vaccine side effects + vaccines side effects) and "inquinamento + cambiamento climatico" (pollution + climate change) queries. In this way, it was possible to evaluate web interest in two very distant topics in terms of health risk and incidence [26-28].

### COVID-19–Related Fake News

To monitor the trend of fake news in Italy after more than a year of the pandemic, we referred to the following: (1) previous studies conducted during both the first and second waves of COVID-19 in Italy [20,29], (2) the Bufale.net website [25], and (3) the official website of the Italian Ministry of Health [30].

The keywords that reached an  concerned the following topics: the creation of COVID-19 in a Chinese laboratory ("coronavirus laboratorio + covid laboratorio – analisi – tampone – tamponi"), vaccine plot ("vaccino calamita + vaccino chip + vaccino microchip + vaccino bill gates"), 5G plot ("coronavirus 5g + covid 5g + corona 5g + virus 5g"), COVID-19 plot ("complotto coronavirus + complotto covid + complotto pandemia + grande reset"), and unproven remedies ("coronavirus vitamina + covid vitamina + coronavirus aglio + covid aglio").

## Statistical Analysis

### Welch t Test

The Welch $t$ test was used independently of the data set distribution, thanks to the central limit theorem (N>30, where N the is the number of measures). Nevertheless, a qualitative graphic control was performed to confirm the absence of too-pronounced skewness. The difference between the two mean values was considered significant indicatively when $t>1.9$.

### Percentage Change

The percentage change, $\Delta_\%$, was calculated through the formula $[y(T_2) - x(T_1)]/x(T_1) \times 100$, where $T_i$ is a specific time-lapse period.

### Shapiro-Wilk Test

The Shapiro-Wilk test was used, together with a qualitative graphic control, to evaluate the distributive normality of the data set in question.

### Mean Values

All mean values were calculated using the standard arithmetic mean and are presented as mean (SEM [standard error of the mean]). When N<30, the Shapiro-Wilk test was performed to assess the goodness of the mean value as a statistical measure.

### Data Series Analysis

All data series were graphed. To signal the presence or absence of trends, augmented Dickey-Fuller (ADF), Mann-Kendall (MK), and Sen slope (SS) tests were adopted. The same tests were used to evaluate the data sets' stationarity. Calculations were performed with Microsoft Excel 2021 software through the Real-Statistics 2021 package (Multimedia Appendix 2). The optimal lag was determined using the Schwert criterion.

### Data Series Comparison

To estimate the effect of COVID-19 on web queries, RSV trends over the last 5 years (April 21, 2016, to April 21, 2021) were analyzed. As shown in a previous paper, Italian netizens showed a marked interest in the COVID-19 pandemic only when it became a direct national problem [31]. Therefore, the time-lapse periods of "April 21, 2016, to February 16, 2020" (period 1) and "February 16, 2020, to April 21, 2021" (period 2) were compared. When period 1 turned out to be stationary or contained a negative trend, $t$ and $\Delta_\%$ were calculated. When period 1 contained a stationary positive trend, the trend slopes of period 1 and a specific subperiod of "February 16, 2020, to $x$" of period 2 were compared by $\Delta_\%$; such a subperiod was selected by observing the region of the graph in which a possible positive level-shift occurred. Period 1 data were then linearly, quadratically, or sigmoidally interpolated, depending on which monotone function minimized the statistical errors. Period 2 data were interpolated through a polynomial function of the 9th degree. Finally, $\Delta_\%$ was calculated between the areas subtended by the two curves after February 16, 2020 ($\Delta A_\%$). These were calculated using a definite integral between weeks 201 and 264. When period 1 contained a positive level-shift but was piecewise quasi-stationary, $t$ and $\Delta_\%$ were calculated considering only the last quasi-stationary subperiod.

### Pearson Correlation

Pearson correlation ($r$) was used only after verifying the distributive normality of the data set through the Shapiro-Wilk test plus a graphical check. No strength thresholds have been adopted. Since all of the samples in which the Pearson correlation was calculated were sufficiently Gaussian, nonparametric correlations were not exploited.

### P Values

Two-tailed $P$ values were used as graded measures of the strength of evidence against the null hypothesis. An indicative threshold has been set at $P=.05$; however, exact $P$ values for the ADF and MK+SS tests are reported in Multimedia Appendix 2. The remaining $P$ values are reported in full in this manuscript.
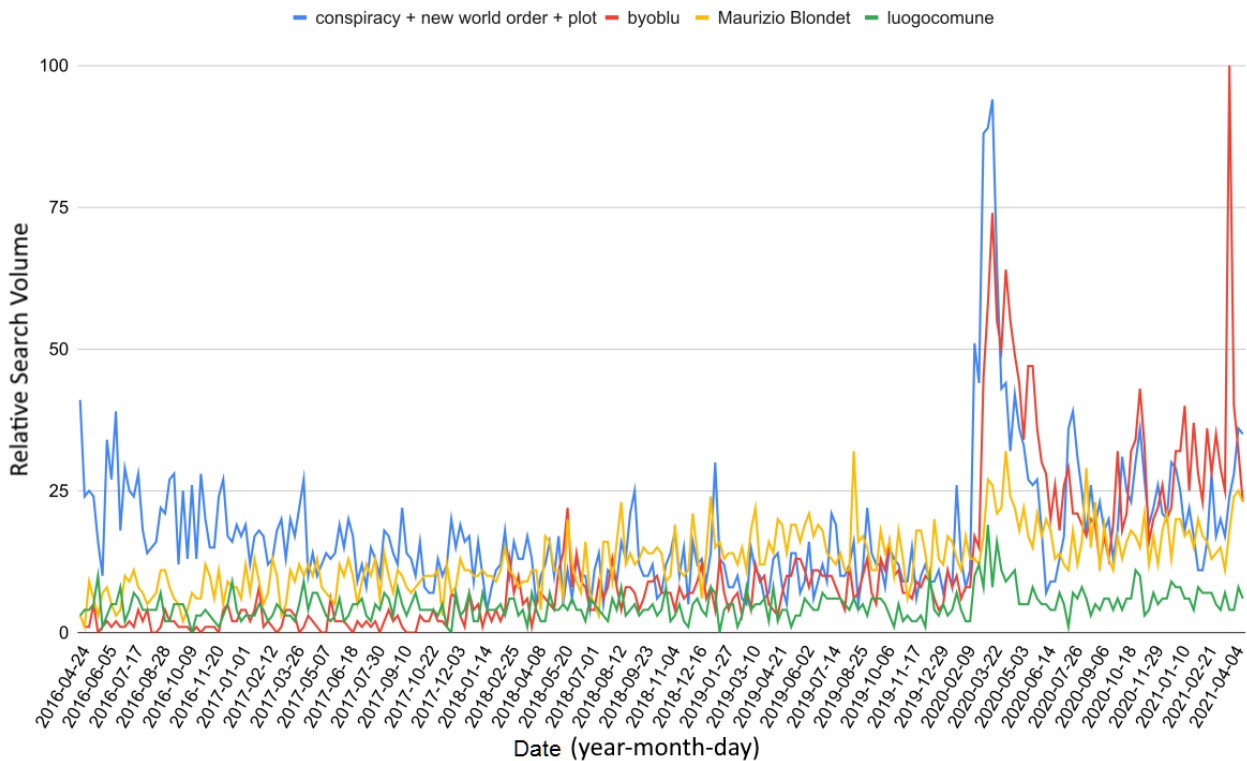
## Results

### Web Interest in Pre-existing Fake News

The impact of COVID-19 on the RSVs of conspiracy-related queries was evident (Figure 1); in particular, all considered infodemic queries underwent a significant level-shift during the Codogno, Italy, outbreak at the end of February 2020, signaling

an immediate increase in fake news with the arrival of the pandemic in Italy ($\Delta_{1\%}$=+102.5, $t_1$=6.3; $\Delta_{2\%}$=+288.2, $t_2$=11.5; $\Delta A_{3\%}$=+10.2; $\Delta_{4\%}$=+60.6, $t_4$=6.2).

It is relevant that the COVID-19 national outbreak has influenced the RSV trend even after the end of the first lockdown (May 2020). Indeed, Figure 1 shows a permanent level-shift for all of the investigated web queries.

**Figure 1.** Relative search volumes of conspiracy-related web search queries from April 21, 2016, to April 21, 2021, in Italy.



Regional interest in conspiracy-related keyword 1 decreased, on average, during the 2016-2020 time-lapse period, from 71 (SEM 4.2) to 52 (SEM 4.1) (Table 1). From 2020 to 2021, the increase in RSV was manifest and common to all regions ($\boxed{\times}$=80 ± 2.8, $t_{min}$=1.8, $\Delta_{min\%}$=+12.4; $t_{max}$=5.7, $\Delta_{max\%}$=+54.6). Over the 2016-2020 period, the percentage standard deviation ranged from 22.1 to 32.6, while it reached an absolute minimum of 14.4 during the COVID-19 pandemic. This fact shows that web interest in conspiracies has also become more homogeneous across regions. No correlation was sought, as data were subject to a strong dependence on the day of collection [24]; however, the mean values and standard deviations never changed significantly ($t_{max}$=0.4).

**Table 1.** Relative search volumes (RSVs) on the web of the keyword "conspiracy + new world order + plot" from 2016 to 2021.

| Variable | Value | | | | |
|---|---|---|---|---|---|
| | 2016-2017 | 2017-2018 | 2018-2019 | 2019-2020 | 2020-2021 |
| **RSV for each Italian region** | | | | | |
| Abruzzo | 100 | 76 | 38 | 53 | 67 |
| Basilicata | N/A[a] | N/A | N/A | N/A | N/A |
| Calabria | 48 | 61 | 80 | 27 | 78 |
| Campania | 62 | 69 | 59 | 44 | 75 |
| Emilia-Romagna | 74 | 70 | 77 | 68 | 80 |
| Friuli-Venezia Giulia | 91 | 68 | 39 | 100 | 100 |
| Lazio | 64 | 85 | 74 | 51 | 84 |
| Liguria | 81 | 100 | 52 | 69 | 89 |
| Lombardia | 85 | 75 | 63 | 43 | 86 |
| Marche | 73 | 68 | 64 | 62 | 91 |
| Molise | N/A | N/A | N/A | N/A | N/A |
| Piemonte | 56 | 75 | 94 | 58 | 88 |
| Puglia | 70 | 78 | 45 | 49 | 67 |
| Sardegna | 78 | 36 | 100 | 41 | 95 |
| Sicilia | 80 | 77 | 62 | 43 | 69 |
| Toscana | 84 | 75 | 72 | 44 | 80 |
| Trentino-Alto Adige | 26 | 41 | 27 | 29 | 54 |
| Umbria | 70 | 82 | 87 | 49 | 77 |
| Valle d'Aosta | N/A | N/A | N/A | N/A | N/A |
| Veneto | 72 | 54 | 72 | 52 | 84 |
| **Other statistics** | | | | | |
| Mean (SD) | 71.4 (17.2) | 70.0 (15.5) | 65.0 (20.3) | 51.9 (16.9) | 80.2 (11.5) |
| SD% | 24.1 | 22.1 | 31.2 | 32.6 | 14.4 |
| SEM (standard error of the mean) | 4.2 | 3.8 | 4.9 | 4.1 | 2.8 |
| SEM% | 5.9 | 5.4 | 7.6 | 7.9 | 3.5 |
| Shapiro-Wilk $P$ value | <.001 | .18 | .97 | .05 | .93 |

[a]N/A: not applicable due to Google Trends detection anomalies.

On the contrary, web interest in the ByobBlu disinformation channel has always been compatible during the 2016-2021 period ($t \in$ [–1.3, 1.1], $\Delta_{\%} \in$ [–16.3, 21.8]; $t_{20-21}=-0.1$, $\Delta_{20-21}=-1.0$; Multimedia Appendix 2).

Figure 1 shows a substantial increase in national searches, and it is evident that the query "byoblu" was searched more in some regions than others. By analyzing the regional trends one by one, it can be observed that web interest in "byoblu" has increased over time, except for in Basilicata and Molise (Figure 2). Although a growing trend was already present, the novel coronavirus seems to have strongly impacted RSVs in Campania ($\Delta_{20-21\%}=+93.3$ vs $\Delta_{19-20\%}=+63.8$), Friuli-Venezia Giulia ($\Delta_{20-21\%}=+187.6$ vs $\Delta_{19-20\%}=+59.0$), Lazio ($\Delta_{20-21\%}=+95.7$ vs $\Delta_{19-20\%}=+15.8$), Trentino-Alto Adige ($\Delta_{20-21\%}=+173.2$ vs $\Delta_{19-20\%}=+4.5$, and Valle d'Aosta ($\Delta_{20-21\%}=+199.2$ vs $\Delta_{19-20\%}=+86.2$).

**Figure 2.** Web interest in the "byoblu" search query by each Italian region from 2016 to 2021.



## Risk Perception

A fraction of the users seemed aware of the danger inherent in COVID-19 fake news circulating on the web and tried to limit its effects by relying on antihoax websites (Figure 3). Nevertheless, the impact of COVID-19 was more incisive for conspiracy hypotheses ($t_1$=11.3 vs $t_2$=4.5; $\Delta_{1\%}$=+157.6 vs $\Delta_{2\%}$=+84.7), so much so that it is possible to observe a greater level-shift in the trend of infodemic web queries. This worsening was homogeneous on a national scale (Figure 4).

With the announcement of the discovery of COVID-19 vaccines, web interest in side effects immediately soared ($t_{27.1}$=10.3, $\Delta_\%$=+905.2 from October 2020 to April 2021). Although health authorities reported rare side effects of adenoviral vector vaccines only, web interest in this topic exceeded interest in pollution and climate change ($SS_{vaccines}$=1.2, $P$<.001 vs $SS_{pollution}$=0.81, $P$<.001; $\Delta_\%$=+44.0; Figure 5), which are much more urgent issues. Queries related to the effects of pollution and climate change returned negligible RSVs.

When considering queries related to vaccine names, the gap between the RSVs further widened; specifically, the discrepancy between the SS values of the two graphs was substantial, which underlines the unjustified disproportion in risk perception between these two topics ($SS_{vaccines}$=0.22, $P$<.001 vs $SS_{pollution}$=0.05, $P$<.001; $\Delta_\%$=+296.4). As observable in Figure 6, the risk perception in vaccines has surpassed that in pollution and climate change in a homogeneous way throughout Italy.

**Figure 3.** Comparison between web interest in conspiracy hypotheses and antihoax services, by keyword search, from 2016 to 2021.



**Figure 4.** Heat maps comparing web interest in conspiracy hypotheses and antihoax services for each time period in the Italian regions. The index shows the percentage of infodemic queries (eg, 70 means 70% conspiracy-related queries vs 30% antihoax-related queries).

**Figure 5.** Web interest in vaccine side effects compared with interest in pollution and climate change from 2016 to 2021.



**Figure 6.** Heat maps comparing web interest in vaccines and web interest in pollution and climate change for each time period in the Italian regions. The index shows the percentage of vaccine-related queries (eg, 70 means 70% queries about vaccine side effects vs 30% queries about pollution and climate change).



## COVID-19–Related Fake News

By temporarily excluding vaccines, web interest in COVID-19–related fake news reached its peak during the first wave of the pandemic and then declined, as of April 2021 ($\Delta_\% \in [-86.1, -73.7]$, $t \in [-5.7, -2.3]$). However, as observable in Figure 7, the trend of keywords related to the engineered novel coronavirus and unproven COVID-19 remedies had stabilized at values significantly far from 0 ($\boxed{\times} = 4.8 \pm 0.4$, $\boxed{\times} = 3.8 \pm 0.4$, respectively). By restricting the domain from January to June 2021—so as to obtain daily RSVs instead of weekly

RSVs—a level-shift of web interest in the manufactured origin of SARS-CoV-2 was evident (comparison between May 1 to 22 and May 23 to June 1; $\Delta_\%=+119.1$, $t_{10.3}=2.5$). Through the iterative comparison of RSVs, it was possible to estimate that, in the last 12 months, this query represented about 0.04% of COVID-19–related web searches. Regarding vaccines, the highest RSV peak was reached in the week of May 16 to 22, 2021. Such a surge was mainly due to the query "vaccino calamita" (vaccine magnet). By comparing the time-lapse periods of January 1 to May 22, 2021, and May 22 to 31, 2021, a 761% increase in infodemic searches on vaccines was found ($t_{9.5}=6.1$).

**Figure 7.** Web interest in COVID-19–related conspiracies over time. The square roots of the relative search volume values have been reported for reasons of readability.



At the regional level (Table 2), from the beginning of the pandemic, web interest in COVID-19–related fake news was not equally distributed (SD% ∈ [21.4, 33.2]), as opposed to that of generic news (SD%=9.4). Furthermore, Tables 2 and 3 testify to the absence of a regional predisposition for fake news in general and, at the same time, prove a diffused interest in specific topics. Indeed, all of the keywords were low or noncorrelated to each other ($|r| \in [0.01, 0.39]$, $P \in [.13, .97]$).

**Table 2.** Relative search volumes (RSVs) on the web of COVID-19–related fake news from January 2020 to June 2021.

| Variable | Value | | | | | |
|---|---|---|---|---|---|---|
| | Laboratory origin | Plot | 5G | Unproven remedies | Vaccines | General news |
| **RSV for each Italian region** | | | | | | |
| Abruzzo | 67 | 50 | 57 | 58 | 92 | 86 |
| Aosta | N/A[a] | N/A | 100 | N/A | N/A | 84 |
| Apulia | 100 | 64 | 77 | 53 | 92 | 80 |
| Basilicata | N/A | 74 | 28 | 78 | N/A | 92 |
| Calabria | 90 | 100 | 61 | 76 | 100 | 90 |
| Campania | 71 | 62 | 80 | 68 | 46 | 73 |
| Emilia-Romagna | 64 | 69 | 91 | 62 | 64 | 80 |
| Friuli-Venezia Giulia | 47 | 47 | 91 | 60 | 97 | 87 |
| Lazio | 84 | 57 | 58 | 68 | 81 | 84 |
| Liguria | 80 | 42 | 77 | 73 | 45 | 79 |
| Lombardy | 63 | 69 | 76 | 98 | 56 | 82 |
| Marche | 72 | 33 | 68 | 92 | 82 | 85 |
| Molise | N/A | 43 | N/A | N/A | N/A | 85 |
| Piedmont | 71 | 78 | 75 | 100 | 71 | 80 |
| Sardinia | 85 | 50 | 65 | 72 | 39 | 88 |
| Sicily | 75 | 69 | 82 | 42 | 50 | 80 |
| Trentino-Alto Adige | 44 | 26 | 40 | 55 | N/A | 66 |
| Tuscany | 65 | 53 | 61 | 80 | 89 | 98 |
| Umbria | 64 | 25 | 71 | 87 | N/A | 100 |
| Veneto | 50 | 68 | 67 | 68 | 95 | 76 |
| **Other statistics** | | | | | | |
| Mean (SD) | 70.1 (15.0) | 56.8 (18.9) | 69.7 (17.3) | 71.7 (15.8) | 73.3 (21.5) | 83.8 (7.9) |
| SD% | 21.4 | 33.2 | 24.7 | 22.0 | 29.4 | 9.4 |
| SEM (standard error of the mean) | 3.6 | 4.3 | 4.0 | 3.7 | 7.6 | 1.8 |
| Shapiro-Wilk $P$ value | .86 | .78 | .59 | .92 | .09 | .76 |

[a]N/A: not applicable due to Google Trends detection anomalies.

**Table 3.** Correlation analysis (relative search volume Pearson $r$ and two-tailed $P$ value) among topics regarding COVID-19–related fake news in the Italian regions.

| Topic | Laboratory origin | Plot | 5G | Unproven remedies | Vaccines[a] | General news |
|---|---|---|---|---|---|---|
| **Laboratory origin** | | | | | | |
| $r$ | 1 | 0.38 | 0.06 | –0.01 | –0.15 | 0.20 |
| $P$ value | __[b] | .13 | .82 | .97 | .59 | .44 |
| **Plot** | | | | | | |
| $r$ | 0.38 | 1 | 0.26 | 0.03 | 0.15 | –0.06 |
| $P$ value | .13 | — | .31 | .91 | .59 | .82 |
| **5G** | | | | | | |
| $r$ | 0.06 | 0.26 | 1 | –0.02 | –0.32 | 0.05 |
| $P$ value | .82 | .31 | — | .94 | .24 | .85 |
| **Unproven remedies** | | | | | | |
| $r$ | –0.01 | 0.03 | –0.02 | 1 | –0.05 | 0.36 |
| $P$ value | .97 | .91 | .94 | — | .86 | .16 |
| **Vaccines[a]** | | | | | | |
| $r$ | –0.15 | 0.15 | –0.32 | –0.05 | 1 | 0.39 |
| $P$ value | .59 | .59 | .24 | .86 | — | .15 |
| **General news** | | | | | | |
| $r$ | 0.20 | –0.06 | 0.05 | 0.36 | 0.39 | 1 |
| $P$ value | .44 | .82 | .85 | .16 | .15 | — |

[a]Only 15 Italian regions were included in this analysis.

[b]Not applicable.

## *Discussion*

### Principal Findings

To the best of the author's knowledge, this is the first study to investigate the impact of COVID-19 on pre-existing fake news and the risk perception of Italian web users. These findings show that the pandemic— understood as a set of different situations, such as a health crisis, an economic crisis, lockdowns, disease, and an infodemic—has significantly increased the phenomenon of conspiracies and interest in them. This influence not only caused a marked initial growth of RSV during the first lockdown (March to May 2020) but also generated a pronounced level-shift in web interest that has persisted until at least April 2021. Regional web interest in conspiracy hypotheses during the pre–COVID-19 2016-2020 period had assumed a clear negative trend and was more noticeable in some areas than in others. However, with the advent of the novel coronavirus, interest has increased to reach the highest level in the last 5 years, becoming even more homogeneous across regions. Due to the high dependence of the RSV on the day of gathering, it was not possible to search for correlations with the regional numbers of COVID-19 cases; nevertheless, the mean values and the sample variances have always remained similar. On the contrary, when analyzing the data year by year, no change in average interest in the infodemic YouTube channel ByoBlu, which had over 500,000 subscribers, was observed between the regions. Since a strong increase was highlighted nationwide,

some regions must have contributed far more than others to the jump in total RSV. Specifically, as shown in Figure 2, Campania, Lazio, Friuli-Venezia Giulia, Trentino-Alto Adige, and Valle d'Aosta experienced a much higher increase than the other regions. Moreover, a growing trend in RSV during the last 5 years involved all regions except Basilicata and Molise. Finally, web interest in fake news sources has increased more than interest in antihoax services. These results are not to be underestimated; indeed, the Ministry of Health, various online platforms such as YouTube and Twitter, and social networks such as Facebook and Instagram have declared war without borders against the rampant infodemic. Specifically, under each video relating to the pandemic, YouTube has affixed a warning bar that offers users the opportunity to read the latest COVID-19 news on the Ministry of Health official website, complete with a button to access it. A similar procedure has been adopted by Facebook and Instagram. All of these companies have banned accounts and channels that are protagonists of the spread of fake news, including ByoBlu [32-36]. This approach is partly consistent with the procedure proposed by the WHO to deal with the infodemic but was not enough to contain disinformation in Italy. Among the problems that have undermined the effectiveness of these strategies, there is the resonance given by newspapers and television channels to unreliable or misleading information [37-40]. Beyond the mere disinformation contribution, this can foster distrust toward mass media, making the information campaign even more complex during times of crisis [41]. Furthermore, despite all of the countermeasures

XSL•FO

RenderX

adopted, social networks and messaging apps, such as WhatsApp, are fake news vehicles [37,38,41].

Alongside this, the influence of newspapers and television news on the risk perception related to COVID-19 vaccines was evident. Although the trend has been on the rise since early October 2020, the headlines of online, printed, and television news publications have often been the subject of criticism from the scientific community as sensationalistic and far from the scientific evidence [42-44]. Distrust of vaccines is a growing issue that raises a serious public health question [45-47]. Although most vaccine-related fake news circulates on social networks, the national mass media must attend to the evidence presented in the scientific literature with appropriate and thoughtful language. In particular, the effect of deliberately misleading titles linked to secondary aspects of the article can have serious consequences [48]. In such an intricate scenario, web interest in vaccine side effects has overtaken interest in pollution and climate change. Notwithstanding that the author of this paper loudly supports the pharmacovigilance process and is aware of the existence of numerous studies on the possible causal link between adenoviral vector COVID-19 vaccines and thrombotic events [49-51], it is necessary to consider that pollution and climate change constitute one of the major global threats today, claiming millions of victims every year [27,52]. Since Janssen and Vaxzevria vaccines have very rare side effects [28,52-55], it is reasonable to conclude that the risk perception of Italian users is distorted and disconnected from the real dangers that menace them. This is even more true when considering the incidence of COVID-19 itself in this type of event [56,57].

Finally, COVID-19 and the crisis it caused have generated fertile ground for new conspiracy hypotheses. While some of these, such as the link between 5G and the spread of the epidemic, have waned over time, others, including the human engineering of the virus in a Wuhan laboratory, phantom infection remedies with no scientific basis, and intentionally altered vaccines, have persisted until today. To further complicate the scenario, the spread of COVID-19–related fake news has not been uniform among the regions; indeed, the RSV groups showed low multicollinearity and vast discrepancies (eg, Abruzzo, Puglia, Calabria, Friuli-Venezia Giulia, and Veneto showed a high interest in the vaccine infodemic and a significantly lower

interest in unproven remedies). As a conclusive consideration, the author of this paper emphasizes that it is correct to use the term "infodemic" to describe news that supports any hypotheses without supporting evidence [58]. At the same time, science needs to continue to investigate any possible leads [59].

## Limitations

There are no guarantees that Google Trends is sufficient for investigating the totality of the interests of the Italian public. In particular, internet penetration in Italy is equal to about 74% of the population [60]. Of this fraction, almost 96% use Google as their default online search engine [61]. Therefore, 29% of the Italian population is not considered in this survey. Furthermore, it is not certain that the keywords used in this research included all the terms related to the topics investigated. Indeed, the algorithm with which Google selects the most relevant related queries is unknown (ie, it may not consider web searches pertinent to the discussion). Finally, some relevant keywords may not have been selected for the analysis. Future research could rely on machine learning algorithms for textual analysis to derive the topics of interest to search for on Google Trends.

## Conclusions

COVID-19 has given a significant boost to web interest in conspiracy hypotheses and has made it more uniform across regions. The pandemic accelerated an already-growing trend in users' interest toward some fake news sources, including the 500,000-subscriber YouTube channel ByoBlu, which was removed from the platform by YouTube for disinformation in March 2021. The risk perception related to COVID-19 vaccines has been so distorted that vaccine side effect–related queries outweighed those relating to pollution and climate change, which are much more urgent issues. Moreover, a large amount of fake news circulated about COVID-19 vaccines, remedies, and origin. Based on these findings, it is recommended that the Italian authorities implement more effective infoveillance systems, and that communication by the mass media be less sensationalistic and more consistent with the available scientific evidence. In this context, Google Trends can be used to monitor the users' response to specific infodemiological countermeasures. Further research is needed to understand the psychological mechanisms that regulate risk perception.

## Conflicts of Interest

None declared.

Multimedia Appendix 1
Analysis of the degree of an infodemic using keywords mined from Bufale.net.
[DOCX File , 509 KB - infodemiology_v1i1e29929_app1.docx ]

Multimedia Appendix 2
Details of the analysis.
[XLSX File (Microsoft Excel File), 149 KB - infodemiology_v1i1e29929_app2.xlsx ]

## References

1.  WHO Coronavirus (COVID-19) Dashboard. World Health Organization. 2020. URL: https://covid19.who.int/ [accessed 2021-04-24]

2.  The COVID-19 infodemic. World Health Organization. 2020. URL: https://www.who.int/health-topics/infodemic/ the-covid-19-infodemic [accessed 2021-04-24]

3.  Eysenbach G. Infodemiology and infoveillance: Framework for an emerging set of public health informatics methods to analyze search, communication and publication behavior on the internet. J Med Internet Res 2009;11(1):e11 [FREE Full text] [doi: 10.2196/jmir.1157] [Medline: 19329408]

4.  Hu Z, Yang Z, Li Q, Zhang A. The COVID-19 infodemic: Infodemiology study analyzing stigmatizing search terms. J Med Internet Res 2020 Nov 16;22(11):e22639 [FREE Full text] [doi: 10.2196/22639] [Medline: 33156807]

5.  Di Simone E, Di Muzio M, Dionisi S, Giannetta N, Di Muzio F, De Gennaro L, et al. Infodemiological patterns in searching medication errors: Relationship with risk management and shift work. Eur Rev Med Pharmacol Sci 2019 Jun;23(12):5522-5529 [FREE Full text] [doi: 10.26355/eurrev_201906_18224] [Medline: 31298407]

6.  Tsao S, Chen H, Tisseverasinghe T, Yang Y, Li L, Butt ZA. What social media told us in the time of COVID-19: A scoping review. Lancet Digit Health 2021 Mar;3(3):e175-e194 [FREE Full text] [doi: 10.1016/S2589-7500(20)30315-0] [Medline: 33518503]

7.  Cross-regional statement on "infodemic" in the context of COVID-19. France Organisation des Nations Unies. 2020. URL: https://onu.delegfrance.org/IMG/pdf/cross-regional_statement_on_infodemic_final_with_all_endorsements.pdf [accessed 2021-04-24]

8.  FAQ about Google Trends data. Google Support. URL: https://support.google.com/trends/answer/4365533?hl=en [accessed 2021-04-24]

9.  Sousa-Pinto B, Anto A, Czarlewski W, Anto JM, Fonseca JA, Bousquet J. Assessment of the impact of media coverage on COVID-19–related Google Trends data: Infodemiology study. J Med Internet Res 2020 Aug 10;22(8):e19611 [FREE Full text] [doi: 10.2196/19611] [Medline: 32530816]

10. Nuti SV, Wayda B, Ranasinghe I, Wang S, Dreyer RP, Chen SI, et al. The use of Google Trends in health care research: A systematic review. PLoS One 2014;9(10):e109583 [FREE Full text] [doi: 10.1371/journal.pone.0109583] [Medline: 25337815]

11. Saladino V, Algeri D, Auriemma V. The psychological and social impact of Covid-19: New perspectives of well-being. Front Psychol 2020;11:577684 [FREE Full text] [doi: 10.3389/fpsyg.2020.577684] [Medline: 33132986]

12. Asgari Mehrabadi M, Dutt N, Rahmani AM. The causality inference of public interest in restaurants and bars on daily COVID-19 cases in the United States: Google Trends analysis. JMIR Public Health Surveill 2021 Apr 06;7(4):e22880 [FREE Full text] [doi: 10.2196/22880] [Medline: 33690143]

13. Effenberger M, Kronbichler A, Shin JI, Mayer G, Tilg H, Perco P. Association of the COVID-19 pandemic with internet search volumes: A Google Trends analysis. Int J Infect Dis 2020 Jun;95:192-197 [FREE Full text] [doi: 10.1016/j.ijid.2020.04.033] [Medline: 32305520]

14. Jimenez AJ, Estevez-Reboredo RM, Santed MA, Ramos V. COVID-19 symptom-related Google searches and local COVID-19 incidence in Spain: Correlational study. J Med Internet Res 2020 Dec 18;22(12):e23518 [FREE Full text] [doi: 10.2196/23518] [Medline: 33156803]

15. Cervellin G, Comelli I, Lippi G. Is Google Trends a reliable tool for digital epidemiology? Insights from different clinical settings. J Epidemiol Glob Health 2017 Sep;7(3):185-189 [FREE Full text] [doi: 10.1016/j.jegh.2017.06.001] [Medline: 28756828]

16. Jun SP, Yoo HS, Choi S. Ten years of research change using Google Trends: From the perspective of big data utilizations and applications. Technol Forecast Soc Change 2018 May;130:69-87 [FREE Full text] [doi: 10.1016/j.techfore.2017.11.009]

17. Sousa-Pinto B, Anto A, Czarlewski W, Anto JM, Fonseca JA, Bousquet J. Assessment of the impact of media coverage on COVID-19–related Google Trends data: Infodemiology study. J Med Internet Res 2020 Aug 10;22(8):e19611 [FREE Full text] [doi: 10.2196/19611] [Medline: 32530816]

18. Benefits of getting a COVID-19 vaccine. Centers for Disease Control and Prevention. 2020. URL: https://www.cdc.gov/ coronavirus/2019-ncov/vaccines/vaccine-benefits.html [accessed 2021-04-24]

19. COVID-19 vaccines: Key facts. European Medicines Agency. 2020. URL: https://www.ema.europa.eu/en/human-regulatory/ overview/public-health-threats/coronavirus-disease-covid-19/treatments-vaccines/vaccines-covid-19/ covid-19-vaccines-key-facts [accessed 2021-04-24]

20. Rovetta A, Bhagavathula AS. COVID-19–related web search behaviors and infodemic attitudes in Italy: Infodemiological study. JMIR Public Health Surveill 2020 May 05;6(2):e19374 [FREE Full text] [doi: 10.2196/19374] [Medline: 32338613]

21. Pluchino A, Biondo AE, Giuffrida N, Inturri G, Latora V, Le Moli R, et al. A novel methodology for epidemic risk assessment of COVID-19 outbreak. Sci Rep 2021 Mar 05;11(1):5304 [FREE Full text] [doi: 10.1038/s41598-021-82310-4] [Medline: 33674627]

22. Pegoraro V, Heiman F, Levante A, Urbinati D, Peduto I. An Italian individual-level data study investigating on the association between air pollution exposure and Covid-19 severity in primary-care setting. BMC Public Health 2021 May 12;21(1):902 [FREE Full text] [doi: 10.1186/s12889-021-10949-9] [Medline: 33980180]

23. Auci S, Vignani D. Climate variability and agriculture in Italy: A stochastic frontier analysis at the regional level. Economia Polit 2020 Jan 20;37(2):381-409. [doi: 10.1007/s40888-020-00172-x]

24. Rovetta A. Reliability of Google Trends: Analysis of the limits and potential of web infoveillance during COVID-19 pandemic and for future research. Front Res Metr Anal 2021;6:670226 [FREE Full text] [doi: 10.3389/frma.2021.670226] [Medline: 34113751]

25. The Black List. Bufale.net. URL: https://www.bufale.net/the-black-list-la-lista-nera-del-web/ [accessed 2021-04-22]

26. Alessandrini ER, Faustini A, Chiusolo M, Stafoggia M, Gandini M, Demaria M, Gruppo collaborativo EpiAir2. Air pollution and mortality in twenty-five Italian cities: Results of the EpiAir2 Project [Article in Italian]. Epidemiol Prev 2013;37(4-5):220-229 [FREE Full text] [Medline: 24293487]

27. Global climate change: Vital signs of the planet. NASA. URL: https://climate.nasa.gov/effects/ [accessed 2021-04-24]

28. COVID-19 vaccine Janssen: EMA finds possible link to very rare cases of unusual blood clots with low blood platelets. European Medicines Agency. 2021 Apr 20. URL: https://www.ema.europa.eu/en/news/ covid-19-vaccine-janssen-ema-finds-possible-link-very-rare-cases-unusual-blood-clots-low-blood [accessed 2021-04-24]

29. Moscadelli A, Albora G, Biamonte MA, Giorgetti D, Innocenzio M, Paoli S, et al. Fake news and Covid-19 in Italy: Results of a quantitative observational study. Int J Environ Res Public Health 2020 Aug 12;17(16):5850 [FREE Full text] [doi: 10.3390/ijerph17165850] [Medline: 32806772]

30. Nuovo coronavirus: Fake news. Ministero della Salute. URL: https://www.salute.gov.it/portale/nuovocoronavirus/ archivioFakeNewsNuovoCoronavirus.jsp [accessed 2021-06-04]

31. Rovetta A, Castaldo L. The impact of COVID-19 on Italian web users: A quantitative analysis of regional hygiene interest and emotional response. Cureus 2020 Sep 29;12(9):e10719 [FREE Full text] [doi: 10.7759/cureus.10719] [Medline: 33150116]

32. Coronavirus: YouTube bans 'medically unsubstantiated' content. BBC News. 2021 Apr 22. URL: https://www.bbc.com/ news/technology-52388586 [accessed 2021-04-25]

33. COVID-19 misleading information policy. Twitter. URL: https://help.twitter.com/en/rules-and-policies/ medical-misinformation-policy [accessed 2021-04-25]

34. Paul K. Facebook bans misinformation about all vaccines after years of controversy. The Guardian. 2021 Feb 08. URL: https://www.theguardian.com/technology/2021/feb/08/facebook-bans-vaccine-misinformation [accessed 2021-04-25]

35. Keeping people informed, safe, and supported on Instagram. Instagram. 2020 Mar 24. URL: https://about.instagram.com/ blog/announcements/coronavirus-keeping-people-safe-informed-and-supported-on-instagram [accessed 2021-04-25]

36. Youtube chiude Byoblu. Il fondatore Messora lancia il crowdfunding: "Compriamo un canale sul digitale". La Repubblica. 2021 Mar 31. URL: https://www.repubblica.it/politica/2021/03/31/news/byoblu_chiusura_youtube_messora-294490746/ [accessed 2021-04-25]

37. Tagliabue F, Galassi L, Mariani P. The "pandemic" of disinformation in COVID-19. SN Compr Clin Med 2020 Aug 01:1-3 [FREE Full text] [doi: 10.1007/s42399-020-00439-1] [Medline: 32838179]

38. Ali S. Hum Arenas 2020 Oct 07:1-16 [FREE Full text] [doi: 10.1007/s42087-020-00139-1]

39. Rovetta A, Castaldo L. The influence of mass media on Italian web users during COVID-19: An infodemiological analysis. SocArXiv. Preprint posted online on March 24, 2021. [FREE Full text] [doi: 10.31235/osf.io/28m6n]

40. Ferreira G, Borges S. Media and misinformation in times of COVID-19: How people informed themselves in the days following the Portuguese declaration of the state of emergency. Journal Media 2020 Dec 02;1(1):108-121 [FREE Full text] [doi: 10.3390/journalmedia1010008]

41. Fernández-Torres MJ, Almansa-Martínez A, Chamizo-Sánchez R. Infodemic and fake news in Spain during the COVID-19 pandemic. Int J Environ Res Public Health 2021 Feb 12;18(4):1781 [FREE Full text] [doi: 10.3390/ijerph18041781] [Medline: 33673095]

42. "AstraZeneca, paura in Europa": Polemiche per il titolo di Repubblica. Today. 2021 Mar 12. URL: https://www.today.it/ rassegna/vaccino-astrazeneca-oggi-repubblica.html [accessed 2021-04-25]

43. I giornali stupiti per i rifiuti al vaccino AstraZeneca dopo giorni di titoloni su reazioni avverse. Bufale.net. 2021 Mar 14. URL: https://www.bufale.net/i-giornali-stupiti-per-i-rifiuti-al-vaccino-astrazeneca-dopo-giorni-di-titoloni-su-reazioni-avverse/ [accessed 2021-04-25]

44. Bucci E. Chi urla "morto dopo il vaccino!" dovrebbe prima dimostrare il nesso di causalità. Il Foglio. 2021 Jan 19. URL: https://www.ilfoglio.it/scienza/2021/01/18/news/ chi-urla-morto-dopo-il-vaccino-dovrebbe-prima-dimostrare-il-nesso-di-causalita--1699059/ [accessed 2021-04-25]

45. Elleray E. Public vaccine distrust. Br Dent J 2021 Jan;230(2):60 [FREE Full text] [doi: 10.1038/s41415-021-2617-8] [Medline: 33483639]

46. Bogart LM, Ojikutu BO, Tyagi K, Klein DJ, Mutchler MG, Dong L, et al. COVID-19 related medical mistrust, health impacts, and potential vaccine hesitancy among Black Americans living with HIV. J Acquir Immune Defic Syndr 2021 Feb 01;86(2):200-207 [FREE Full text] [doi: 10.1097/QAI.0000000000002570] [Medline: 33196555]

47. Caudal H, Briend-Godet V, Caroff N, Moret L, Navas D, Huon JF. Vaccine distrust: Investigation of the views and attitudes of parents in regard to vaccination of their children. Ann Pharm Fr 2020 Jul;78(4):294-302. [doi: 10.1016/j.pharma.2020.03.003] [Medline: 32434681]

48.  Ecker UKH, Lewandowsky S, Chang EP, Pillai R. The effects of subtle misinformation in news headlines. J Exp Psychol Appl 2014 Dec;20(4):323-335. [doi: 10.1037/xap0000028] [Medline: 25347407]

49.  Greinacher A, Thiele T, Warkentin TE, Weisser K, Kyrle PA, Eichinger S. Thrombotic thrombocytopenia after ChAdOx1 nCov-19 vaccination. N Engl J Med 2021 Jun 03;384(22):2092-2101 [FREE Full text] [doi: 10.1056/NEJMoa2104840] [Medline: 33835769]

50.  Schultz NH, Sørvoll IH, Michelsen AE, Munthe LA, Lund-Johansen F, Ahlen MT, et al. Thrombosis and thrombocytopenia after ChAdOx1 nCoV-19 vaccination. N Engl J Med 2021 Jun 03;384(22):2124-2130 [FREE Full text] [doi: 10.1056/NEJMoa2104882] [Medline: 33835768]

51.  Marks P, Schuchat A. US Food and Drug Administration. 2021 Apr 13. URL: https://www.fda.gov/news-events/press-announcements/joint-cdc-and-fda-statement-johnson-johnson-covid-19-vaccine [accessed 2021-04-25]

52.  Manisalidis I, Stavropoulou E, Stavropoulos A, Bezirtzoglou E. Environmental and health impacts of air pollution: A review. Front Public Health 2020;8:14 [FREE Full text] [doi: 10.3389/fpubh.2020.00014] [Medline: 32154200]

53.  Østergaard SD, Schmidt M, Horváth-Puhó E, Thomsen RW, Sørensen HT. Thromboembolism and the Oxford-AstraZeneca COVID-19 vaccine: Side-effect or coincidence? Lancet 2021 Apr 17;397(10283):1441-1443 [FREE Full text] [doi: 10.1016/S0140-6736(21)00762-5] [Medline: 33798498]

54.  Vogel G, Kupferschmidt K. Side effect worry grows for AstraZeneca vaccine. Science 2021 Apr 02;372(6537):14-15. [doi: 10.1126/science.372.6537.14] [Medline: 33795437]

55.  Hunter PR. Thrombosis after COVID-19 vaccination. BMJ 2021 Apr 14;373:n958. [doi: 10.1136/bmj.n958] [Medline: 33853865]

56.  Muñoz-Rivas N, Abad-Motos A, Mestre-Gómez B, Sierra-Hidalgo F, Cortina-Camarero C, Lorente-Ramos RM, Infanta Leonor Thrombosis Research Group. Systemic thrombosis in a large cohort of COVID-19 patients despite thromboprophylaxis: A retrospective study. Thromb Res 2021 Mar;199:132-142 [FREE Full text] [doi: 10.1016/j.thromres.2020.12.024] [Medline: 33503547]

57.  Porfidia A, Valeriani E, Pola R, Porreca E, Rutjes AWS, Di Nisio M. Venous thromboembolism in patients with COVID-19: Systematic review and meta-analysis. Thromb Res 2020 Dec;196:67-74 [FREE Full text] [doi: 10.1016/j.thromres.2020.08.020] [Medline: 32853978]

58.  Shi Z. Origins of SARS-CoV-2: Focusing on science. Infect Dis Immun 2021 Apr 20;1(1):3-4 [FREE Full text] [doi: 10.1097/ID9.0000000000000008]

59.  Bloom JD, Chan YA, Baric RS, Bjorkman PJ, Cobey S, Deverman BE, et al. Investigate the origins of COVID-19. Science 2021 May 14;372(6543):694. [doi: 10.1126/science.abj0016] [Medline: 33986172]

60.  Statista. 2021 Apr 29. URL: https://www.statista.com/topics/4217/internet-usage-in-italy/ [accessed 2021-06-03]

61.  Market share held by the leading search engines in Italy as of February 2021. Statista. 2021 Apr. URL: https://www.statista.com/statistics/623043/search-engines-ranked-by-market-share-in-italy/ [accessed 2021-06-03]

## Abbreviations

**ADF:** augmented Dickey-Fuller
**HCoV-EMC/2012:** human coronavirus–Erasmus Medical Center/2012
**MK:** Mann-Kendall
**RSV:** relative search volume
**SEM:** standard error of the mean
**SS:** Sen slope
**WHO:** World Health Organization

Original Paper

# Public Interest and Behavior Change in the United States Regarding Colorectal Cancer Following the Death of Chadwick Boseman: Infodemiology Investigation of Internet Search Trends Nationally and in At-Risk Areas

Nicholas B Sajjadi[1*], BSc; Kaylea Feldman[1*], BSc; Samuel Shepard[1*], BSc; Arjun K Reddy[1*], BA; Trevor Torgerson[1*], BSc; Micah Hartwell[1,2*], PhD; Matt Vassar[1,2*], PhD

[1]Office of Medical Student Research, College of Osteopathic Medicine, Oklahoma State University Center for Health Sciences, Tulsa, OK, United States

[2]Department of Psychiatry and Behavioral Sciences, College of Osteopathic Medicine, Oklahoma State University Center for Health Sciences, Tulsa, OK, United States

*all authors contributed equally

Corresponding Author:
Nicholas B Sajjadi, BSc
Office of Medical Student Research
College of Osteopathic Medicine
Oklahoma State University Center for Health Sciences
1111 W 17th Street
Tulsa, OK, 74107
United States
Phone: 1 9185821972
Email: nicholas.sajjadi@okstate.edu

## Abstract

**Background:** Colorectal cancer (CRC) has the third highest cancer mortality rate in the United States. Enhanced screening has reduced mortality rates; however, certain populations remain at high risk, notably African Americans. Raising awareness among at-risk populations may lead to improved CRC outcomes. The influence of celebrity death and illness is an important driver of public awareness. As such, the death of actor Chadwick Boseman from CRC may have influenced CRC awareness.

**Objective:** We sought to assess the influence of Chadwick Boseman's death on public interest in CRC in the United States, evidenced by internet searches, website traffic, and donations to prominent cancer organizations.

**Methods:** We used an auto-regressive integrated moving average model to forecast Google searching trends for the topic "Colorectal cancer" in the United States. We performed bivariate and multivariable regressions on state-wise CRC incidence rate and percent Black population. We obtained data from the American Cancer Society (ACS) and the Colon Cancer Foundation (CCF) for information regarding changes in website traffic and donations.

**Results:** The expected national relative search volume (RSV) for colorectal cancer was 2.71 (95% CI 1.76-3.66), reflecting a 3590% (95% CI 2632%-5582%) increase compared to the expected values. With multivariable regression, the statewise RSV increased for each percent Black population by 1.09 (SE 0.18, $P<.001$), with 42% of the variance explained ($P<.001$). The American Cancer Society reported a 58,000% increase in CRC-related website traffic the weekend following Chadwick Boseman's death compared to the weekend before. The Colon Cancer Foundation reported a 331% increase in donations and a 144% increase in revenue in the month following Boseman's death compared to the month prior.

**Conclusions:** Our results suggest that Chadwick Boseman's death was associated with substantial increases in awareness of CRC. Increased awareness of CRC may support earlier detection and better prognoses.

**KEYWORDS**

XSL·FO
RenderX

## Introduction

Colorectal cancer (CRC) currently has the third highest mortality rate among cancers in the United States [1]. Implementing enhanced screening, namely by colonoscopy, has led to a significant decline in CRC mortality rates in the older population; [2] however, some populations remain at disproportionately high risk. African Americans have the highest incidence and mortality rates for CRC of any ethnic group; however, they are less likely to receive appropriate screening [3]. Additionally, in 2020, Rogers et al [4] identified regions in the United States that are associated with higher rates of early onset CRC, specifically among African American males, highlighting the need to focus on improving outcomes in this population. Raising awareness of CRC disparity among stakeholders and those at increased risk is a necessary component for achieving improved outcomes.

Infodemiology is the scientific study of distributions, determinants, and characteristics of information in an electronic medium, specifically the internet, or in a population, with the ultimate goal of informing public health and public policy [5]. Infodemiologic frameworks have provided a methodological means of analyzing public interest and awareness of medical conditions [6]. A few established infodemiologic metrics include aggregated data sets revealing patterns of information-seeking or information utility on websites and social media [7], the discourse and discussion found in web-based forums or blogs, and a population's activities on search engines over time [5,8]. Using infodemiologic metrics for medical research provides a real-time data stream reflecting the dynamics of information prevalence and utility that may be difficult to capture with traditional methodologies. Using infodemiologic methods may allow researchers to gauge the public interest in and awareness of CRC in at-risk populations.

As efforts are needed to raise CRC awareness, the entertainment industry may be well positioned to exert a positive influence on public awareness. The role of celebrity influence in health communications has been well studied by communications researchers, and communications scholarship is a useful source for understanding web-based searching and engagement behaviors related to publicized celebrity health information [9,10]. The influence of celebrity health on the public's awareness of medical conditions has been expanded upon by medical infodemiology research. In one infodemiologic study, media coverage of public figures disclosing a cancer diagnosis was shown to generate substantial public interest in various types of cancer. For example, large spikes in internet searches for *pancreatic cancer* were associated with actor Patrick Swayze's public announcement of his pancreatic cancer diagnosis. Similar spikes were observed when Steve Jobs took medical leave from Apple after being diagnosed with pancreatic cancer. Additional large spikes in searches for *pancreatic cancer* were observed following both of their deaths. In some cases, this coverage led to greater measurable increases in awareness than that of traditional awareness campaigns [11]. Although the study by Kaleem et al [11] examined numerous spikes in interest for many cancers, it did not show isolated peaks for CRC due to paucity of public figures announcing they had colon cancer.

The influence of celebrity illness and death from CRC may have important implications for increasing awareness of CRC.

On August 28, 2020, Chadwick Boseman, star of the Marvel movie *Black Panther*, died at age 43 of colorectal cancer [12]. Boseman's role in this film helped to normalize African heritage and culture in the United States and is considered to have significant cultural importance in the African American community [13-15]. As Boseman was a prominent figure in the United States and in the African American community, his death presents a unique opportunity to study the influence prominent public figures have on public awareness and behaviors concerning CRC in the United States. Thus, we examined internet searching data before and after Boseman's death to examine the potential influence on public interest regarding CRC, nationally and in states containing regions at high risk for CRC. Additionally, we contacted the American Cancer Society (ACS) and the Colon Cancer Foundation (CCF) to inquire about differences in website traffic and donation revenue around the time Boseman's death was disclosed to the public. Findings from this study may strengthen our understanding of the influence that public figures have on public health and may help to raise awareness of CRC among at-risk communities.

## Methods

### Data Sources

Three data sources were used for gauging public interest in CRC surrounding Boseman's death: Google Trends (GT), Wikipedia, and Twitter. We used Google Trends [16] because it is useful for identifying regional population interests and behaviors regarding medical information [8]. GT is an open database that presents population-based Google search trends over specified time periods, allowing for nearly real-time observation. GT reports regional search volumes for selected topics over time. The relative search volume (RSV) for a topic or search term is represented as a value ranging from 0 to 100, with 100 indicating peak popularity during the designated time frame. As such, GT data may serve as a proxy for the relative interest in or awareness of a given topic in a specific region over a specific time [11]. GT data is a particularly useful tool in the field of oncology, as internet searching for cancer-related topics over time has been shown to significantly correspond with the statewise incidence and mortality rates of certain cancers in the United States, including CRC [17]. On September 13, 2020, we collected national RSV data for the topic *colorectal cancer* in the United States. We used the date range of June 13 through September 11, 2020, to observe trends prior to, during, and shortly after Boseman's death, aiming to minimize confounding from any related events.

Additionally, we collected statewise RSV data for the topic *colorectal cancer* in the United States during the peak interest time. The peak interest time was designated as 1 week before Boseman's death to 1 week after (August 21 through September 4, 2020) to capture the immediate public response associated with Chadwick Boseman's death. Data were collected on September 13, 2020, under the "subregion" category for the United States. Each state's RSV was then paired with its most recently reported CRC incidence rate [18] as well as its percent

Black population [19]. Pairing these data allowed us to explore associations between a state's RSV for *colorectal cancer* during peak interest with its CRC incidence rate and its percent Black population.

Wikipedia is the most frequently used source for seeking medical information on the internet, and it is therefore a valuable source for assessing public interest in medical topics [20]. We used the Pageviews Analysis [21] tool to acquire the number of visits to the Wikipedia pages for "Colorectal cancer" and "Chadwick Boseman" over the same time period. Pageviews Analysis was used by Brigo et al [22] to provide evidence of associations with celebrity appearances and increased Wikipedia searching that suggested increased public knowledge of multiple sclerosis. Twitter is also a useful source for infodemiologic studies, and it has been used to assess the impact of awareness campaigns and analyze public engagement with medical information [23-26]. Sprout Social [27] was used to acquire the number of tweets containing the text "colon cancer" and "chadwick boseman". We selected similar terms for Wikipedia and Twitter data to ensure content uniformity across platforms. The temporal trends for Wikipedia and Twitter were worldwide trends and were observed over the same time period as the US national GT data.

### Analyses

With the national GT data, we used an autoregressive integrated moving average [28] (ARIMA) model to forecast the expected relative search volume for *colorectal cancer* had Boseman's death from CRC not been publicly disclosed, comparing expected values to observed values. Using the statewise RSV data from the peak interest time, we employed bivariate regression models on the statewise CRC incidence rate and on the percent Black population. We then used multivariable models considering both parameters. We performed statewise regression analysis on all 50 states and the District of Columbia. We then performed a subanalysis on the 19 states shown by Rogers et al [4] to contain so-called CRC hot spots—regions with disproportionately high rates of CRC among African Americans, particularly young African American males.

To explore the behavior change associated with Boseman's death, we obtained data from the ACS regarding the percent increase in colon cancer–related traffic on their website the weekend following Boseman's death compared to the weekend prior. We also obtained the number of daily donations received by the CCF and the percent increase in revenue seen 1 month following Boseman's death compared to the month prior. Additionally, both organizations provided estimates of percent increases in total revenue 1 month following Boseman's death compared to the same time period in 2019.

All statistical analyses were performed using R, version 4.0.2 (R Foundation for Statistical Computing) [29]. Statistical significance was defined as $P<.05$. The Oklahoma State University Center for Health Science Institutional Review Board determined that this project did not qualify as human subject research as defined in 45 CFR 46.102(d) and (f); therefore, it was not subject to further oversight.

## Results

The auto.arima function in R established parameters for the ARIMA model to forecast values based on the historical mean RSV. The expected national RSV for *colorectal cancer* was 2.71 (95% CI 1.76-3.66). The observed peak RSV (100) occurring on August 29, 2020, reflects a 3590% (95% CI 2632%-5582%) increase in national RSV compared to expected values. Large spikes in Wikipedia searches and Twitter keywords were observed during the peak interest time. The GT forecast comparison and temporal trends for Wikipedia and Twitter can be found in Figure 1.

The statewise bivariate regression models for all 50 states and the District of Columbia showed that for each 1 point increase in incidence rate, the RSV would decrease by 0.25 (SE=0.4, $P=.53$), a nonsignificant finding; however, for each 1% increase in the Black population, the RSV increased by 1.06 (SE=0.18, $P<.001$). With the multivariable model, the statewise RSV further decreased per point of incidence to –0.47 (SE=0.3, $P=.13$) and increased for each percent Black population to 1.09 (SE=0.18, $P<.001$), with 42% of variance explained by the model ($P<.001$). The coefficients in the subanalysis adjusted model for the 19 states containing hot spot regions were in the same direction, but to a lesser degree, and they accounted for 33% of the variability in RSV ($P=.01$). Full results of the regression models are displayed in Table 1.

The American Cancer Society reported a 58,000% increase in colon cancer–related website traffic the weekend following Boseman's death compared to the weekend before. The CCF reported receiving 595 donations from July 29 to August 28, 2020, and 2565 donations from August 28 to September 30, 2020, representing a 331% increase in the number of donations the month following Boseman's death. The ACS reported a 35.4% increase in total revenue the week following Boseman's death compared to the same time in 2019, and the CCF reported a nearly 500% increase in total revenue the week following Boseman's death compared to the same time period in 2019. Additionally, the CCF reported a 144% increase in total revenue the month following Boseman's death compared to the month prior, and they stated that they anticipated additional daily revenue influx as companies continued to run campaigns on the CCF's behalf following Boseman's death.

**Figure 1.** (A) Relative search volumes for *colorectal cancer* in the United States before and after Chadwick Boseman's death on August 28, 2020 (indicated by the asterisk). The expected forecast from the autoregressive integrated moving average model is shown by the red line. (B) The total number of visits (worldwide) to the "Colorectal cancer" and "Chadwick Boseman" Wikipedia pages before and after Chadwick Boseman's death. The vertical axis is on a logarithmic scale, with each large tick mark representing an order of magnitude. (C) The number of tweets containing the text "colon cancer" or "chadwick boseman". The vertical axis is on a logarithmic scale.

**Table 1.** Correlation of relative search volume (RSV) with colorectal cancer (CRC) incidence and percent Black population.

| | All 50 US states and DC[a] | | | | CRC hot spot states[b] (n=19) | | | |
|---|---|---|---|---|---|---|---|---|
| | Coefficient | SE | *P* value | $R^2$ | Coefficient | SE | *P* value | $R^2$ |
| **Bivariate model** | | | | | | | | |
| Incidence of CRC | –0.25 | 0.40 | .53 | <.01 | –1.33 | 0.39 | *.003* [c] | 0.38 |
| Percent Black population | 1.06 | 0.18 | *<.001* | 0.40 | 0.63 | 0.28 | *.04* | 0.18 |
| **Multivariable model** | | | | | | | | |
| Incidence of CRC | –0.47 | 0.30 | 0.132 | 0.42 | –0.56 | 0.32 | .10 | 0.33 |
| Percent Black population | 1.09 | 0.18 | *<.001* | 0.42 | 0.75 | 0.24 | *.007* | 0.33 |

[a]DC: District of Columbia.

[b]States with hot spot counties: Alabama, Arkansas, Florida, Georgia, Illinois, Indiana, Kentucky, Louisiana, Maryland, Mississippi, North Carolina, Ohio, Oklahoma, South Carolina, Tennessee, Texas, Virginia, West Virginia.

[c]Italic text indicates statistical significance.

## Discussion

### Principal Findings

Our results suggest that Chadwick Boseman's death was associated with substantial increases in the national interest in and awareness of CRC in the United States. Although a state's RSV for *colorectal cancer* during the peak interest time was not statistically associated with its CRC incidence rate, our results suggest that the RSV for *colorectal cancer* during peak interest time was significantly associated with a state's percentage of Black residents. Significant increases in searches for CRC-related topics in states with higher percentages of African American residents suggest that Boseman's death effectively increased awareness among at-risk populations. This is an important finding given that Rogers et al [4] found young (ages 20-49) non-Hispanic Black men to have the lowest early onset CRC survival rates among all ethnic groups in 232 hot spot counties—one of which is Anderson County, South Carolina, Boseman's birthplace and home county. Numerous studies have shown that raising awareness of factors for heightened risk of CRC corresponds with increased willingness to undergo CRC screening procedures and improved attitudes toward CRC [30-32]. Increased interest in and awareness of CRC among at-risk populations, namely African Americans, may lead to earlier detections and better prognoses. The RSVs for *colorectal cancer* can be visualized in Figure 2. For comparison, we have included a US map showing the African American population density by state population for comparison [33].

**Figure 2.** (Top) RSV on Google Trends for the term *colorectal cancer* during peak interest surrounding the death of Chadwick Boseman. (Bottom) Percentage of Americans who identified as Black/African American on the 2015 American Community Survey [33]. CRC: colorectal cancer; RSV: relative search volume.

The spikes in Wikipedia searches for *colorectal cancer* and in tweets containing "colorectal cancer" coinciding with Boseman's death were the all-time largest recorded volumes for their respective platforms, putting into perspective the massive magnitude of increased searching for CRC topics following Boseman's death. Boseman's death was announced on his personal Twitter account shortly after his death from CRC, and the announcement is now the single most "liked" tweet in Twitter history [34]. Our findings support the use of Wikipedia and Twitter data as reliable indicators of public interest resulting from public disclosure of celebrity illness and necessitate further exploration of these platforms as research tools.

We found substantial increases in website traffic and donations to two prominent cancer organizations following Boseman's

death, suggesting increased public financial support of CRC research and awareness campaigns. Although these findings represent website traffic and donation behavior over a short period of time, they are compelling and promising. Chadwick Boseman's death may serve as a catalyst toward increasing public awareness of CRC, leading to increased financial support for CRC research and awareness campaigns. Increasing CRC research funding is necessary, as it is disproportionately underfunded [35]. In 2015, CRC caused 50,620 deaths and received US $18 million in funding from non-profits, whereas breast cancer caused 41,070 deaths and received US $460 million [36]. Funding from the National Cancer Institute was similarly distributed [35,37]. Long-term improvements in health outcomes are best supported by policy, evidence-based medicine, and public health initiatives; however, events such

as Boseman's death have the potential to overcome cultural and societal barriers to positive health behavior in ways that the aforementioned factors cannot.

As mentioned previously, communications research [9,10] and infodemiologic GT studies have demonstrated the significant impact public figures can have on the public's awareness of various medical conditions. Importantly, we assert that *increased awareness not resulting in measurable positive behavior change* is, while still admirable and necessary, only a partial victory. Until recently, the impact of a public figure's disclosure of illness regarding CRC was unknown. In line with our own findings, one recently published study by Naik et al [38] examined internet search interests in CRC-related topics following Chadwick Boseman's death and found increases in searches on Google and Wikipedia comparable to our own. Moreover, the study found proportionally higher increases among areas with higher proportions of Black Americans, a novel and important finding that should be used to increase CRC education among this at-risk community. Further, novel components within our study include additional infodemiologic parameters (Twitter) to measure increased interest and awareness, showing significant increases in activity and the inclusion of data from the ACS and the CCF demonstrating increased donations following Boseman's death. Our study complements that by Naik et al [38] and suggests positive behavior changes regarding CRC following Boseman's death. Together, these findings contribute to the international literature base by providing evidence of positive associated behavior change that goes beyond increased internet searching, awareness, or interest following the death of the actor. Further research regarding the direct impact of Chadwick Boseman's death on CRC screenings and attitudes toward CRC screening among African Americans is warranted, as it cannot be assessed herein.

## Strengths and Limitations

Our study has several strengths. To our knowledge, our study is the first of its kind to provide evidence of positive behavior change in the form of donations associated with Chadwick Boseman's death. Our methodology was adapted from published literature that used GT and ARIMA models to provide evidence for increased awareness of suicide prevention resources following a nationally publicized event [39]. The data provided by the ACS and CCF reflect increases in public interest and support for CRC research, possibly related to Boseman's death, and these data are unique to our study. Limitations of this study include the limited time frame for observing RSV, website traffic and donations, use of few search terms, and lack of access to donation data entailing actual dollar amounts. The limited observation period limits the generalizability of our findings to the time periods in which we searched, and the long-term influence of Boseman's death cannot be ascertained from our findings. It is also not possible to make causal claims based on our data, and our results should be interpreted accordingly. Other events could have influenced our outcomes; thus, the potential for confounding factors is present.

## Conclusion

National public interest in colorectal cancer substantially increased relative to predicted values following Chadwick Boseman's death. A state's RSV for *colorectal cancer* immediately surrounding Boseman's death was significantly associated with its percentage of Black residents, possibly suggesting increased CRC awareness among this population. Increased interest among at-risk populations associated with Chadwick Boseman's death may lead to improved health outcomes and attitudes regarding CRC. Website traffic and revenue to prominent cancer organizations increased following Chadwick Boseman's death. Increased public awareness of CRC associated with Chadwick Boseman's death may lead to increased support for CRC research and awareness campaigns.

## References

1.  Cancer facts and figures. Cancer.org. URL: https://www.cancer.org/research/cancer-facts-statistics/all-cancer-facts-figures/cancer-facts-figures-2020.html [accessed 2020-09-30]
2.  Siegel RL, Fedewa SA, Anderson WF, Miller KD, Ma J, Rosenberg PS, et al. Colorectal cancer incidence patterns in the United States, 1974-2013. J Natl Cancer Inst 2017 Aug 01;109(8) [FREE Full text] [doi: 10.1093/jnci/djw322] [Medline: 28376186]

3.    Warren Andersen S, Blot WJ, Lipworth L, Steinwandel M, Murff HJ, Zheng W. Association of race and socioeconomic status with colorectal cancer screening, colorectal cancer risk, and mortality in southern US adults. JAMA Netw Open 2019 Dec 02;2(12):e1917995 [FREE Full text] [doi: 10.1001/jamanetworkopen.2019.17995] [Medline: 31860105]

4.    Rogers CR, Moore JX, Qeadan F, Gu LY, Huntington MS, Holowatyj AN. Examining factors underlying geographic disparities in early-onset colorectal cancer survival among men in the United States. Am J Cancer Res 2020;10(5):1592-1607 [FREE Full text] [Medline: 32509399]

5.    Eysenbach G. Infodemiology and infoveillance: framework for an emerging set of public health informatics methods to analyze search, communication and publication behavior on the Internet. J Med Internet Res 2009 Mar 27;11(1):e11 [FREE Full text] [doi: 10.2196/jmir.1157] [Medline: 19329408]

6.    Eysenbach G. Infodemiology and infoveillance tracking online health information and cyberbehavior for public health. Am J Prev Med 2011 May;40(5 Suppl 2):S154-S158. [doi: 10.1016/j.amepre.2011.02.006] [Medline: 21521589]

7.    Bernardo TM, Rajic A, Young I, Robiadek K, Pham MT, Funk JA. Scoping review on search queries and social media for disease surveillance: a chronology of innovation. J Med Internet Res 2013 Jul 18;15(7):e147 [FREE Full text] [doi: 10.2196/jmir.2740] [Medline: 23896182]

8.    Mavragani A, Ochoa G. Google Trends in infodemiology and infoveillance: methodology framework. JMIR Public Health Surveill 2019 May 29;5(2):e13439 [FREE Full text] [doi: 10.2196/13439] [Medline: 31144671]

9.    Kosenko KA, Binder AR, Hurley R. Celebrity influence and identification: a test of the Angelina Effect. J Health Commun 2016;21(3):318-326. [doi: 10.1080/10810730.2015.1064498] [Medline: 26192626]

10.   Kresovich A, Noar SM. The power of celebrity health events: meta-analysis of the relationship between audience involvement and behavioral intentions. J Health Commun 2020 Jun 02;25(6):501-513. [doi: 10.1080/10810730.2020.1818148] [Medline: 32990198]

11.   Kaleem T, Malouff TD, Stross WC, Waddle MR, Miller DH, Seymour AL, et al. Google search trends in oncology and the impact of celebrity cancer awareness. Cureus 2019 Aug 10;11(8):e5360 [FREE Full text] [doi: 10.7759/cureus.5360] [Medline: 31608195]

12.   Ugwu R, Levenson M. Black Panther star Chadwick Boseman dies of cancer at 43. The New York Times. 2020 Aug 29. URL: https://www.nytimes.com/2020/08/28/movies/chadwick-boseman-dead.html [accessed 2021-02-24]

13.   King S. Black Panther is one of the most important cultural moments in American history. Medium. 2018. URL: https://medium.com/@ShaunKing/black-panther-is-one-of-the-most-important-moments-in-american-history-1fc9166a0972 [accessed 2020-10-01]

14.   Smith J. The revolutionary power of Black Panther. Time. URL: https://time.com/black-panther/ [accessed 2021-02-24]

15.   Manish PA. Chadwick Boseman 'was a real-life black superhero'. BBC. URL: https://www.bbc.com/news/newsbeat-53989167 [accessed 2021-02-24]

16.   Google Trends. URL: https://trends.google.com/trends [accessed 2020-09-13]

17.   Wehner MR, Nead KT. Can Google help us fight cancer? Lancet Oncol 2018 Jul;19(7):867. [doi: 10.1016/S1470-2045(18)30296-1] [Medline: 30084368]

18.   United States Cancer Statistics: Data Visualizations. US Centers for Disease Control and Prevention. URL: https://gis.cdc.gov/Cancer/USCS/DataViz.html [accessed 2020-09-13]

19.   State population by race, ethnicity data. Governing. URL: https://www.governing.com/gov-data/census/state-minority-population-data-estimates.html [accessed 2020-09-12]

20.   Smith DA. Situating Wikipedia as a health information resource in various contexts: a scoping review. PLoS One 2020;15(2):e0228786 [FREE Full text] [doi: 10.1371/journal.pone.0228786] [Medline: 32069322]

21.   Pageviews Analysis. URL: http://pageviews.toolforge.org [accessed 2020-09-13]

22.   Brigo F, Lattanzi S, Bragazzi N, Nardone R, Moccia M, Lavorgna L. Why do people search Wikipedia for information on multiple sclerosis? Mult Scler Relat Disord 2018 Feb;20:210-214. [doi: 10.1016/j.msard.2018.02.001] [Medline: 29428464]

23.   Lyles CR, López A, Pasick R, Sarkar U. "5 mins of uncomfyness is better than dealing with cancer 4 a lifetime": an exploratory qualitative analysis of cervical and breast cancer screening dialogue on Twitter. J Cancer Educ 2013 Mar;28(1):127-133. [doi: 10.1007/s13187-012-0432-2] [Medline: 23132231]

24.   Viguria I, Alvarez-Mon MA, Llavero-Valero M, Asunsolo Del Barco A, Ortuño F, Alvarez-Mon M. Eating disorder awareness campaigns: thematic and quantitative analysis using Twitter. J Med Internet Res 2020 Jul 14;22(7):e17626 [FREE Full text] [doi: 10.2196/17626] [Medline: 32673225]

25.   Robinson P, Turk D, Jilka S, Cella M. Measuring attitudes towards mental health using social media: investigating stigma and trivialisation. Soc Psychiatry Psychiatr Epidemiol 2019 Jan;54(1):51-58. [doi: 10.1007/s00127-018-1571-5] [Medline: 30069754]

26.   Berry N, Lobban F, Belousov M, Emsley R, Nenadic G, Bucci S. #WhyWeTweetMH: understanding why people use Twitter to discuss mental health problems. J Med Internet Res 2017 Apr 05;19(4):e107 [FREE Full text] [doi: 10.2196/jmir.6173] [Medline: 28381392]

27.   Sprout Social: social media management solutions. URL: http://sproutsocial.com [accessed 2020-09-13]

28.   Hyndman RJ, Khandakar Y. Automatic time series forecasting: the forecast package for R. J. Stat. Soft 2008;27(3):N/A. [doi: 10.18637/jss.v027.i03]

29.   Ripley B. The R project in statistical computing. MSOR Connections: The Newsletter of the LTSN Maths. 2001. URL: http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.430.3979&rep=rep1&type=pdf [accessed 2020-09-13]

30.   Knudsen MD, Hoff G, Tidemann-Andersen I, Bodin GE, Øvervold S, Berstad P. Public awareness and perceptions of colorectal cancer prevention: a cross-sectional survey. J Cancer Educ 2020 Feb 28:N/A. [doi: 10.1007/s13187-020-01721-5] [Medline: 32112366]

31.   McCaffery K, Wardle J, Waller J. Knowledge, attitudes, and behavioral intentions in relation to the early detection of colorectal cancer in the United Kingdom. Prev Med 2003 May;36(5):525-535. [doi: 10.1016/s0091-7435(03)00016-1] [Medline: 12689797]

32.   Gimeno-García AZ, Quintero E, Nicolás-Pérez D, Jiménez-Sosa A. Public awareness of colorectal cancer and screening in a Spanish population. Public Health 2011 Sep;125(9):609-615. [doi: 10.1016/j.puhe.2011.03.014] [Medline: 21794885]

33.   American Community Survey (ACS). US Census Bureau. URL: https://www.census.gov/programs-surveys/acs [accessed 2021-08-10]

34.   Most liked Tweet ever. A tribute fit for a King. #WakandaForever. @Twitter. 2020 Aug 29. URL: https://twitter.com/Twitter/status/1299808792322940928 [accessed 2020-09-30]

35.   Carter AJ, Nguyen CN. A comparison of cancer burden and research spending reveals discrepancies in the distribution of research funding. BMC Public Health 2012 Jul 17;12(1):N/A. [doi: 10.1186/1471-2458-12-526]

36.   Kamath SD, Kircher SM, Benson AB. Comparison of cancer burden and nonprofit organization funding reveals disparities in funding across cancer types. J Natl Compr Canc Netw 2019 Jul 01;17(7):849-854. [doi: 10.6004/jnccn.2018.7280] [Medline: 31319386]

37.   Funding for research areas. National Cancer Institute. URL: https://www.cancer.gov/about-nci/budget/fact-book/data/research-funding [accessed 2020-09-16]

38.   Naik H, Johnson MDD, Johnson MR. Internet interest in colon cancer following the death of Chadwick Boseman: infoveillance study. J Med Internet Res 2021 Jun 15;23(6):e27052 [FREE Full text] [doi: 10.2196/27052] [Medline: 34128824]

39.   Torgerson T, Khojasteh J, Vassar M. Public awareness for a sexual assault hotline following a Grey's Anatomy episode. JAMA Intern Med 2020 Mar 01;180(3):456-458 [FREE Full text] [doi: 10.1001/jamainternmed.2019.5280] [Medline: 31790540]

## Abbreviations

**ACS:** American Cancer Society
**ARIMA:** autoregressive integrated moving average
**CCF:** Colon Cancer Foundation
**CRC:** colorectal cancer
**GT:** Google Trends
**RSV:** relative search volume

XSL•FO
**RenderX**

Original Paper

# Online Search Behavior Related to COVID-19 Vaccines: Infodemiology Study

Lawrence An[1,2], MD; Daniel M Russell[3], PhD; Rada Mihalcea[4], PhD; Elizabeth Bacon[1], MPH; Scott Huffman[3], PhD; Ken Resnicow[1,5], PhD

[1]Center for Health Communications Research, Rogel Cancer Center, University of Michigan, Ann Arbor, MI, United States

[2]Division of General Medicine, School of Medicine, University of Michigan, Ann Arbor, MI, United States

[3]Google, Mountain View, CA, United States

[4]Computer Science and Engineering Division, College of Engineering, University of Michigan, Ann Arbor, MI, United States

[5]Department of Health Behavior & Health Education, University of Michigan School of Public Health, Ann Arbor, MI, United States

**Corresponding Author:**
Lawrence An, MD
Center for Health Communications Research
Rogel Cancer Center
University of Michigan
North Campus Research Complex, Building 16
2800 Plymouth Rd.
Ann Arbor, MI, 48109
United States
Phone: 1 734 763 6099
Email: lcan@med.umich.edu

## Abstract

**Background:** Vaccination against COVID-19 is an important public health strategy to address the ongoing pandemic. Examination of online search behavior related to COVID-19 vaccines can provide insights into the public's awareness, concerns, and interest regarding COVID-19 vaccination.

**Objective:** The aim of this study is to describe online search behavior related to COVID-19 vaccines during the start of public vaccination efforts in the United States.

**Methods:** We examined Google Trends data from January 1, 2021, through March 16, 2021, to determine the relative search volume for vaccine-related searches on the internet. We also examined search query log data for COVID-19 vaccine-related searches and identified 5 categories of searches: (1) general or other information, (2) vaccine availability, (3) vaccine manufacturer, (4) vaccine side-effects and safety, and (5) vaccine myths and conspiracy beliefs. In this paper, we report on the proportion and trends for these different categories of vaccine-related searches.

**Results:** In the first quarter of 2021, the proportion of all web-based search queries related to COVID-19 vaccines increased from approximately 10% to nearly 50% of all COVID-19–related queries ($P<.001$). A majority of COVID-19 vaccine queries addressed vaccine availability, and there was a particularly notable increase in the proportion of queries that included the name of a specific pharmacy (from 6% to 27%; $P=.01$). Queries related to vaccine safety and side-effects (<5% of total queries) or specific vaccine-related myths (<1% of total queries) were uncommon, and the relative frequency of both types of searches decreased during the study period.

**Conclusions:** This study demonstrates an increase in online search behavior related to COVID-19 vaccination in early 2021 along with an increase in the proportion of searches related to vaccine availability at pharmacies. These findings are consistent with an increase in public interest and intention to get vaccinated during the initial phase of public COVID-19 vaccination efforts.

**KEYWORDS**

## Introduction

We are currently in the midst of a global pandemic caused by COVID-19. At all times, and particularly during a pandemic, it is critical for the public to have access to timely and accurate health information [1-3]. The internet is a major source of such health information [4-7]. Analysis of health-related web search behavior can provide critical insight into the public's awareness and interest in specific health issues and their health concerns and information needs, health experiences, and health-related intentions and behaviors [8-11].

Infodemiology is the scientific study of the "distribution and determinants of information in an electronic medium, specifically the internet, or in a population, with the ultimate aim to inform public health and health policy" [8]. To date, a number of infodemiology studies related to the COVID-19 pandemic have been published. Many of these studies identified an association between general COVID-19–related or symptom-specific search trends and COVID-19 case incidence and associated deaths [12-18]. In several cases, search behavior appeared to be an effective predictor of disease trends. Several studies have examined the occurrence and spread of COVID-19–related misinformation, which has been a public health challenge during the pandemic [19-25]. Another group of studies has examined search behavior to gain insights into public awareness, interest, attitudes, and behaviors related to COVID-19 [21,26-28]. Husain et al [27] found that countries that demonstrated a more rapid increase in public search interest regarding COVID-19 also tended to be more effective in their control of the pandemic.

Vaccination against COVID-19 is a major public health strategy in the effort to end the pandemic [29]. As of Spring 2021, a number of effective COVID-19 vaccines have been available that substantially reduce the risk of COVID-19–related illness, hospitalizations, and death [30]. Understanding public awareness and interest in COVID-19 vaccines and willingness to vaccinate are critical to help guide vaccination efforts. Unfortunately, vaccine hesitancy is common and poses a major barrier to successful vaccination efforts [31]. Infodemiologic approaches can provide potentially important insights into the public's awareness, interests, concerns, and intentions related to COVID-19 vaccination. Thus far, few infodemiological studies have focused on COVID-19 vaccines [32,33]. To address this critical gap, we examined the relative search volume using Google Trends data and also search query logs capturing users' online search behaviors related to COVID-19 vaccines in the first quarter of 2021.

## Methods

### Study Design

This study describes users' online search behavior via Google search engine for searches related to COVID-19 vaccines in the first quarter of 2021. Google is the dominant search engine in the United States, accounting for approximately 89% of the total search volume in the country as of January 2021 [34]. We focused on the time period from January 1, 2021, through March 16, 2021, which follows the initial emergency-use authorization by the US Food and Drug Administration (FDA) for the Pfizer and Moderna COVID-19 vaccines and the start of public vaccination efforts.
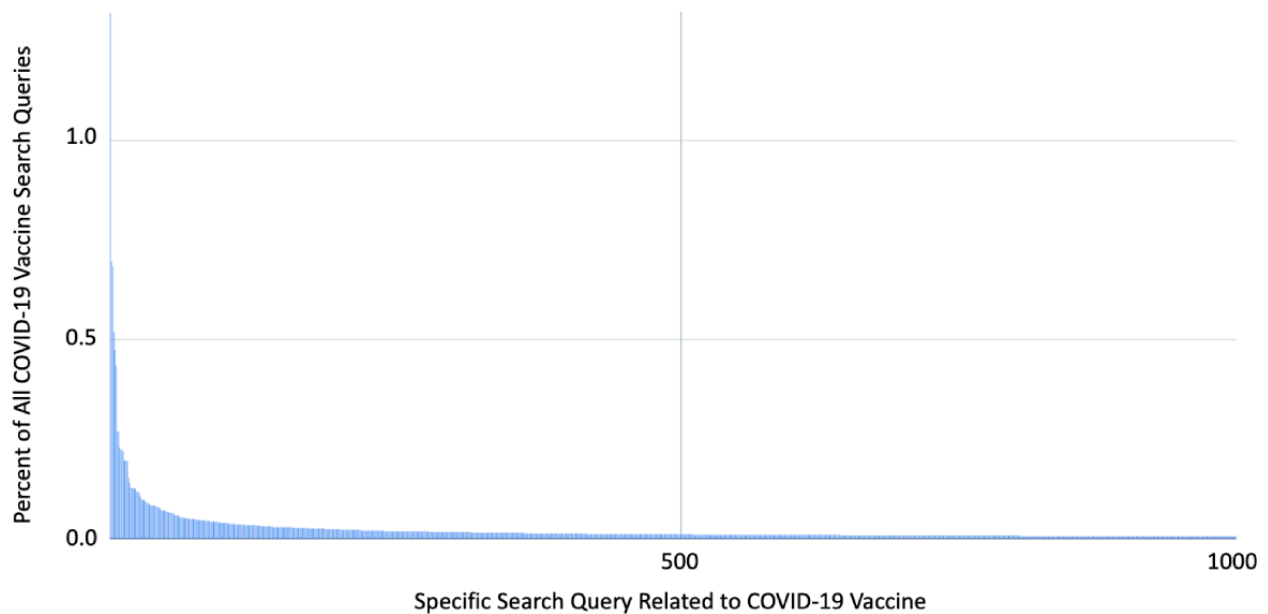
### Data Sources

#### Google Trends

Google Trends provides open access to time-series data related to Google search engine search volumes for specific terms [35]. Search query volume was normalized to a percentage scale (0% to 100%) to provide a measure of relative search volume (RSV), with 100% corresponding to the peak in search volume in any given time frame for that specific topic. By searching for multiple terms simultaneously, we were able to compare the RSV for different terms. For this report, we focus on Google Trends data for COVID-19 vaccine-related searches in the United States.

#### Search Query Logs

Search engine query logs record the specific language that users employ when conducting online searches and can provide insight into the users' information needs and interests and how these change over time [36-38]. To examine COVID-19 vaccine-related search behavior in greater detail, we compiled anonymized data from Google search query logs. That is, we examined a complete sample of English-language queries conducted in the United States during the search period (January 1, 2021, to March 16, 2021). We collected only queries that contain both the terms "COVID" and "vaccine." This data set comprises over 45.4 million queries during the sampling period, which suggests that when people search for information about COVID-19 vaccines, they use a fairly limited number of common queries. For example, the top 150 most common queries related to COVID-19 account for approximately half of all queries in the data set. The distribution of search queries by volume is shown in Figure 1. This figure shows that a small number of queries accounts for a large proportion of the overall query volume. This query volume distribution curve essentially becomes asymptotic after the top 1000 most common queries. To create this search query log data set, we collected the top 5000 most common COVID-19 vaccine-related search queries for each day during the study period.

**Figure 1.** Frequency distribution of the 1000 most common search queries related to "COVID-19 vaccine.".



## Metrics

### Overview

Metrics of interest for this study are based upon both Google Trends and the search query log data set. Specific Google Trend metrics include (1) comparison of the RSV for any searches related to COVID-19 and those related to COVID-19 vaccine, (2) general RSV for the term "vaccine" since 2005, and (3) comparison of RSV for specific vaccine myths and conspiracy beliefs identified in the search query logs.

Our main interest with the search query log data set is to examine the distribution and trends of different types of COVID-19 vaccine-related searches. We developed a COVID-19 vaccine search query classifier using the following steps.

### Step 1: Identify Search Categories

Two study authors (DMR and LA) performed independent manual review of a random sample of 1000 queries to identify common themes, as well as unique terms associated with each category. The two authors met to review and resolve any discrepancies and reached complete agreement on both categories and associated terms. Based on this review, we identified the following categories of search queries: (1) vaccine availability, (2) vaccine maker or manufacturer, (3) vaccine side effects or safety, (4) vaccine myth or conspiracy beliefs, and (5) general or other vaccine-related searches. A definition of each of these categories of searches, associated search terms, and examples is shown in Table 1. Because pharmacies were a major channel for distribution of COVID-19 vaccines, we also created a subcategory of vaccine availability queries that asked about COVID-19 vaccines in relation to pharmacies (eg, included the name of specific pharmacy chains).

**Table 1.** Types of COVID-19 vaccine–related search queries.

| Category | Definition | Associated terms | Examples of specific queries |
|---|---|---|---|
| Availability | Query that included a term or phrase identifying locations where or time when COVID-19 vaccines might be available | Names of US states, counties, or cities, names of organizations or specific locations that provide COVID-19 vaccines (eg, pharmacies, hospitals or health systems, vaccination sites), when or where to get COVID-19 vaccines | "ny covid vaccine", "covid vaccine california", "florida covid vaccine", "covid vaccine near me", "where to get covid vaccine", "cvs covid vaccine", "covid vaccine rite-aid", "covid vaccine appointment", "when can I get covid vaccine" |
| Maker or manufacturer | Query that included the name of a COVID-19 vaccine maker or manufacturer | Names of different COVID-19 vaccines, names of companies or organizations that developed or manufactured different vaccines | "pfizer vaccine", "moderna vaccine", "johnson vaccine", "j&j vaccine" |
| Side effects or safety | Query that included general or specific terms associated with side effects or safety of COVID-19 vaccines | Side effects, safety, specific vaccine-related worries and concerns | "covid vaccine side effects", "covid vaccine safety", "reaction to covid vaccine", "pregnant women covid vaccine", "covid vaccine blood clot", "problems with covid vaccine", "covid vaccine fever", "covid vaccine allergy" |
| Myths or conspiracies | Query that included general or specific terms associated with COVID-19 vaccine myths or conspiracy beliefs | Specific myths or conspiracy beliefs | "covid vaccine infertility", "does covid vaccine change dna", "covid vaccine microchip", "can I get covid from vaccine", "covid vaccine 5G" |
| General or other | Query related to COVID-19 vaccine that included no additional terms or terms not associated with any of the above categories | COVID-19 vaccine or vaccination or other topics other than identified above | "covid vaccine", "covid-19 vaccine", "coronavirus vaccine", "covid vaccine update", "covid vaccination rates" |

### Step 2: Identify Terms Associated With Each Search Category

One study author (LA) then manually reviewed an additional random sample of 5000 queries to identify any additional terms that might be associated with each search category. The results of this review were discussed with additional authors (DMR, KR, and RM) to create a final list of unique search query terms associated with each search category.

### Step 3: Create and Apply Search Query Classifier

We created a rules-based classifier that assigned a query to one or more of the 5 categories based on the presence of a unique set of associated terms while accounting for common variations in spelling. Some search queries contained terms associated with multiple categories, and these queries were counted separately in each appropriate category. For example, a search for "Pfizer covid vaccine CVS" would be counted as a query related to vaccine availability (given the presence of the name of a specific pharmacy) and also as a query related to vaccine manufacturer H (given the presence of the name of a specific vaccine manufacturer).

### Step 4: Evaluate Performance of the Search Query Classifier

This classifier was able to classify 90% of all queries in the entire COVID-19 vaccine search query log dataset. The remaining 10% of unclassified queries represented searches for additional vaccine-related information (eg, "covid vaccination rates," or "how long does the covid vaccine last") and are labeled as "other" and included as part of the category of "general or other" searches. After application of the classifier, an additional random sample of 1000 search queries with classifier results was reviewed separately by 2 authors (DMR

and LA) to assess the accuracy of the classifier. These authors met to resolve any discrepancies in manual review and the results of this review were used to calculate the classifier's precision and recall for each of the search categories. The classifier performed well with precision of 99.8% to 100% and recall over 99.5% across all search query categories.

### Analysis

Based on the search query log data set, we employed linear regression to examine the time trends for the proportion of different types of COVID-19 vaccine-related searches over time. For each week during the study period, we calculated a proportion corresponding to each of our categories of interest. The proportions examined in these analyses correspond to the above-described categories. The proportion of COVID-19 vaccine-related searches in each of these categories serves as the dependent variable in separate linear regression models. In each of these models, time (ie, week number since start of the study period) serves as the independent variable to examine the significance of trends over time.

### Ethics Review

This was reviewed by the University of Michigan Institutional Review Board and judged to be exempt based upon its use of open access and anonymized aggregate data.

## Results

### Overview

The relative search volumes for any searches related to COVID-19 and those related to COVID-19 vaccines are shown in Figure 2. This figure shows a clear increase in the relative volume of COVID-19 vaccine–related searches over the study period.

XSL•FO

RenderX

In the beginning of January 2021, approximately 10% of all COVID-19–related queries were about vaccines. By March 2021, nearly 50% of all COVID-19–related searches were vaccine related. A linear regression was calculated to predict the fraction of queries about COVID-19 vaccines based on daily change during the sample period. A significant linear regression model was found (df=103; $R^2$=0.76; beta coefficient for time=.31; $P$<.001), indicating that the RSV for COVID-19 vaccine queries increased over the study period.

Figure 3 provides a broader historical context for the level of vaccine-related search interest. This figure shows the relative search volume for the term "vaccine" from January 2005 through the first quarter of 2021. The small peak in vaccine-related search volume occurring in October 2009 coincides with the H1N1 influenza epidemic [39]. The peak in vaccine-related searches in early 2021 is several fold higher than this prior peak in 2009.

A breakdown of the proportion and trends for different types of COVID-19 vaccine-related queries based upon search query log data is shown in Figure 4. Trends for the proportions of different categories of COVID-19 vaccine-related searches are described below.

**Figure 2.** Relative search volume (RSV) for the terms "COVID" (blue) and "COVID vaccine" (red) from September 2020 through March 2021.
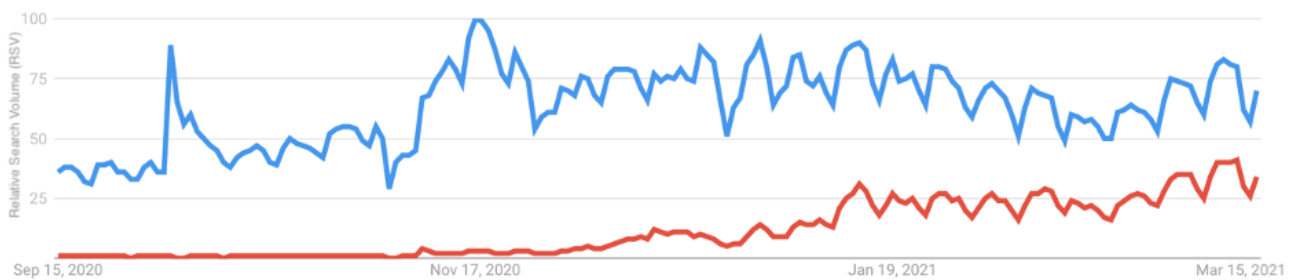


**Figure 3.** Relative search volume for the term "vaccine" from January 2005 through March 2021.
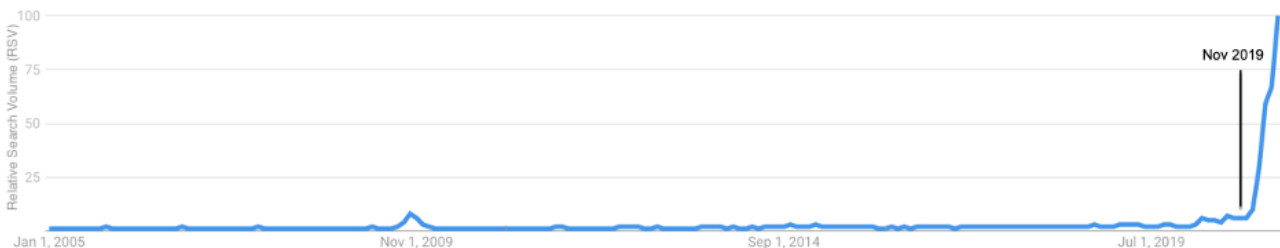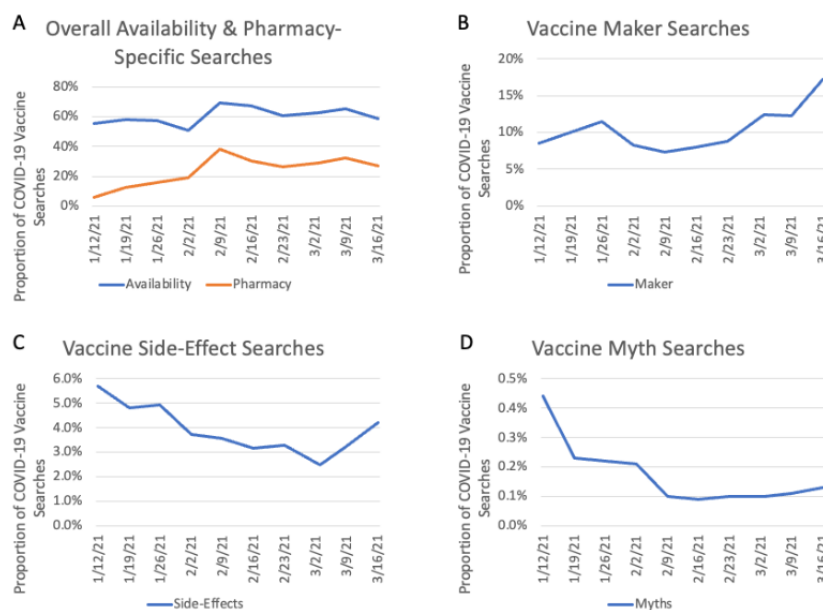


**Figure 4.** Trends for different categories of COVID-19 vaccine search queries: (A) overall availability and pharmacy, (B) vaccine manufacturer, (C) side effects and safety, and (D) myths and conspiracy beliefs.

## Vaccine Availability

During the study period, a majority of searches were classified as related to vaccine availability, with the specific proportion ranging from 55% to 69% each week (Figure 4A). The high proportion of vaccine availability searches was consistent over the study period. Linear regression showed that the time trend for the proportion of queries related to vaccine availability during the study period was not significant (df=8; $R^2$=0.19; beta coefficient for time=.81; $P$=.20).

During the study period, there was a substantial increase in the subcategory of searches related to pharmacies. The proportion of COVID-19 vaccine-related queries that included a specific pharmacy name increased from 5.9% at the start of the study period to 27.2% at the end of the study period (Figure 4A). Linear regression showed the time trend for this change was positive and significant (df=8; $R^2$=0.56; beta coefficient for time=2.51; $P$=.01).

## Vaccine Manufacturers

Over the same period, the proportion of vaccine manufacturer-related searches (eg, Pfizer, Moderna, Janssen or Johnson & Johnson) averaged 10.4%. Linear regression shows that the time trend for the proportion of vaccine manufacturer-related searches during the study period was not significant (df=8; $R^2$=0.38; beta coefficient for time=.61; $P$=.06).
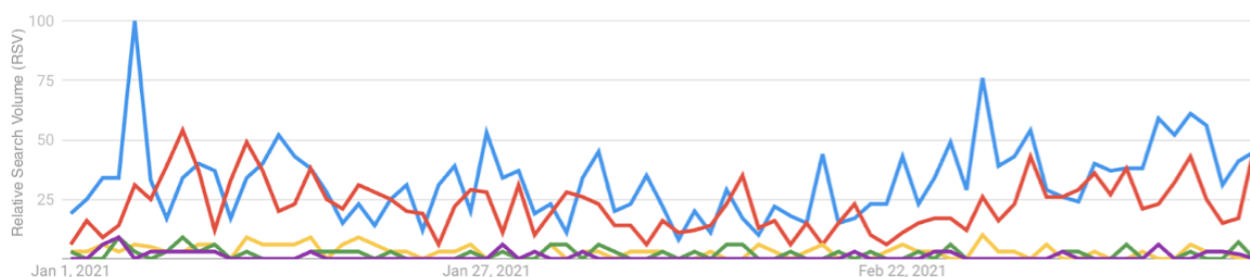
## Side Effects and Safety

The proportion of searches related to side effects or safety of COVID-19 vaccines was quite small, and the proportion actually decreased slightly (from 5.7% to 4.2%) over the study period. Linear regression showed a negative time trend for the proportion of queries related to vaccine side effects or safety during the study period (df=8; $R^2$=0.51; beta coefficient for time=–0.23; $P$=.02).

## Myths

During the study period, the overall proportion of COVID-19 vaccine-related queries that included mention of myths or conspiracies related to COVID-19 vaccines was quite low. This proportion actually decreased slightly (from 0.4% to 0.1%) over the study period. Linear regression showed a negative time trend for the proportion of queries related to vaccine myths or conspiracy beliefs during the study period (df=8; $R^2$=0.58; beta coefficient for time=–0.01; $P$=.01). Searches related to specific myths or conspiracy beliefs included searches related to the COVID-19 vaccine and (1) infertility, (2) potential to cause change in DNA, (3) 5G and the vaccine, (4) microchips, and (5) contracting COVID-19 from the vaccine itself. The Google Trends RSV for searches related to these specific myths and conspiracy topics is presented in Figure 5. This figure reveals that searches related to COVID-19 vaccine and "infertility" and "DNA" were the most common vaccine myth-related searches.

**Figure 5.** Google Trends relative search volume for queries related to specific COVID-19 vaccine myths and conspiracy beliefs (red: vaccine + DNA, blue: vaccine + fertility, green: vaccine + microchip, gold: vaccine + 5G, and purple: "COVID from vaccine").



# *Discussion*

## Principal Findings

This study reports on online search behavior related to COVID-19 vaccines in the United States, from the first quarter of 2021. During this period, there was a clear increase in the volume of online searches for information about COVID-19 vaccines with a consistently high proportion of searches related to vaccine availability. Online search behavior is influenced strongly by external events and associated media coverage [40,41]. Critical events that are likely drivers of the observed patterns in online search behavior include emergency-use authorizations by the US FDA for 3 different vaccines (ie, Pfizer on December 11, 2020, Moderna on December 18, 2020, Johnson & Johnson on February 26, 2021) and the beginning of public vaccination efforts. During this time, web-based registration was also one of the major means to obtain a vaccine

appointment. A particularly notable rise in the proportion of COVID-19 vaccine searches that include the names of specific pharmacies is consistent with the US national strategy that featured pharmacies as vaccination sites.

We interpret these patterns of online search behavior related to COVID-19 vaccines, particularly the rise in pharmacy-related searches, as a sign of increased readiness and intentions to vaccinate among the US population during the study period. This interpretation is consistent with the findings of national tracking surveys in the United States that demonstrate a similar increase in intentions to vaccinate over the study period. Specifically, the Kaiser Family Foundation COVID-19 Tracking Survey shows that the proportion of US adults that either had been vaccinated or would want to get vaccinated as soon as possible increased from under 40% to over 60% during the study period [42]. The increase in online searches related to COVID-19 vaccines is also consistent with national survey

XSL•FO
**RenderX**

findings that reported an increase in the proportion of US adults who reported they had "enough information about when or where to get the vaccine" also over the same period [42].

It is interesting to consider how to interpret our findings regarding online search behavior related to information about COVID-19 vaccine side effects and safety. National surveys conducted during the study period show that among the majority of adults who had not yet been vaccinated, a substantial proportion were concerned about the long-term effects of COVID-19 vaccines (68%), the potential for serious side effects from COVID-19 (59%), or that the vaccines may not be safe (55%) [43]. Given the prevalence of these concerns about COVID-19 vaccine side effects and safety, the relatively low proportion and decreasing time trend for vaccine-related searches that addressed these topics is somewhat surprising. The low proportion of COVID-19 vaccine-related searches pertaining to side effects or safety could be due to a relatively low rate of active information seeking about these aspects among those who are hesitant to get the vaccine, active searching on the web for vaccine appointments among those individuals highly motivated to obtain the vaccine, or some combination of these factors. In such a situation (where some segments of the population are highly motivated to search actively for information while others are not), the relative frequency of searches related to different topics does not appear to provide a good representation of the level of public concern or interest in these topics.

Our findings regarding the frequency of online searches related to COVID-19 vaccine myths or conspiracy theories can be similarly interpreted. Belief, or at least uncertainty, regarding COVID-19 vaccine myths is unfortunately common. During the same time period as this study, national surveys show that 34% of adults in the United States who had not been vaccinated either believed or were unsure about one or more common COVID-19 vaccine-related myths [43]. Nevertheless, at the same time, we found less than 1% of online searches related to COVID-19 vaccines addressed these topics and that this proportion actually decreased over the study period. These findings suggest that many individuals who either believe or are unsure about COVID-19 vaccine myths are not actively seeking additional information online to help clarify their understanding. Although making inferences regarding the population prevalence of a particular COVID-19 vaccine myth may be difficult based upon analysis of search behavior, it is possible that relative search volume might be useful in helping to assess which particular myths are more or less common. For example, the results of our RSV comparison for specific vaccine myths, direct evaluation of COVID-19 vaccine misinformation on the web, and national population surveys identify concerns about infertility related to COVID-19 vaccines as among the most common [33,43].

## Limitations

Several limitations need to be considered while interpreting the results of this study. First, it is important to acknowledge these findings are based on an analysis of search behavior in the United States and are likely to be influenced by specific vaccine approvals and distribution plans in this country. Future work is needed to determine how COVID-19–related vaccine search behavior might differ across countries. Second, these results apply only to search queries conducted using the Google search engine. Although Google is the dominant search engine in the United States, future work is needed to understand how search behavior described here is similar or different for other search engines or on other platforms. For example, social media has been identified as a major source of exposure to COVID-19–related misinformation. Our finding of relatively low rates of active searching for COVID-19 vaccine myths and conspiracy beliefs might or might not apply to search behavior within social media platforms. Third, it is important to acknowledge the subjective nature of our approach for identifying and defining search themes and categories. Other teams could have certainly chosen to identify or define search categories (or subcategories) in other ways. Fourth, of relevance, we also acknowledge the uncertainty regarding our interpretation of the observed patterns of online search behavior representing an increase in interest or intention among the study population to take the COVID-19 vaccine. Although we believe the major study finding that a large increase specifically in pharmacy-related COVID-19 vaccine queries strongly suggests active searching on the web for the vaccine, future work that directly assesses users' information needs (eg, near-time or real-time surveys) would be needed to confirm this interpretation. Finally, it is important to note that the results presented here are based on the aggregate search volume measured at the population level. We, therefore, are not able to determine the degree to which changing patterns in online search behavior are due to changes in the number of individuals performing a specific search, or due to an increase in the number of searches performed by specific individuals, or some combination of these factors.

## Conclusions

Despite these limitations, the findings presented here provide important information about the use of an infodemiologic approach to assess COVID-19 vaccine-related interest and intentions. During the study period, online search behavior related to COVID-19 vaccines suggested a possible historic high in public interest of vaccines. Furthermore, the specific type of vaccine related searches (eg, increased searches related to specific pharmacies and decreased searches related to vaccine side effects) is consistent with reduced vaccine hesitancy and greater intention to vaccinate. The relatively low occurrence of some types of searches (eg, COVID-19 vaccine myths and conspiracy beliefs) suggests that many individuals who lack or are uncertain about critical vaccine-related information are not engaged in active online search to address their information needs. Encouraging more active information-seeking, along with critical appraisal of health information on the web, could be an important strategy to combat misinformation about COVID-19 vaccines and increase vaccine confidence and intention to vaccinate among the general population.

## Conflicts of Interest

None declared.

## References

1.  Brørs G, Norman CD, Norekvål TM. Accelerated importance of eHealth literacy in the COVID-19 outbreak and beyond. Eur J Cardiovasc Nurs 2020 Aug 15;19(6):458-461 [FREE Full text] [doi: 10.1177/1474515120941307] [Medline: 32667217]

2.  Chong YY, Cheng HY, Chan HYL, Chien WT, Wong SYS. COVID-19 pandemic, infodemic and the role of eHealth literacy. Int J Nurs Stud 2020 Aug;108:103644 [FREE Full text] [doi: 10.1016/j.ijnurstu.2020.103644] [Medline: 32447127]

3.  Eysenbach G. How to fight an infodemic: the four pillars of infodemic management. J Med Internet Res 2020 Jun 29;22(6):e21820 [FREE Full text] [doi: 10.2196/21820] [Medline: 32589589]

4.  Rice RE. Influences, usage, and outcomes of Internet health information searching: multivariate results from the Pew surveys. Int J Med Inform 2006 Jan;75(1):8-28. [doi: 10.1016/j.ijmedinf.2005.07.032] [Medline: 16125453]

5.  Bangerter LR, Griffin J, Harden K, Rutten LJ. Health information-seeking behaviors of family caregivers: analysis of the health information national trends survey. JMIR Aging 2019 Jan 14;2(1):e11237 [FREE Full text] [doi: 10.2196/11237] [Medline: 31518309]

6.  Kontos E, Blake KD, Chou WS, Prestin A. Predictors of eHealth usage: insights on the digital divide from the Health Information National Trends Survey 2012. J Med Internet Res 2014 Jul 16;16(7):e172 [FREE Full text] [doi: 10.2196/jmir.3117] [Medline: 25048379]

7.  Sherman LD, Patterson MS, Tomar A, Wigfall LT. Use of digital health information for health information seeking among men living with chronic disease: data from the Health Information National Trends Survey. Am J Mens Health 2020 Jan 23;14(1):1557988320901377 [FREE Full text] [doi: 10.1177/1557988320901377] [Medline: 31973642]

8.  Eysenbach G. Infodemiology and infoveillance: framework for an emerging set of public health informatics methods to analyze search, communication and publication behavior on the Internet. J Med Internet Res 2009 Mar 27;11(1):e11 [FREE Full text] [doi: 10.2196/jmir.1157] [Medline: 19329408]

9.  Stephens-Davidowitz S. Google searches can help us find emerging COVID-19 outbreaks. The New York Times. 2020 Apr 05. URL: https://www.nytimes.com/2020/04/05/opinion/coronavirus-google-searches.html [accessed 2021-10-28]

10. Vasconcellos-Silva PR, Carvalho DBF, Trajano V, de La Rocque LR, Sawada ACMB, Juvanhol LL. Using Google Trends data to study public interest in breast cancer screening in Brazil: why not a pink February? JMIR Public Health Surveill 2017 Apr 06;3(2):e17 [FREE Full text] [doi: 10.2196/publichealth.7015] [Medline: 28385679]

11. Ginsberg J, Mohebbi MH, Patel RS, Brammer L, Smolinski MS, Brilliant L. Detecting influenza epidemics using search engine query data. Nature 2009 Feb 19;457(7232):1012-1014. [doi: 10.1038/nature07634] [Medline: 19020500]

12. Badell-Grau RA, Cuff JP, Kelly BP, Waller-Evans H, Lloyd-Evans E. Investigating the prevalence of reactive online searching in the covid-19 pandemic: infoveillance study. J Med Internet Res 2020 Oct 27;22(10):e19791 [FREE Full text] [doi: 10.2196/19791] [Medline: 32915763]

13. Higgins TS, Wu AW, Sharma D, Illing EA, Rubel K, Ting JY, Snot Force Alliance. Correlations of online search engine trends with coronavirus disease (COVID-19) incidence: infodemiology study. JMIR Public Health Surveill 2020 May 21;6(2):e19702 [FREE Full text] [doi: 10.2196/19702] [Medline: 32401211]

14. Jimenez AJ, Estevez-Reboredo RM, Santed MA, Ramos V. COVID-19 symptom-related google searches and local covid-19 incidence in spain: correlational study. J Med Internet Res 2020 Dec 18;22(12):e23518 [FREE Full text] [doi: 10.2196/23518] [Medline: 33156803]

15. Mattiuzzi C, Lippi G. Analysis of Google Searches for COVID-19 and its symptoms for predicting disease epidemiology in the United States. Acta Biomed 2020 Dec 04;92(1):e2021064 [FREE Full text] [doi: 10.23750/abm.v92i1.11070] [Medline: 33682813]

16. Mavragani A. Tracking COVID-19 in Europe: infodemiology approach. JMIR Public Health Surveill 2020 Apr 20;6(2):e18941 [FREE Full text] [doi: 10.2196/18941] [Medline: 32250957]

17. Rajan A, Sharaf R, Brown RS, Sharaiha RZ, Lebwohl B, Mahadev S. Association of search query interest in gastrointestinal symptoms with COVID-19 diagnosis in the United States: infodemiology study. JMIR Public Health Surveill 2020 Jul 17;6(3):e19354 [FREE Full text] [doi: 10.2196/19354] [Medline: 32640418]

18. Effenberger M, Kronbichler A, Shin JI, Mayer G, Tilg H, Perco P. Association of the COVID-19 pandemic with internet search volumes: a Google trends analysis. Int J Infect Dis 2020 Jun;95:192-197 [FREE Full text] [doi: 10.1016/j.ijid.2020.04.033] [Medline: 32305520]

19. Chen B, Chen X, Pan J, Liu K, Xie B, Wang W, et al. Dissemination and refutation of rumors during the COVID-19 outbreak in China: infodemiology study. J Med Internet Res 2021 Feb 15;23(2):e22427 [FREE Full text] [doi: 10.2196/22427] [Medline: 33493124]

20. Gerts D, Shelley CD, Parikh N, Pitts T, Watson Ross C, Fairchild G, et al. "Thought I'd share first" and other conspiracy theory tweets from the COVID-19 infodemic: exploratory study. JMIR Public Health Surveill 2021 Apr 14;7(4):e26527 [FREE Full text] [doi: 10.2196/26527] [Medline: 33764882]

21.    Hou Z, Du F, Zhou X, Jiang H, Martin S, Larson H, et al. Cross-country comparison of public awareness, rumors, and behavioral responses to the COVID-19 epidemic: infodemiology study. J Med Internet Res 2020 Aug 03;22(8):e21143 [FREE Full text] [doi: 10.2196/21143] [Medline: 32701460]

22.    Nsoesie EO, Cesare N, Müller M, Ozonoff A. COVID-19 misinformation spread in eight countries: exponential growth modeling study. J Med Internet Res 2020 Dec 15;22(12):e24425 [FREE Full text] [doi: 10.2196/24425] [Medline: 33264102]

23.    Rovetta A, Bhagavathula AS. COVID-19-related web search behaviors and infodemic attitudes in Italy: infodemiological study. JMIR Public Health Surveill 2020 May 05;6(2):e19374 [FREE Full text] [doi: 10.2196/19374] [Medline: 32338613]

24.    Rovetta A, Bhagavathula AS. Global infodemiology of COVID-19: analysis of Google web searches and Instagram hashtags. J Med Internet Res 2020 Aug 25;22(8):e20673 [FREE Full text] [doi: 10.2196/20673] [Medline: 32748790]

25.    Sajjadi NB, Nowlin W, Nowlin R, Wenger D, Beal JM, Vassar M, et al. United States internet searches for "infertility" following COVID-19 vaccine misinformation. J Osteopath Med 2021 Apr 12;121(6):583-587 [FREE Full text] [doi: 10.1515/jom-2021-0059] [Medline: 33838086]

26.    Hu Z, Yang Z, Li Q, Zhang A. The COVID-19 infodemic: infodemiology study analyzing stigmatizing search terms. J Med Internet Res 2020 Nov 16;22(11):e22639 [FREE Full text] [doi: 10.2196/22639] [Medline: 33156807]

27.    Husain I, Briggs B, Lefebvre C, Cline DM, Stopyra JP, O'Brien MC, et al. Fluctuation of Public Interest in COVID-19 in the United States: retrospective analysis of Google Trends search data. JMIR Public Health Surveill 2020 Jul 17;6(3):e19969 [FREE Full text] [doi: 10.2196/19969] [Medline: 32501806]

28.    Husnayain A, Shim E, Fuad A, Su EC. Understanding the community risk perceptions of the COVID-19 outbreak in South Korea: infodemiology study. J Med Internet Res 2020 Sep 29;22(9):e19788 [FREE Full text] [doi: 10.2196/19788] [Medline: 32931446]

29.    COVID-19 Vaccination. Centers for Disease Control and Prevention. URL: https://www.cdc.gov/vaccines/covid-19/index.html [accessed 2021-07-07]

30.    COVID-19 Vaccine Product Information. Centers for Disease Control and Prevention. URL: https://www.cdc.gov/vaccines/covid-19/info-by-product/index.html [accessed 2021-07-07]

31.    Wood S, Schulman K. Beyond politics — promoting Covid-19 vaccination in the United States. N Engl J Med 2021 Feb 18;384(7):e23. [doi: 10.1056/nejmms2033790]

32.    Pullan S, Dey M. Vaccine hesitancy and anti-vaccination in the time of COVID-19: a Google Trends analysis. Vaccine 2021 Apr 01;39(14):1877-1881 [FREE Full text] [doi: 10.1016/j.vaccine.2021.03.019] [Medline: 33715904]

33.    Islam MS, Kamal AM, Kabir A, Southern DL, Khan SH, Hasan SMM, et al. COVID-19 vaccine rumors and conspiracy theories: the need for cognitive inoculation against misinformation to improve vaccine adherence. PLoS One 2021 May 12;16(5):e0251605 [FREE Full text] [doi: 10.1371/journal.pone.0251605] [Medline: 33979412]

34.    Search Engine Market Share in United State of America-June 2021. Global Stats - Stat Counter. URL: https://gs.statcounter.com/search-engine-market-share/all/united-states-of-america [accessed 2021-07-07]

35.    Arora VS, McKee M, Stuckler D. Google Trends: opportunities and limitations in health and health policy research. Health Policy 2019 Mar;123(3):338-341. [doi: 10.1016/j.healthpol.2019.01.001] [Medline: 30660346]

36.    Silvestri F. Mining query logs: turning search usage data into knowledge. Foundations and Trends in Information Retrieval 2009 Nov 29;4(1–2):1-174. [doi: 10.1561/1500000013]

37.    Kulkarni A, Teevan J, Svore K, Dumais ST. Understanding temporal query dynamics. In: WSDM '11.: Association for Computing Machinery; 2011 Presented at: Proceedings of the Fourth ACM International Conference on Web Search and Data Mining; February 2011; Hong Kong, China p. 167-176. [doi: 10.1145/1935826.1935862]

38.    Dumais S, Jeffries R. Understanding user behavior through log data analysis. In: Olson J, Kellogg W, editors. Ways of Knowing in Human-Computer Interaction. New York: Springer-Verlag; 2014:349-372.

39.    Malik MT, Gumel A, Thompson LH, Strome T, Mahmud SM. "Google flu trends" and emergency department triage data predicted the 2009 pandemic H1N1 waves in Manitoba. Can J Public Health 2011 Jul 1;102(4):294-297. [doi: 10.1007/bf03404053]

40.    Sousa-Pinto B, Anto A, Czarlewski W, Anto JM, Fonseca JA, Bousquet J. Assessment of the impact of media coverage on COVID-19-related Google Trends data: infodemiology study. J Med Internet Res 2020 Aug 10;22(8):e19611 [FREE Full text] [doi: 10.2196/19611] [Medline: 32530816]

41.    Huynh Dagher S, Lamé G, Hubiche T, Ezzedine K, Duong TA. The influence of media coverage and governmental policies on Google queries related to COVID-19 cutaneous symptoms: infodemiology study. JMIR Public Health Surveill 2021 Feb 25;7(2):e25651 [FREE Full text] [doi: 10.2196/25651] [Medline: 33513563]

42.    KFF COVID-19 Vaccine Monitor: March 2021. Kaiser Family Foundation. 2021. URL: https://www.kff.org/coronavirus-covid-19/poll-finding/kff-covid-19-vaccine-monitor-march-2021/ [accessed 2021-10-07]

43.    KFF COVID-19 Vaccine Monitor: January 2021. Kaiser Family Foundation. 2021 Jan 21. URL: https://www.kff.org/report-section/kff-covid-19-vaccine-monitor-january-2021-vaccine-hesitancy/ [accessed 2021-10-07]

**Abbreviations**

   **FDA:** Food and Drug Administration

**RSV:** relative search volume

XSL·FO
**RenderX**

Original Paper

# Temporal Variations and Spatial Disparities in Public Sentiment Toward COVID-19 and Preventive Practices in the United States: Infodemiology Study of Tweets

Alexander Kahanek[1], BSc; Xinchen Yu[1], MSc; Lingzi Hong[1], PhD; Ana Cleveland[1], PhD; Jodi Philbrick[1], PhD

College of Information, University of North Texas, Denton, TX, United States

**Corresponding Author:**
Lingzi Hong, PhD
College of Information
University of North Texas
E292
3940 N Elm St
Denton, TX, 76203
United States
Phone: 1 9192607578
Email: lingzi.hong@unt.edu

## Abstract

**Background:** During the COVID-19 pandemic, US public health authorities and county, state, and federal governments recommended or ordered certain preventative practices, such as wearing masks, to reduce the spread of the disease. However, individuals had divergent reactions to these preventive practices.

**Objective:** The purpose of this study was to understand the variations in public sentiment toward COVID-19 and the recommended or ordered preventive practices from the temporal and spatial perspectives, as well as how the variations in public sentiment are related to geographical and socioeconomic factors.

**Methods:** The authors leveraged machine learning methods to investigate public sentiment polarity in COVID-19–related tweets from January 21, 2020 to June 12, 2020. The study measured the temporal variations and spatial disparities in public sentiment toward both general COVID-19 topics and preventive practices in the United States.

**Results:** In the temporal analysis, we found a 4-stage pattern from high negative sentiment in the initial stage to decreasing and low negative sentiment in the second and third stages, to the rebound and increase in negative sentiment in the last stage. We also identified that public sentiment to preventive practices was significantly different in urban and rural areas, while poverty rate and unemployment rate were positively associated with negative sentiment to COVID-19 issues.

**Conclusions:** The differences between public sentiment toward COVID-19 and the preventive practices imply that actions need to be taken to manage the initial and rebound stages in future pandemics. The urban and rural differences should be considered in terms of the communication strategies and decision making during a pandemic. This research also presents a framework to investigate time-sensitive public sentiment at the county and state levels, which could guide local and state governments and regional communities in making decisions and developing policies in crises.

## Introduction

### Background

The COVID-19 pandemic has had worldwide economic and mortality impacts, with more than 118 million confirmed cases and over 2.6 million deaths globally as of March 12, 2021 [1].

Since the initial outbreak of COVID-19, many public health professionals and authoritative organizations, such as the Centers for Disease Control and Prevention (CDC) and the World Health Organization, have recommended that people change their fundamental behaviors of daily life to prevent the virus from spreading, for example, wearing masks, social distancing, and restricting travel [2]. However, the effectiveness of these

measures in reducing the spread hinges on compliance by the public. The level of compliance varies among citizens in following the suggested practices. In the United States, there are divergent opinions about the preventive practices, which have existed from the onset of the CDC guidelines.

## Prior Work

It is critical to gauge public sentiment and responses to the preventive practices for effective communication strategies, decisions, and policies, as disparities in practices may affect the spread of the disease and delay society's recovery from the pandemic. Social media has been widely adopted by people to acquire information and share opinions in crises, which provides time-sensitive opportunities for governments and public institutions to understand public opinions. Social media data have been used as crowd sources of information to understand citizens' issues of concern [3,4], response to policies [5,6], and emotional consequences [7] in crises. Several recent studies have used Twitter and Facebook data for closer-to-real-time infodemiology studies, for example, to analyze emotions concerning the lockdown [8] and reopening [9] and to understand COVID-19 discussions and the associated sentiments [10]. However, these studies usually rely on an implicit assumption that strategies based on the understanding of the whole society at a time or during a time range work for all. Some studies have investigated the evolvement of public responses as the crisis unfolded, for example, the content analysis of crisis-related tweets before, during, and after the crisis [11]; temporal variations of public sentiment toward COVID-19 in China [12]; and changes in risk perception of COVID-19 in the United States in the early stage of the pandemic [13]. Several studies have examined the spatial differences. For example, Ntompras et al [14] conducted comparisons of the content of Twitter posts related to the COVID-19 pandemic across nations. They found several topics were triggered by local events, which implies that social media data can act as political, economic, and social monitoring in pandemics. Cuomo et al [15] performed a more granular analysis and investigated the longitudinal and geospatial relationships between volumes of self-reporting COVID-19 cases and elevated risks of virus spreading in the United States at the county level. Similar studies have found geolocated tweets on COVID-19 symptoms, concerns, and experiences are indicative of officially reported COVID-19 cases at the county level in the United States [16] and volumes of misinformation are related to increased COVID-19 cases at the state and county level in the United States [17]. Currently, few have investigated the temporal variations in public sentiment in a high geospatial resolution. Hou et al [18] found that mobility behaviors differ in communities during COVID-19, which could be related to various socioeconomic and cultural factors. Schmelz [19]

conducted a survey study in Germany and found that people with different levels of trust in the government or with different political identities may have varying reactions in their response to government policies during COVID-19. These studies suggest it is important to consider the heterogeneity of the population in public health decision making. Methods on the time-sensitive understanding of a crisis with social media data have seldomly considered the geographical disparities and the associated socioeconomic factors. This study aimed to address this gap by proposing a social media data analysis framework for a longitudinal investigation of US public sentiment about COVID-19 and the preventive practices on different spatial scales.

## Goal of This Study

The focus of the study was to identify the variations in public sentiment toward COVID-19 and the preventive practices from the temporal and spatial perspectives and to investigate how the variations are related to geographical and socioeconomic factors in the United States. Specifically, we analyzed discussions of COVID-19 on Twitter in the United States to answer the following questions:

1. Research question 1: Are there temporal variations in public sentiment toward overall COVID-19 issues and preventive practices?
2. Research question 2: Are there spatial disparities in public sentiment toward overall COVID-19 issues and preventive practices?
3. Research question 3: What geographic factors may be related to the differences in public sentiment toward COVID-19 issues and preventive practices?
4. Research question 4: What socioeconomic factors may be related to the differences in public sentiment toward COVID-19 issues and preventive practices?

Exploring these 4 questions could offer rich insight into public sentiment about COVID-19 issues and preventive practices, with fine temporal and spatial granularity. This allows policy makers to explicitly consider these variations in developing communication strategies or adjusting enforcement policies for efficient coordination in pandemics or crises like COVID-19. This study sets the groundwork for analyzing, comparing, and potentially predicting public sentiment in future crises.

## *Methods*

The method was composed of 3 parts (Figure 1): data collection, data preparation, and data analysis. In this study, we collected and analyzed Twitter data on COVID-19 and the preventive practices.

**Figure 1.** Data analysis framework. API: application programming interface; PP: preventive practices; SE: socioeconomic.



## Data Collection

COVID-19 cases were first reported in Wuhan, China in late December 2019. The disease was fast-spreading and led to increasing infections and deaths globally. Starting in January 2020, other countries started to report confirmed cases of COVID-19. To retrieve online discussions on COVID-19, we collected a Twitter data set with about 160,000,000 tweets containing COVID-19–related keywords starting from January 21, 2020. The list of keywords includes "Coronavirus," "Corona," "CDC," "Covid19," "Covid19," "Sarscov2," "pandemic," "epidemic," and their variants [20]. The data were collected using Python through the Twitter's streaming application programming interface (API). The Twitter streaming API returns 1% of the total Twitter volume, with multilingual tweets posted from around the world. As we were focused on the public sentiment in the United States, only tweets in English were kept.

## Ethics Statement

In compliance with Twitter policy, we removed identifiers from the data before analysis to avoid potential profiling or targeting of individuals. We only present aggregated analyses. To support reproducibility, the tweet IDs, processing code, and intermediate results will be available upon request to the corresponding author.

## Data Preparation

Data preparation was composed of 4 parts (Figure 1): the geographical projection of tweets, the identification of posts from individual users, the subsetting of topics on preventative practices, and sentiment detection. All the data preparation was implemented with Python.

### Geographical Projection

The collected tweets only satisfied the condition of semantic relevance to COVID-19, of which some had embedded geolocations, such as a point location, a bounding box defined with geographical coordinates, or user-entered location tags. Although many tweets contain location tags, such tags often vary in geographical scale or do not refer to real locations. Therefore, only tweets with geographical coordinates, either as a point location or a bounding box, were used. We used the GeoPandas package in Python for all geographical data

processing. After calculating the center of the bounding box, they were projected to the coordinate system of the Shapefile map of the United States [21] and then matched with the geographical units at the county and state levels. If the location of a tweet fell in a county, we assigned the tweet with the associated county and state, together with the aggregated socioeconomic information from the US Census Bureau [22]. In addition, we used the urban/rural map Shapefile to identify whether a tweet was posted from an urban or rural area [23]. Urban areas included "Urbanized Areas of 50,000 or more people and Urban Clusters of at least 2,500 and less than 50,000 people" [24]. Other areas were classified as rural. After geographical projection and filtering for tweets only posted in the English language that were located within the United States, there was a total of 344,227 tweets.

### Identification of Posts From Individual Users

Different types of users are on Twitter, including media outlets, accounts of government authorities and organizations, social bots, and individual accounts. The quality of data and the generated insights may be impacted by the activities of bots and official accounts [7]. The first step of geographical projection left tweets that were highly probable to be from individual users. To assure that the tweets used for analysis were mainly from individual users, we applied the traditional approach by checking the social relations of the authors, assuming that media outlets and bots usually have a high ratio between their followers and friends:



Specifically, we identified users whose numbers of followers and followees were 2 SDs larger than the average values as nonpersons [25]. We found 9 tweets by media outlets or bots. The result confirmed the filtered geolocated tweets to be mainly from individual users. After filtering for individual users, the COVID-19 data set included 344,218 tweets.

### Tweets on COVID-19 Preventive Practices

We were specifically interested in the public sentiment toward COVID-19 prevention practices, as people's compliance to these practices highly affect the spread of the disease. To identify the potential keywords that describe COVID-19 prevention practices, we collected all the guidelines released by the CDC

[2]. Three graduate research assistants read through the documents and identified keywords and phrases that were relevant to the preventive behaviors for reducing the spread of the disease. Specifically, 4 categories of practices were collected, including physical or social distancing, personal protective equipment (PPE), disinfection, and other. Physical or social distancing included social distancing, social distance, physical distance, 6-feet, stay-at-home, school isolation, isolation, stay home, avoid touching. PPE included mask, covering, face shield, wear a mask, surgical mask, N95 respirator, wearing gloves, face shields, facial covering, skin protection, eye protection, PPE. Disinfection included wash hands, hand sanitizer, disinfect, clean, detergent, handwash, hand hygiene, prevention hygiene, sprays, concentrates, wipes, routine cleaning, bleach solution. Others included test, business closure.
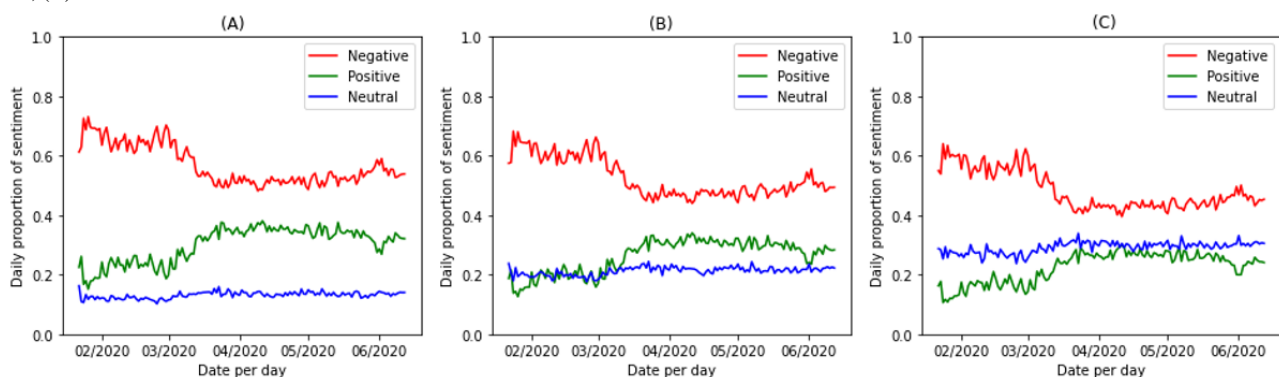
These keywords and phrases were used to identify if a tweet was about COVID-19 preventive practices and the category of preventive practices. As the language used in social media posts may have syntax or typographical errors, using the formal keywords and phrases from CDC guidelines may affect the recall of tweets on preventive practices. Therefore, we applied token normalization for both tweets and the keywords and phrases. After the normalization of each tweet, we checked if any tokens in a tweet matched the normalized keywords or phrases. Tweets containing these keywords or phrases were aggregated to form the subset on preventive practices (shortened to CDC subset in the following analysis), which had a total of 53,272 tweets. Based on the tokens, the tweets were further categorized as discussions on 1 of the 4 categories. The top keywords found in the COVID-19 data set were mask, stay home, social distancing, test, and PPE. These individual keywords had more than 8000 occurrences.

## Sentiment Detection

We used a pretrained deep learning model, FLAIR, to detect the sentiment contained in each tweet [26]. The model was constructed with the recurrent neural network architecture, which enables the capture of semantic and syntactic information of words and the surrounding context for the prediction of the sentiment of input text. As the model was designed to capture different meanings for polysemous words and handle rare and misspelled words with ease, it works well for a Twitter corpus, where words are often misspelled and have ambiguous meanings. The model had state-of-the-art performance in sentiment classification, with an accuracy of 89.5% and an F1 score of 0.89 on a separate data set [27].

For each input tweet, the output was a sentiment of 1 of 2 categories: positive or negative with the associated confidence of the model's prediction. However, not all tweets include the expression of sentiment. In fact, about 25% of crisis-related Twitter data do not contain subjective information [28]. Tweets with a low confidence to a sentiment category are probable to be neutral or objective. As each sentiment category has a confidence between 0 and 1, we explored 3 thresholds (0.8, 0.9, and 0.95) to understand whether the choice of thresholds would affect the temporal variations of sentiment. Figure 2 shows the ratio of COVID-19 tweets that are positive, negative, and neutral on a daily level, when the confidence thresholds of 0.8, 0.9, and 0.95 were used to define neutral tweets. We found that the choice of threshold would not significantly impact the temporal patterns of positive or negative sentiment in the COVID-19 data set. To obtain more samples for analysis, we chose the confidence level of 0.8 and considered tweets with a confidence level of positive or negative less than 0.8 as neutral.

**Figure 2.** The daily proportion of sentiment in the COVID-19 data set when neutral tweets were detected with different confidence thresholds: (A) 0.8, (B) 0.9, (C) 0.95.



## Summary of the Data Set

Table 1 shows the summary statistics of the COVID-19 data set and the CDC subset for analysis. Both data sets have tweets from 50 states plus Washington DC and other territories of the United States in the time range from January 21, 2020 to June 12, 2020.

**Table 1.** Descriptive summary of the COVID-19 data set and the Centers for Disease Control and Prevention (CDC) subset.

| Characteristics | COVID-19 | CDC |
|---|---|---|
| Tweets, n | 344,218 | 53,272 |
| Negative sentiment, n (%) | 195,166 (56.7) | 32,408 (60.8) |
| Positive sentiment, n (%) | 103,698 (30.1) | 13,411 (25.2) |
| Neutral sentiment, n (%) | 45,354 (13.2) | 7453 (14.0) |
| Tweets per day, mean (SD) | 2424 (1488.96) | 375 (294.55) |
| Tweets per week, mean (SD) | 16,391 (6935.53) | 6935 (1529.67) |

## Analysis of Temporal Variations and Spatial Disparities

We conducted temporal analysis to answer research question 1. First, we computed the ratio of tweets with positive, negative, or neutral sentiment separately in the granularity of a day for the COVID-19 data set and the CDC subset and weekly for each category of the preventive practices (ie, physical or social distancing, PPE, disinfection, and others). The time series of public sentiment were analyzed with an algorithm that helped to detect the turning points when the sentiment patterns started to change. The turning points and the nearby dates were investigated to explore what events might be related to the significant changes in public sentiment polarity.

The enforcement policies issued by state and local governments and the dates of intervention could be different in the United States. The enforcement may trigger changes in public sentiment to preventive practices [29]. We conducted county- and state-level analyses to examine the spatial disparities. We aggregated the tweets by states and generated sentiment polarity maps for the COVID-19 data set and CDC subset. Further, 4 representative states were analyzed to investigate the dynamics of public sentiment at a finer spatial granularity, which enabled analysis of whether the changes in public sentiment related to state-level events or policies.

The variations in public sentiment in different regions could be related to the heterogeneity of the population, as existing studies have shown the response behaviors in COVID-19 are related to cultural, socioeconomic, and political factors [19,29,30]. Two types of analysis were implemented to answer research questions 3 and 4. The analysis was conducted at the aggregated county level. We did not investigate aggregation at a finer spatial granularity, such as by census tracts or census block groups, due to the sparsity of tweets with geolocations. Using a smaller geographical unit means fewer tweet samples in each unit, which could be easily affected by sentiment detection errors.
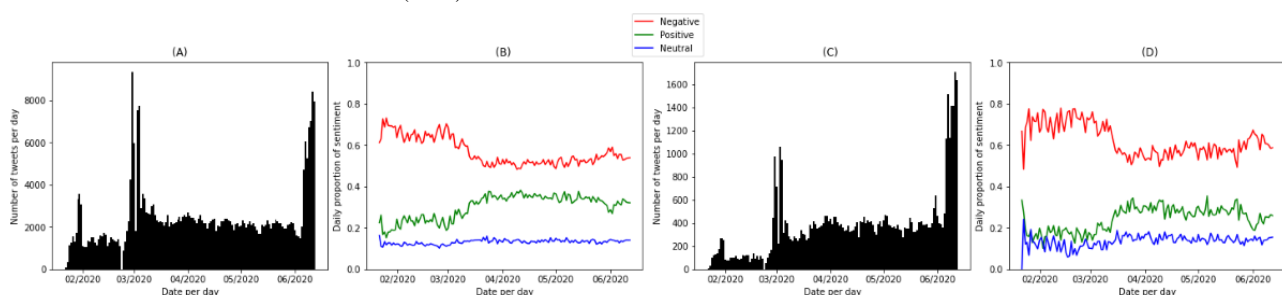
First, we compared the public sentiment polarity in the urban and rural areas of counties to answer research question 3. For each county, we obtained the ratio of tweets with positive, negative, and neutral sentiments separately for urban and rural areas. We ran a *t* test between the sentiment polarity in the urban and rural areas of counties to find out if the urban/rural factor would explain the variances in the public sentiment to COVID-19 and preventive practices. Second, we examined if the sentiment polarity to COVID-19 and preventive practices were statistically related to the socioeconomic factors for research question 4. The socioeconomic information was obtained from the *U.S. Census Bureau Indicators of the 2017 American Community Survey 5-year Estimate* [22].

## Results

### Temporal Variations in Public Sentiment Toward COVID-19 and Preventive Practices

Figure 3 presents the visualization of volumes and sentiment polarity separately for tweets in the COVID-19 data set and the CDC subset. Tweets on preventive practices represented about 15.5% (53,272/344,218) of COVID-19 tweets. The 2 timelines on volumes showed a common pattern and had 2 large spikes at similar time points: one in the beginning of March 2020 and another in the middle of June 2020. The timing of these 2 spikes corresponded to the turning points when public sentiment polarity started to change. Since the beginning of March 2020, the negative sentiment about COVID-19 issues and preventive practices started to decrease, although the second spike in June 2020 was associated with increasing negative sentiment in both data sets.

**Figure 3.** The numbers of tweets and proportions of sentiment per day in (A) and (B), respectively, the COVID-19 data set and (C) and (D), respectively, the Centers for Disease Control and Prevention (CDC) subset.

XSL•FO

**RenderX**

The daily proportion of neutral tweets had little variation over the studied time range; the time series of positive sentiment was almost mirrored to that of the negative sentiment. Therefore, we focused on the analysis of the negative sentiment. Figure 4 presents the visualization of the 4 stages indicated by 4 different colors. The turning points about COVID-19 were March 6, 2020; March 29, 2020; and April 30, 2020. Table 2 and Table 3 show the summary statistics of the 4 stages for the COVID-19 data set and the CDC subset.

The dynamics of negative sentiment in the COVID-19 data set and CDC subset shared similar patterns, except that the timing of turning points varied. The COVID-19 data set and the CDC subset both had a high proportion of negative sentiment in Stage 1. The mean daily proportions of negative tweets were 66.6% (59,805/89,757) in the COVID-19 data set and 70.7% (8107/11,475) in the CDC subset. In Stage 2, there was a consistent decline in the negative sentiment in both the COVID-19 data set and the CDC subset, although the turning point of the COVID-19 time series (March 5, 2020) came earlier than that of the CDC subset (March 15, 2020). After a certain amount of time, the decreasing trend stopped and reached another turning point. In Stage 3, the negative proportion remained stable in the COVID-19 data set. The average negative proportion (37,350/72,849, 51.3%) was lower than in Stage 1 (59,805/89,757, 66.6%) and Stage 2 (32,062/58,010, 55.3%) in the COVID-19 data set. Comparatively, there were more variations in the negative sentiment in the CDC subset. People showed increasing negative sentiment toward preventive practices in Stage 3 (4945/8610, 57.4%) and Stage 4 (9937/16,328, 60.9%) after Stage 2 (9419/16,859, 55.9%). There was also an increasing trend in negative sentiment toward general COVID-19 issues in Stage 4 (65,954/123,602, 53.4%). In all stages, the sentiment polarities in the CDC subset were higher than those in the COVID-19 data set.

There were similar trends by categories. The negative sentiment polarities were the highest at the beginning of COVID-19 from January. Then, the proportion of negative tweets decreased until later May 2020 when negative sentiment started to rebound. There were also noticeable differences. For example, the sentiment for the disinfection topic had a small spike in April 2020, which could be related to the criticism of "disinfectant injection." Overall, public sentiment to the PPE topic (11,157/19,640, 56.8%) was more negative than to physical or social distancing (10,684/19,466, 54.9%) and disinfection (1706/3061, 55.8%).

To further investigate public sentiment toward different categories of preventive practices, we generated the timelines of public sentiment polarity separately for tweets in the 4 categories (ie, 19,640 tweets on PPE, 19,466 tweets on physical or social distancing, 3061 tweets on disinfection, and 16,425 tweets on other measurements). As we investigated more detailed granularity, there were fewer representative samples at the daily level, which led us to adjust the aggregation from the daily level to the weekly level. Figure 5 shows the visualization of weekly volumes and sentiment polarity of tweets in the 4 categories of preventive practices.

**Figure 4.** Stage splitting for the (A) COVID-19 data set and (B) Centers for Disease Control and Prevention (CDC) subset.
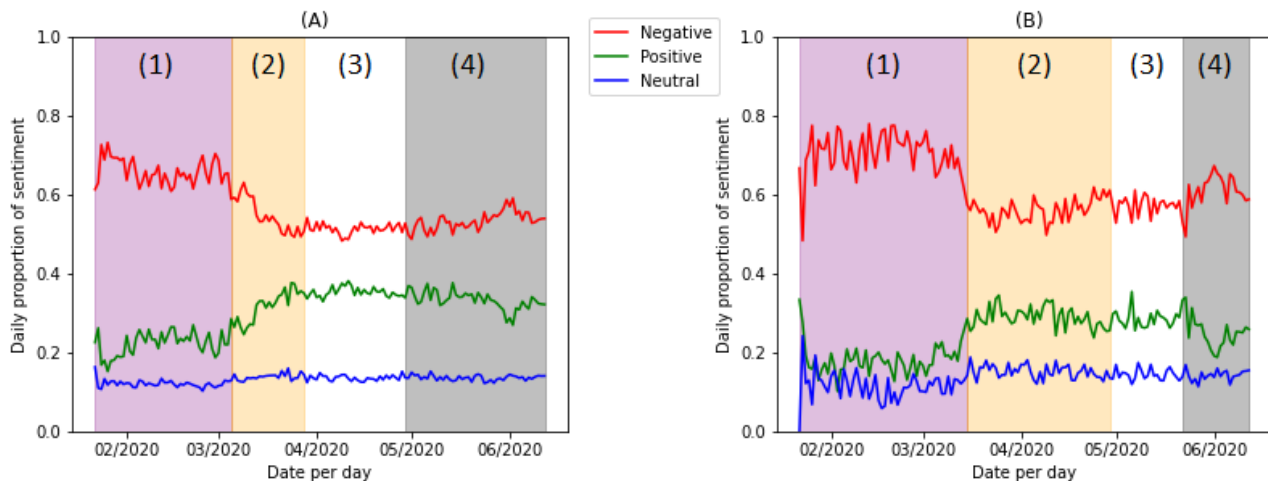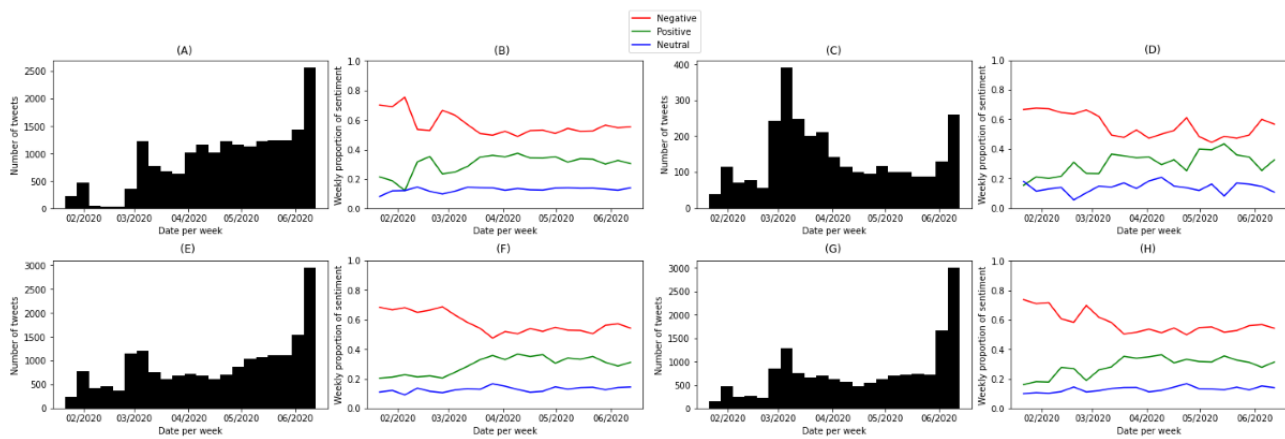


**Table 2.** Summary statistics for the 4 stages of the COVID-19 data set.

| Stage | Date range | Volume | Negative sentiment | Positive sentiment | Neutral sentiment |
|---|---|---|---|---|---|
| 1 | January 21, 2020 to March 5, 2020 | 89,757 | 0.6663 | 0.2140 | 0.1197 |
| 2 | March 6, 2020 to March 28, 2020 | 58,010 | 0.5527 | 0.3102 | 0.1371 |
| 3 | March 29, 2020 to April 29, 2020 | 72,849 | 0.5127 | 0.3521 | 0.1352 |
| 4 | April 30, 2020 to June 12, 2020 | 123,602 | 0.5336 | 0.3304 | 0.1360 |
| Total | January 21, 2020 to June 12, 2020 | 344,218 | 0.5669 | 0.3013 | 0.1318 |

**Table 3.** Summary statistics for the 4 stages of the Centers for Disease Control and Prevention (CDC) subset.

| Stage | Date range | Volume | Negative sentiment | Positive sentiment | Neutral sentiment |
|---|---|---|---|---|---|
| 1 | January 21, 2020 to March 5, 2020 | 11,475 | 0.7065 | 0.1780 | 0.1155 |
| 2 | March 6, 2020 to March 28, 2020 | 16,859 | 0.5587 | 0.2895 | 0.1518 |
| 3 | March 29, 2020 to April 29, 2020 | 8610 | 0.5743 | 0.2825 | 0.1432 |
| 4 | April 30, 2020 to June 12, 2020 | 16,328 | 0.6086 | 0.2484 | 0.1430 |
| Total | January 21, 2020 to June 12, 2020 | 53,272 | 0.6084 | 0.2517 | 0.1399 |

**Figure 5.** The number of tweets and proportion of sentiment per week in the Centers for Disease Control and Prevention (CDC) subset, split by themes: (A) and (B), respectively, physical or social distancing; (C) and (D), respectively, disinfection; (E) and (F), respectively, personal protective equipment; (G) and (H), respectively, others.
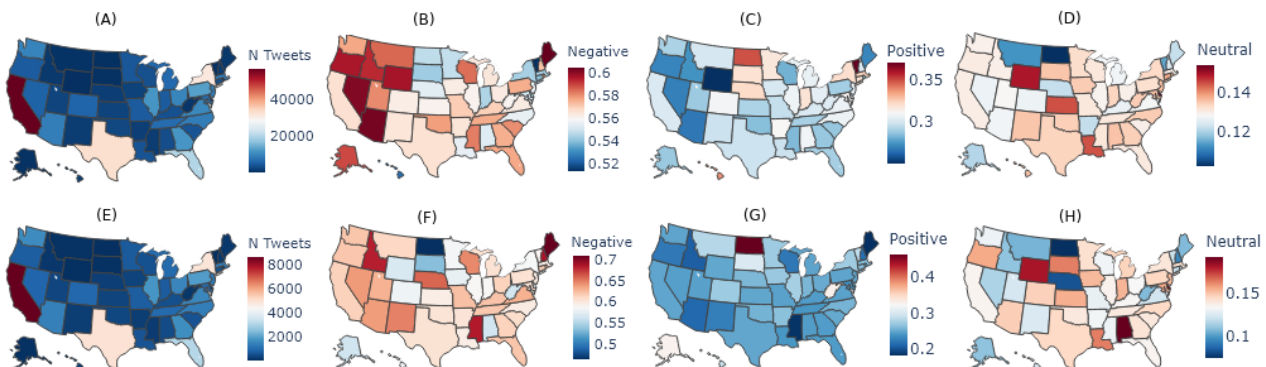


## Spatial Disparities at the State Level

Figure 6 shows the number of tweets on COVID-19 and preventive practices at the state level in the United States; 4 states had the largest number of tweets on COVID-19 and preventive practices: California, New York, Texas, and Florida. These 4 states are the most populated states in the United States according to the US Census Bureau [22]. The counts of tweets, per state, were proportionally similar between the COVID-19 and CDC data sets, ensuring that changes in sentiment between the 2 data sets were not due to geographical sampling differences. The sentiment polarity map of the COVID-19 data set showed that negative sentiment was highest in Maine and some of the states in the Pacific Western area including Arizona, Nevada, Wyoming, Oregon, and Idaho. More research needs to be done to investigate why the negative sentiment presented such a geographic pattern. On the other hand, states with the most negative sentiment on CDC were geographically dispersed. The top 3 states included Maine, New Hampshire, and Mississippi.

**Figure 6.** The number of tweets and proportion of sentiment for each state across the entire timeline: (A) count of tweets in the COVID-19 data set; (B) proportion of negative sentiment in the COVID-19 data set; (C) proportion of positive sentiment in the COVID-19 data set; (D) proportion of neutral sentiment in the COVID-19 data set; (E) count of tweets in the Centers for Disease Control and Prevention (CDC) subset; (F) proportion of negative sentiment in the CDC subset; (G) proportion of positive sentiment in the CDC subset; (H) proportion of neutral sentiment in the CDC subset.



Further, we chose 4 states with the highest volumes of tweets (ie, California, n=56,188; Texas, n=32,890; New York, n=31,178; and Florida, n=19,965) for temporal analysis. Figure 7 shows the volume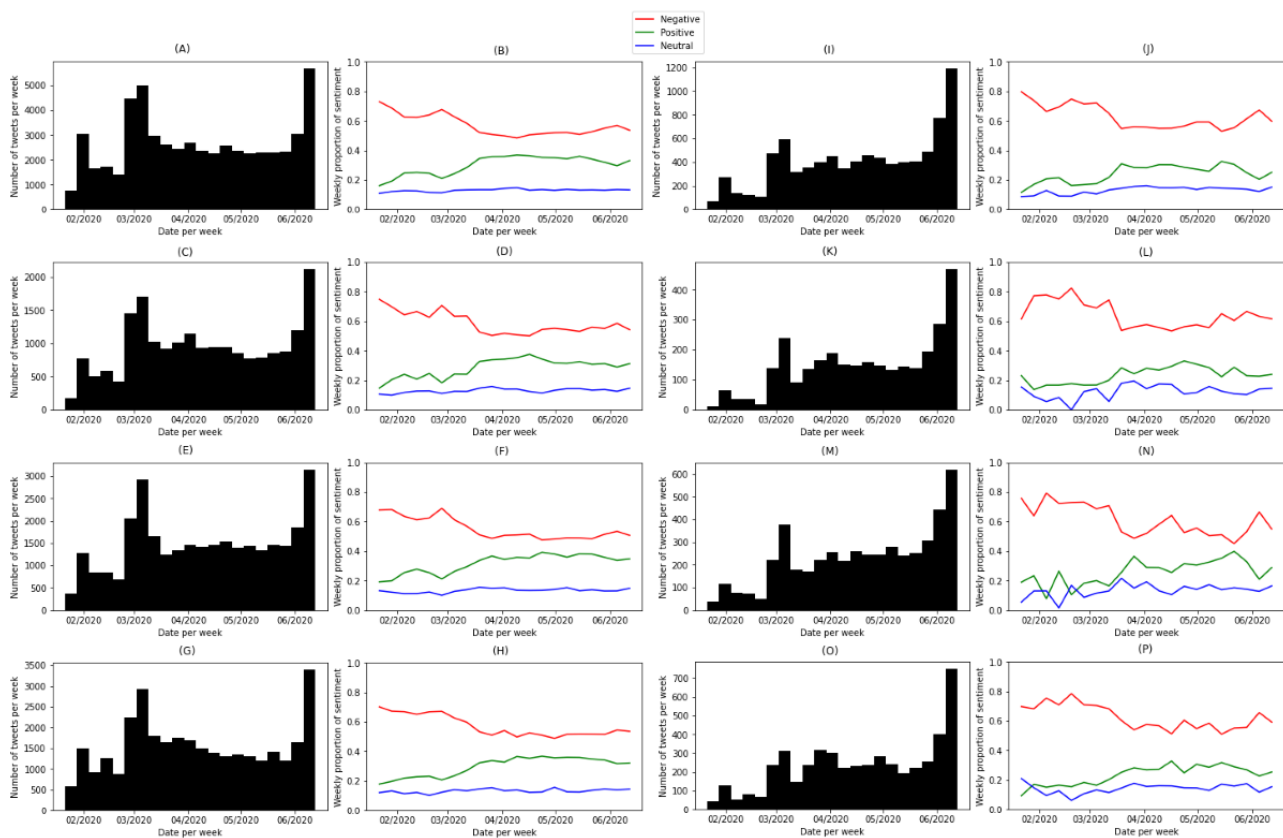s and sentiment polarity toward COVID-19 issues and preventive practices at the weekly level for the 4 states. The timelines of the 4 states demonstrated similar patterns and were close to the general trend in the United States. We observed state differences in many places. Florida

(11,554/19,965, 57.9%) showed more negative sentiment to COVID-19 issues than the other 3 states: California (31,926/56,188, 56.8%), Texas (18,682/32,890, 56.8%), and New York (17,020/31,178, 54.6%). The starting point of Stage 4, when negative sentiment started to increase, came early in Florida (approximately late April 2020), while for New York and California, Stage 4 started in the middle of May 2020.

There were more variations in public sentiment toward preventive practices among the 4 states. Overall, there was a higher proportion of negative tweets in Florida (1916/3087, 62.1%) than in California (5284/8588, 61.5%), Texas (2991/4949, 60.4%), and New York (2850/4865, 58.6%) in the CDC subset. California and Florida shared similar trends, where the timeline started with high ratios of negative tweets, which lasted until the middle of March 2020 and stayed relatively low

and increased at the later stage. New York was different in that the public sentiment to preventive practices seemed to vary greatly in the timeline. The proportion of negative sentiment decreased to almost 40% in the middle of March 2020 and started to increase to about 60%, then decreased to converge to almost 40% and increased at the later stage. There was a spike in negative sentiment in the middle of April 2020. After checking the term frequency-inverse document frequency (TFIDF) of keywords, we found keywords related to political figures, the Black Lives Matter movement, and various current events. This shows that topics were not solely related to COVID-19 nor preventative practices and tweets' sentiments may have additional influence from other topics. In Texas, there was decreasing negative sentiment until the middle of April 2020, when there was a spike in negative sentiment. Following that, the negative sentiment increased gradually.

**Figure 7.** The number of tweets and proportion of sentiment, respectively, per week in the COVID-19 data set in (A) and (B) California, (C) and (D) Florida, (E) and (F) New York, (G) and (H) Texas and in the Centers for Disease Control and Prevention (CDC) subset split in (I) and (J) California, (K) and (L) Florida, (M) and (N) New York, (O) and (P) Texas.



## Sentiment Polarity in Urban and Rural Areas

We conducted *t* tests to compare the public sentiment toward COVID-19 and preventive practices in the urban and rural areas of counties. The results are shown in Table 4. Counties were split into their respective urban and rural areas. After filtering urban and rural counties that did not have at least 15 tweets, 830 counties with urban areas and 182 counties with rural areas remained in the COVID-19 data set, and 355 urban and 52 rural counties remained in the CDC subset.

The *t* tests showed that there were no significant differences toward COVID-19–related issues. However, public discussions on preventive practices (CDC subset) were significantly more negative in rural areas (mean 0.6543, SD 0.0785) than in urban areas (mean 0.6112, SD 0.0891; $t_{405}$=−3.6332, $P<.001$). Additionally, we observed more positive sentiment for people in urban (mean 0.2454, SD 0.0822) than in rural areas (mean 0.2173, SD 0.0716; $t_{405}$=2.5976, $P=.01$) and a higher proportion of neutral posts in urban areas (mean 0.1433, SD 0.0565) than in rural areas (mean 0.1283, SD 0.0411; $t_{405}$=2.313, $P=.02$).

XSL•FO

RenderX

**Table 4.** Comparison of sentiment polarity between urban and rural areas.

| Sentiment | COVID-19 data set | | | | CDC[a] subset | | | |
|---|---|---|---|---|---|---|---|---|
| | Urban, mean (SD) | Rural, mean (SD) | $t$ value | $P$ value | Urban, mean (SD) | Rural, mean (SD) | $t$ value | $P$ value |
| Negative | 0.5768 (0.0897) | 0.5833 (0.1064) | –0.7691 | .44 | 0.6112 (0.0891) | 0.6543 (0.0785) | –3.6332 | <.001 |
| Positive | 0.2921 (0.0827) | 0.2918 (0.0995) | 0.0397 | .97 | 0.2454 (0.0822) | 0.2173 (0.0716) | 2.5976 | .01 |
| Neutral | 0.1310 (0.0505) | 0.1248 (0.0588) | 1.3211 | .19 | 0.1433 (0.0565) | 0.1283 (0.0411) | 2.3213 | .02 |

[a]CDC: Centers for Disease Control and Prevention.

### Sentiment Polarity and Socioeconomic Factors

The variances in public sentiment toward COVID-19 and preventive practices were then examined with the socioeconomic indicators of poverty and unemployment rates, as well as the median household income. The normality of these 3 variables and sentiment polarity was checked. Table 5 presents the distributions of all socioeconomic factors, the sentiment values, and the average tweet populations for counties.

Table 6 presents the Pearson correlation results. The unemployment rate was positively correlated with the proportion of negative sentiment ($r_{907}=0.0982$, $P=.003$) and negatively correlated with the proportion of positive sentiment ($r_{907}=-0.1407$, $P<.001$) in the COVID-19 data set. It means counties with higher unemployment rates had higher negative sentiment polarity toward COVID-19 issues. Similarly, counties with higher poverty rates tended to have a lower proportion of positive discussions on COVID-19 issues ($r_{907}=-0.0836$, $P=.01$). Finally, median household income was negatively correlated with the proportion of negative sentiment ($r_{907}=-0.1322$, $P<.001$) and positively correlated with the proportion of positive sentiment ($r_{907}=0.1554$, $P<.001$) in the COVID-19 data set. No significant correlations were found between any socioeconomic factors and public sentiment toward preventive practices.

**Table 5.** Mean (SD) for all socioeconomic and sentiment variables.

| Variable | COVID-19 data set (909 counties), mean (SD) | CDC[a] subset (413 counties), mean (SD) |
|---|---|---|
| Number of tweets | 371.9417 (1052.9472) | 119.2421 (229.7825) |
| Poverty rate | 12.4608 (4.6809) | 11.8521 (4.5293) |
| Unemployment rate | 3.7809 (1.2576) | 3.6608 (1.2024) |
| Median household income (US $) | 61489.36 (16394.21) | 66516.33 (17888.07) |
| Proportion of negative sentiment | 0.5769 (0.0877) | 0.6118 (0.0886) |
| Proportion of positive sentiment | 0.2922 (0.0823) | 0.2446 (0.0817) |
| Proportion of neutral sentiment | 0.1309 (0.0494) | 0.1436 (0.0571) |

[a]CDC: Centers for Disease Control and Prevention.

**Table 6.** Associations among socioeconomic factors and negative, positive, and neutral polarities.

| Variable | COVID-19 data set | | CDC[a] subset | |
|---|---|---|---|---|
| | Coefficient | P value | Coefficient | P value |
| **Negative** | | | | |
| Poverty rate | 0.0461 | .17 | –0.0261 | .60 |
| Unemployment rate | 0.0982 | .003 | 0.0770 | .12 |
| Household income | –0.1322 | <.001 | –0.0203 | .68 |
| **Positive** | | | | |
| Poverty rate | –0.0836 | .01 | –0.0039 | .94 |
| Unemployment rate | –0.1407 | <.001 | –0.0661 | .18 |
| Household income | 0.1554 | <.001 | 0.0329 | .51 |
| **Neutral** | | | | |
| Poverty rate | 0.0574 | .08 | 0.0460 | .35 |
| Unemployment rate | 0.0599 | .07 | –0.0249 | .61 |
| Household income | –0.0242 | .47 | –0.0155 | .75 |

[a]CDC: Centers for Disease Control and Prevention.

## Discussion

### Principal Findings

We conducted 4 types of analysis to answer the 4 research questions. The time series analysis revealed the 4 stages of change in public sentiment toward COVID-19 and the preventive practices for research question 1. People showed high negativity in the initial stage from late January 2020 to the beginning of March 2020, when the COVID-19 risks were not widely recognized in the United States. Wise et al [13] identified that the first week of the COVID-19 pandemic in the United States was March 11, 2020 to March 16, 2020. The first stage mainly reflected how the US population viewed COVID-19 in other countries. Starting from the week of March 11, 2020, when COVID-19 was identified as affecting the United States, people demonstrated growing awareness of the risks associated with COVID-19 and were more engaged in preventative behaviors [13]. Our findings, based on Twitter data, showed similar patterns in that people started to have fewer negative discussions on COVID-19 issues and showed more positive attitudes toward preventive practices in Stage 2. However, the decreasing trend in negative sentiment did not persist. In Stage 3, the proportion of negative sentiment remained stable and lasted for a month. After that, people started to show increasing negative sentiment to both COVID-19 issues and preventive practices, which was not a good sign at a time when the pandemic was far from over. These findings illustrated several challenges in the communication strategies of public health authorities and in government policy making. The first challenge is how to inform people about the disease and its potential risks as well as to convince people to take actions to prevent the spread of the virus in the initial stage when the risks are not geographically close. The second challenge is that people may change their attitudes toward preventive practices after they have experienced the pandemic and obtained information about the disease. It is important to understand what led to the change

in their attitudes and behaviors and how long it takes for people to adapt to or get tired of the changing behaviors.

We analyzed the dynamics of public sentiment at the state level and presented the results from 4 states—California, Florida, New York, and Texas—which showed similar patterns but differed in the timing of the turning points and sentiment polarity to answer research question 2. Our findings were consistent with some existing studies. For example, Hung et al [31] found that Florida was one of the states that expressed the most negative sentiment in COVID-19–related discussions, and our study showed that Floridians were generally more negative in their discussions on general COVID-19 topics and preventive practices.

For research question 3, our study further revealed that people in rural areas generally have more negative sentiment toward COVID-19 issues and preventive practices suggested by the CDC. Czeisler et al [30] conducted representative panel surveys and found that people in New York City and Los Angeles, which are large urban clusters, had more agreement on the stay-at-home orders, business closures, self-isolation, and wearing facial masks in public than the general US population. These findings could be helpful in guiding public authorities in decision making and policy development, for example, to consider the urban and rural differences in communication strategies and guidance.

Further, median household income, as well as poverty and unemployment rates, were not associated with differences in public sentiment to preventive practices; however, higher unemployment rate was positively correlated with negative polarity to COVID-19–related topics, which addresses research question 4. The finding was different from that in the survey study by Czeisler et al [30], which showed people who were unemployed had more agreement on social distancing, wearing masks, stay-at-home orders, and business closures and were less likely to accept the reopening of the United States. The

differences might be caused by the sampling method used in the survey. Combined with the results of the urban/rural analysis, we suggest that different policies or communication strategies may be considered more from the urban/rural perspective than based on socioeconomic differences in pandemics similar to COVID-19.

## Limitations

This study relied on geolocated Twitter data to estimate sentiment polarity at different levels of temporal and spatial granularity. We used the followers-to-followees ratio to remove accounts that were potentially nonindividual users such as bots, which may not be fully accurate. For future work, we believe a bot detection algorithm incorporating more user information may provide more accurate user filtering. Twitter users who have geolocated posts are profiled to be of the younger generation with higher socioeconomic levels who may not represent the whole population in the United States. Considering that the proportion of such users in the population is similar across counties or states, the comparative directions with sampled Twitter users can be representative. To avoid biased interpretation, our findings focused more on the directions and significant level of relationships rather than how large the differences or the correlation coefficients were. Studies with social media data are valuable as they could provide time-sensitive knowledge at different spatial scales, which are difficult to achieve with survey studies in a cost-effective way. Notably, survey methods are irreplaceable to collect attitudes of people who do not go online.

Another limitation came from the algorithm we used for the detection of sentiment. Although the pretrained deep learning model has state-of-the-art sentiment classification accuracy, it may generate wrong sentiment classifications for posts. When the data are scarce, the errors caused by the detection algorithm may lead to large variances in the aggregated sentiment polarity. That is why we adjusted the temporal granularity in the computation of public sentiment for preventive practices and for states. Given more scarce data in the study of other topics, the choice of aggregation level should be more coarsely grained.

Finally, we focused on the sentiment and classified posts as positive, negative, or neutral. There is a need for a deeper understanding and assessment of Twitter content to accurately characterize reaction in multiple dimensions, such as support, hope, and happy that belong to the positive sentiment and fear, despair, and hate that belong to the negative sentiment [32].

## Comparison With Prior Work

Many researchers have studied online discussions, specifically public sentiment, and popular topics, during COVID-19 for timely situational awareness. For example, Xiang et al [33] examined discussions related to older adults on Twitter between January 23, 2020 and May 20, 2020. They identified the lockdown theme was the most popular one where "fear" and "sadness" were the prevalent sentiments. Wang et al [12] analyzed the topics and associated sentiment of social media posts about COVID-19 in China. There were increasing negative emotions expressed from January 20, 2020. Worries about production activity, such as "go to work" and "resume work,"

started to grow from January 26, 2020. In our study, we focused on topics related to preventative COVID-19 practices on Twitter. Although they have been studied locally with survey methods [13,30], few have systematically investigated the topics through social media analysis.

Studies have been done to analyze public sentiment from the perspectives of temporal variations and spatial disparities in COVID-19. Xi et al [34] used Weibo data to understand concerns of the elderly during COVID-19 in China. They identified 3 temporal stages from January 20, 2020 to April 28, 2020, with "older adults contributing to the community" in the first stage and "older patients in hospital" in the second and third stages. Zhou et al [6] tracked the sentiment dynamics of tweets on COVID-19 in Australia regarding topics such as lockdown and social distancing. The overall sentiment polarity toward these policies changed at different stages. Positive sentiment played a dominant role initially but decreased over time. Li et al [8] analyzed English tweets from March 25, 2020 to April 7, 2020. Their results showed a high variation in sadness, anger, and anticipation in tweets containing the term "mask" and disgust and sadness in tweets containing the term "lockdown." A temporal analysis on COVID-19 tweets from January 2020 to June 2020 in 4 countries showed that negative sentiment increased following the lockdown policy enforced by the government of these countries [35].

Several studies have leveraged the geolocation information in social media data to examine public sentiment in different administrative units. Han et al [36] analyzed microblogs in China and showed that the topic "Government response" was the most prominent in Beijing, Sichuan, and Xi'an, while in the surrounding areas of Wuhan, negative sentiment and the topic "Seeking help" were trending in early 2020. Nilima et al [37] investigated the psychosocial factors associated with COVID-19 and the lockdown in India. They detected a clustering of places with similar reaction patterns and found people in different states have different concerns. Imran et al [38] found people's reactions to COVID-19 were culturally different, as people in Pakistan and India showed different sentiment patterns from people in the United States and Canada. Not many studies have specifically examined the discussions in the United States. Van et al [29] investigated public attitudes to social distancing in the United States and found there were geographical variances, which can be partially explained by political ideologies. Chun et al [4] collected tweets in one week of March about government enforcement for the spreading of COVID-19 and calculated the citizens' concern index for different measures. It showed that school closing–related tweets contained the highest levels of concern. Our findings contribute to the knowledge of public sentiment and public opinions related to COVID-19 on social media platforms in the United States. We conducted a comprehensive study to analyze temporal changes from the initial stage when COVID-19 was yet to spread in the United States to the stage when people started to show rebound in negative sentiment or resistance to preventive practices.

Additionally, we explored the association between public sentiment polarity and other geographical and socioeconomic factors to identify factors that were related to the spatial disparities. The findings could be helpful to guide public health

authorities in decision making and policy development in a similar pandemic in the future.

In this study, we applied a data analysis framework to investigate public sentiment toward COVID-19 and the preventive practices suggested by public health authorities in the United States. The data processing framework can be applied to the analysis of discussions on other topics such as vaccination and reopening evaluation in COVID-19 or provide useful solutions for future crises.

## Conclusions

This study used a data-driven method to understand public sentiment to the COVID-19 issues and preventive practices with geolocated Twitter data. We first used a deep learning model to acquire the sentiment of each tweet. These tweets were then aggregated into different temporal and geographical units to measure the polarity of public sentiment.

In the temporal analysis, we discovered 4 stages of change that were evident in discussions on both COVID-19 issues and preventive practices, demonstrating a common pattern between the 2 topics. Based on the examination of our sample of 4 states with the largest volume of tweets across the time period studied, Florida had more negative sentiment to COVID 19 issues and CDC preventive practices than California, Texas, and New York. We analyzed the spatial disparities and explored whether the variations in public sentiment were associated with geographical factors and discovered that there were significant differences in sentiment polarity to preventive practices between urban and rural areas. Socioeconomic factors such as median household income, as well as poverty and unemployment rates, were significantly related to sentiment polarity to COVID-19 issues but not to preventive practices.

The insight gained from the study could be helpful for public health authorities and governments to adjust and differentiate the communication strategies and policies throughout the stages of a pandemic. Communication strategies and policies should be considered based on urban and rural differences more than socioeconomic differences.

## Conflicts of Interest

None declared.

## References

1. WHO coronavirus disease (covid-19) dashboard. World Health Organization. URL: https://covid19.who.int/ [accessed 2021-12-14]
2. Guidance documents. Centers for Disease Control and Prevention (CDC). URL: https://www.cdc.gov/coronavirus/2019-ncov/communication/guidance-list.html [accessed 2021-12-14]
3. Hong L, Fu C, Wu J, Frias-Martinez V. Information needs and communication gaps between citizens and local governments online during natural disasters. Inf Syst Front 2018 Mar 3;20(5):1027-1039. [doi: 10.1007/s10796-018-9832-0]
4. Chun SA, Li ACY, Toliyat A, Geller J. Tracking citizen's concerns during covid-19 pandemic. 2020 Presented at: 21st Annual International Conference on Digital Government Research; June 15-19, 2020; Seoul, Republic of Korea. [doi: 10.1145/3396956.3397000]
5. Sesagiri Raamkumar A, Tan SG, Wee HL. Measuring the outreach efforts of public health authorities and the public response on Facebook during the COVID-19 pandemic in early 2020: cross-country comparison. J Med Internet Res 2020 May 19;22(5):e19334 [FREE Full text] [doi: 10.2196/19334] [Medline: 32401219]
6. Zhou J, Yang S, Xiao C, Chen F. Examination of Community Sentiment Dynamics due to COVID-19 Pandemic: A Case Study from A State in Australia. Cornell University. 2021. URL: https://arxiv.org/abs/2006.12185 [accessed 2021-12-14]
7. Samuel J, Rahman MM, Ali GGMN, Samuel Y, Pelaez A, Chong PHJ, et al. Feeling positive about reopening? New normal scenarios rrom COVID-19 US reopen sentiment analytics. IEEE Access 2020;8:142173-142190. [doi: 10.1109/access.2020.3013933]
8. Li I, Li Y, Li T, Alvarez-Napagao S, Garcia-Gasulla D, Suzumura T. What Are We Depressed About When We Talk About COVID-19: Mental Health Analysis on Tweets Using Natural Language Processing. In: Bramer M, Ellis R, editors. Artificial Intelligence XXXVII. Cham, Switzerland: Springer International Publishing; 2020.
9. Samuel J, Ali GGMN, Rahman MM, Esawi E, Samuel Y. COVID-19 public sentiment insights and machine learning for Tweets classification. Information 2020 Jun 11;11(6):314. [doi: 10.3390/info11060314]
10. Xue J, Chen J, Hu R, Chen C, Zheng C, Su Y, et al. Twitter discussions and emotions about the COVID-19 pandemic: machine learning approach. J Med Internet Res 2020 Nov 25;22(11):e20550 [FREE Full text] [doi: 10.2196/20550] [Medline: 33119535]
11. Hong L, Fu C, Torrens P, Frias-Martinez V. Understanding Citizens' and Local Governments' Digital Communications During Natural Disasters: The Case of Snowstorms. 2017 Presented at: ACM on Web Science Conference; June 25-28, 2017; Troy, NY. [doi: 10.1145/3091478.3091502]
12. Wang T, Lu K, Chow KP, Zhu Q. COVID-19 sensing: negative sentiment analysis on social media in China via BERT model. IEEE Access 2020;8:138162-138169. [doi: 10.1109/access.2020.3012595]

XSL•FO

**RenderX**

13. Wise T, Zbozinek T, Michelini G, Hagan CC, Mobbs D. Changes in risk perception and self-reported protective behaviour during the first week of the COVID-19 pandemic in the United States. R Soc Open Sci 2020 Sep;7(9):200742 [FREE Full text] [doi: 10.1098/rsos.200742] [Medline: 33047037]

14. Ntompras C, Drosatos G, Kaldoudi E. A high-resolution temporal and geospatial content analysis of Twitter posts related to the COVID-19 pandemic. J Comput Soc Sc 2021 Oct 20:1-43. [doi: 10.1007/s42001-021-00150-8]

15. Cuomo RE, Purushothaman V, Li J, Cai M, Mackey TK. Sub-national longitudinal and geospatial analysis of COVID-19 tweets. PLoS One 2020 Oct 28;15(10):e0241330 [FREE Full text] [doi: 10.1371/journal.pone.0241330] [Medline: 33112922]

16. Cuomo RE, Purushothaman V, Li J, Cai M, Mackey TK. A longitudinal and geospatial analysis of COVID-19 tweets during the early outbreak period in the United States. BMC Public Health 2021 Apr 24;21(1):793 [FREE Full text] [doi: 10.1186/s12889-021-10827-4] [Medline: 33894745]

17. Forati AM, Ghose R. Geospatial analysis of misinformation in COVID-19 related tweets. Appl Geogr 2021 Aug;133:102473 [FREE Full text] [doi: 10.1016/j.apgeog.2021.102473] [Medline: 34103772]

18. Hou X, Gao S, Li Q, Kang Y, Chen N, Chen K, et al. Intracounty modeling of COVID-19 infection with human mobility: Assessing spatial heterogeneity with business traffic, age, and race. Proc Natl Acad Sci U S A 2021 Jun 15;118(24):e2020524118 [FREE Full text] [doi: 10.1073/pnas.2020524118] [Medline: 34049993]

19. Schmelz K. Enforcement may crowd out voluntary support for COVID-19 policies, especially where trust in government is weak and in a liberal society. Proc Natl Acad Sci U S A 2021 Jan 05;118(1):e2016385118 [FREE Full text] [doi: 10.1073/pnas.2016385118] [Medline: 33443149]

20. Chen E, Lerman K, Ferrara E. Tracking social media discourse about the COVID-19 pandemic: development of a public coronavirus Twitter data set. JMIR Public Health Surveill 2020 May 29;6(2):e19273 [FREE Full text] [doi: 10.2196/19273] [Medline: 32427106]

21. Cartographic Boundary Files - Shapefile. United States Census Bureau. 2018. URL: https://www.census.gov/geographies/mapping-files/time-series/geo/carto-boundary-file.html [accessed 2021-12-14]

22. Explore Census Data. United States Census Bureau. URL: https://data.census.gov/cedsci/ [accessed 2021-12-14]

23. TIGER/Line® Shapefiles. United States Census Bureau. URL: https://www.census.gov/cgi-bin/geo/shapefiles/ [accessed 2021-12-14]

24. Urban and Rural. United States Census Bureau. URL: https://www.census.gov/programs-surveys/geography/guidance/geo-areas/urban-rural.html [accessed 2021-12-14]

25. Hong L, Frias-Martinez V. Modeling and predicting evacuation flows during hurricane Irma. EPJ Data Sci 2020 Sep 29;9(1):1. [doi: 10.1140/epjds/s13688-020-00247-6]

26. Akbik A, Blythe D, Vollgraf R. Contextual string embeddings for sequence labeling. Proceedings of the 27th International Conference on Computational Linguistics 2018:1638-1649 [FREE Full text]

27. Akbik A, Bergmann T, Blythe D, Rasul K, Schweter S, Vollgraf R. FLAIR: an easy-to-use framework for state-of-the-art NLP. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations) 2019:54-59 [FREE Full text] [doi: 10.18653/v1/n19-4010]

28. Xie Z, Jayanth A, Yadav K, Ye G, Hong L. Multi-faceted Classification for the Identification of Informative Communications during Crises: Case of COVID-19. 2021 Presented at: 45th Annual Computers, Software, and Applications Conference (COMPSAC); July 12-16, 2021; Madrid, Spain. [doi: 10.1109/compsac51774.2021.00125]

29. Van Loon A, Stewart S, Waldon B, Lakshmikanth SK, Shah I, Guntuku SC, et al. Explaining the Trump Gap in Social Distancing Using COVID Discourse. Proceedings of the 1st Workshop on NLP for COVID-19 (Part 2) at EMNLP 2020 2020:1 [FREE Full text] [doi: 10.18653/v1/2020.nlpcovid19-2.10]

30. Czeisler ME, Tynan MA, Howard ME, Honeycutt S, Fulmer EB, Kidder DP, et al. Public attitudes, behaviors, and beliefs related to COVID-19, stay-at-home orders, nonessential business closures, and public health guidance - United States, New York City, and Los Angeles, May 5-12, 2020. MMWR Morb Mortal Wkly Rep 2020 Jun 19;69(24):751-758 [FREE Full text] [doi: 10.15585/mmwr.mm6924e1] [Medline: 32555138]

31. Hung M, Lauren E, Hon ES, Birmingham WC, Xu J, Su S, et al. Social network analysis of COVID-19 sentiments: application of artificial intelligence. J Med Internet Res 2020 Aug 18;22(8):e22590 [FREE Full text] [doi: 10.2196/22590] [Medline: 32750001]

32. Gaspar R, Pedro C, Panagiotopoulos P, Seibt B. Beyond positive or negative: Qualitative sentiment analysis of social media reactions to unexpected stressful events. Computers in Human Behavior 2016 Mar;56:179-191. [doi: 10.1016/j.chb.2015.11.040]

33. Xiang X, Lu X, Halavanau A, Xue J, Sun Y, Lai PHL, et al. Modern senicide in the face of a pandemic: an examination of public discourse and sentiment about older adults and COVID-19 using machine learning. J Gerontol B Psychol Sci Soc Sci 2021 Mar 14;76(4):e190-e200 [FREE Full text] [doi: 10.1093/geronb/gbaa128] [Medline: 32785620]

34. Xi W, Xu W, Zhang X, Ayalon L. A thematic analysis of Weibo topics (Chinese Twitter hashtags) regarding older adults during the COVID-19 outbreak. J Gerontol B Psychol Sci Soc Sci 2021 Aug 13;76(7):e306-e312 [FREE Full text] [doi: 10.1093/geronb/gbaa148] [Medline: 32882029]

35. Mansoor M, Gurumurthy K, Prasad VRB. Global sentiment analysis of covid-19 tweets over time. Cornell University. 2020. URL: https://arxiv.org/abs/2010.14234 [accessed 2021-12-14]

36.  Han X, Wang J, Zhang M, Wang X. Using social media to mine and analyze public opinion related to COVID-19 in China. Int J Environ Res Public Health 2020 Apr 17;17(8):2788 [FREE Full text] [doi: 10.3390/ijerph17082788] [Medline: 32316647]

37.  Nilima N, Kaushik S, Tiwary B, Pandey PK. Psycho-social factors associated with the nationwide lockdown in India during COVID- 19 pandemic. Clin Epidemiol Glob Health 2021 Jan;9:47-52 [FREE Full text] [doi: 10.1016/j.cegh.2020.06.010] [Medline: 32838060]

38.  Imran AS, Daudpota SM, Kastrati Z, Batra R. Cross-cultural polarity and emotion detection using sentiment analysis and deep learning on COVID-19 related tweets. IEEE Access 2020;8:181074-181090. [doi: 10.1109/access.2020.3027350]

## Abbreviations

**API:** application programming interface
**CDC:** Centers for Disease Control and Prevention
**PPE:** personal protective equipment
**TFIDF:** term frequency-inverse document frequency

<u>Original Paper</u>

# Desensitization to Fear-Inducing COVID-19 Health News on Twitter: Observational Study

Hannah R Stevens[1], BA; Yoo Jung Oh[1], MA; Laramie D Taylor[1], PhD

Department of Communication, University of California, Davis, Davis, CA, United States

**Corresponding Author:**
Hannah R Stevens, BA
Department of Communication
University of California, Davis
1 Shields Ave
Davis, CA, 95616
United States
Phone: 1 530 752 1011
Email: hrstevens@ucdavis.edu

**Related Article:**

This is a corrected version. See correction statement: https://infodemiology.jmir.org/2021/1/e32231

## Abstract

**Background:** As of May 9, 2021, the United States had 32.7 million confirmed cases of COVID-19 (20.7% of confirmed cases worldwide) and 580,000 deaths (17.7% of deaths worldwide). Early on in the pandemic, widespread social, financial, and mental insecurities led to extreme and irrational coping behaviors, such as panic buying. However, despite the consistent spread of COVID-19 transmission, the public began to violate public safety measures as the pandemic got worse.

**Objective:** In this work, we examine the effect of fear-inducing news articles on people's expression of anxiety on Twitter. Additionally, we investigate desensitization to fear-inducing health news over time, despite the steadily rising COVID-19 death toll.

**Methods:** This study examined the anxiety levels in news articles (n=1465) and corresponding user tweets containing "COVID," "COVID-19," "pandemic," and "coronavirus" over 11 months, then correlated that information with the death toll of COVID-19 in the United States.

**Results:** Overall, tweets that shared links to anxious articles were more likely to be anxious (odds ratio [OR] 2.65, 95% CI 1.58-4.43, $P<.001$). These odds decreased (OR 0.41, 95% CI 0.2-0.83, $P=.01$) when the death toll reached the third quartile and fourth quartile (OR 0.42, 95% CI 0.21-0.85, $P=.01$). However, user tweet anxiety rose rapidly with articles when the death toll was low and then decreased in the third quartile of deaths (OR 0.61, 95% CI 0.37-1.01, $P=.06$). As predicted, in addition to the increasing death toll being matched by a lower level of article anxiety, the extent to which article anxiety elicited user tweet anxiety decreased when the death count reached the second quartile.

**Conclusions:** The level of anxiety in users' tweets increased sharply in response to article anxiety early on in the COVID-19 pandemic, but as the casualty count climbed, news articles seemingly lost their ability to elicit anxiety among readers. Desensitization offers an explanation for why the increased threat is not eliciting widespread behavioral compliance with guidance from public health officials. This work investigated how individuals' emotional reactions to news of the COVID-19 pandemic manifest as the death toll increases. Findings suggest individuals became desensitized to the increased COVID-19 threat and their emotional responses were blunted over time.

XSL·FO

**RenderX**

## Introduction

### Background

The COVID-19 outbreak has spread worldwide, affecting most countries. Since the outbreak of COVID-19, the number of confirmed cases and the death toll have steadily risen. According to Johns Hopkins University, as of May 9, 2021, more than 157.9 million cases of COVID-19 and 3.2 million deaths have been reported worldwide [1]. Among the countries affected by COVID-19, the United States has had 32.7 million cases (23.5% of confirmed cases worldwide) and 580,000 deaths (17.7% of deaths worldwide). The overabundance of information, misinformation, and disinformation surrounding COVID-19 on social media in the United States has fueled a COVID-19 infodemic, which has jeopardized public health policy aimed at mitigating the pandemic [2], raising questions about the cognitive processes underlying public responses to COVID-19 health information.

Extreme safety precautions (eg, statewide lockdowns, travel bans) have impacted individuals' physical and mental health in the United States. People experienced intense psychological frustration and anxiety regarding the virus and strict safety measures (eg, stay-at-home measures), especially during the early stages of the COVID-19 pandemic [3-5]. Social, financial, and mental insecurities have even led to extreme and irrational coping behaviors, such as panic buying from January to March 2020 [5]. However, throughout the pandemic, the public became desensitized to reports of COVID-19's health threat, and the rising number of confirmed cases and death toll began to lose impact [6,7]. As a result, segments of the public began violating public safety measures as the pandemic progressed, despite the consistent spread of COVID-19 [8-10].

From such observations, two key considerations arise. First, fear-eliciting health messages have a significant effect on eliciting motivation to take action to control the threat. However, repeated exposure to these messages over long periods results in desensitization to those stimuli. In this work, we examine the effect of fear-inducing news articles on people's expression of anxiety on Twitter. Additionally, we investigate how people are desensitized by fear-inducing news articles over time, despite the steadily rising COVID-19 death toll.

### Effect of Fear-Inducing Messages on Public Anxiety

The current pandemic has fueled rapidly evolving news cycles and shaped public sentiment [11,12]. Public health experts' recommendations to mitigate the COVID-19 threat, including widespread business shutdowns and physical distancing guidelines, have proven psychologically and emotionally taxing [13], inducing intense psychological frustration and anxiety among the public [3-5]. Previous literature suggests that fear-inducing messages influence emotions and behaviors when individuals perceive the message to be relevant (ie, they feel susceptible to the threat) and serious (ie, the threat is severe). That is, the heightened threat induces fear and anxiety that, in turn, motivate people to take action [14-17].

In the context of the COVID-19 pandemic, the efficacy of fear-inducing messages on behavioral compliance with public health officials is consequential. Reports of increased COVID-19 transmission and the rising death toll may elicit anxiety about the virus, consequently motivating behaviors intended to manage the problem. For instance, a national survey examining the mental health consequences of COVID-19 fear among US adults during March 2020 found that respondents generally expressed moderate to high COVID-19 fear and anxiety (7 on a scale of 10), and increased anxiety was most prevalent in areas with the highest reported COVID-19 cases [3]. Subsequently, the fear and anxiety induced by COVID-19–related threats can lead people to seek more health-related protective strategies. For example, one study found that as the threat of COVID-19 increased, people expressed more fear-related emotions and they were subsequently increasingly motivated to search for preventative behaviors and information online [5].

### Desensitization to Fear-Inducing Messages

Although fear-based health messages have been shown to motivate changes in behavior, repeated exposure to even highly arousing stimuli—such as news of the rising death toll from COVID-19—may eventually result in desensitization to those stimuli [6,18]. Desensitization refers to the process by which cognitive, emotional, and physiological responses to a stimulus are reduced or eliminated over protracted or repeated exposure [19]. It can play an important adaptive role in allowing individuals to function in difficult circumstances that might otherwise result in overwhelming and persistent anxiety or fear. For example, one analysis of Twitter messages from a region of Mexico with then-rising violence found the expressions of negative emotions declined [20]. Although increasing anxiety and fear might prompt security-seeking behavior, these emotions may also be paralyzing; some measure of desensitization can facilitate continuing with necessary everyday tasks.

Numerous studies have demonstrated desensitization to media content. Research has often focused on fictional depictions of violence [21,22]; however, desensitization has also been demonstrated in response to repeated exposure to violent news stories [23], hate speech [24], and sexually explicit internet content [25], although this last finding has mixed support [26].

Researchers studying social media data have explored the possibility that news messages can result in desensitization. Li and colleagues [27] analyzed a large sample of Twitter data, examining posts linked to guns and shootings for emotional language. They observed that across 3 years of mass shootings and school shootings in the United States, the frequency of negative emotional words used in shooting-related tweets declined; they argued that this reflected desensitization to gun violence.

In the context of the COVID-19 pandemic, news audiences have been repeatedly exposed to highly arousing messages related to COVID-19–related deaths—messages that inherently communicate explicit and implicit threats of serious illness and death to readers. Fundamentally, the biological response to threat communicated through text is similar to threats communicated in other ways [28]. Over time, as the death toll has increased, the cognitive, emotional, and physiological responses to threatening COVID-19 news may have been blunted. Individuals may have become desensitized to

XSL•FO

**RenderX**

threatening COVID-19 information and experienced diminished anxiety over time, even in the face of an increasing threat.

## Rationale and Aims

The public relies heavily on news disseminated through social media for information about the spread of the virus [29]. Twitter, in particular, is a popular outlet for sharing news [30] and has become a forum for individuals to communicate their feelings about COVID-19 [11]. Social media text analysis has emerged as a particularly effective way to assess sentiment dynamics surrounding public health crises; consider, for example, the Zika outbreak [31]. This study uses social media text analysis to examine the anxiety levels in news articles and related tweets over 11 months, then considers those levels in the context of deaths from COVID-19 on the day the post was shared [32].

The general hypothesis guiding this research is that audiences will have become desensitized to COVID-19 deaths over the course of the pandemic, decreasing the level of anxiety elicited by fearful COVID-19 health information reported in the news. To the best of our knowledge, ours is the first study to investigate whether, as the objective threat and harm of COVID-19 has increased, individuals have become desensitized to news reports of cautionary COVID-19 health information.

# Methods

## Overview

This study examined how anxiety levels in news articles predicted users' tweet anxiety levels over 11 months, then correlated that information with the total death toll of COVID-19 in the United States as reported to the Centers for Disease Control and Prevention (CDC) on the day the post was shared [32]. Employing semantic analysis procedures to analyze anxiety in the full news articles and their corresponding user tweets allowed us to examine how fear elicited by COVID-19 health news manifests as individuals become desensitized to news of COVID-19–related deaths.
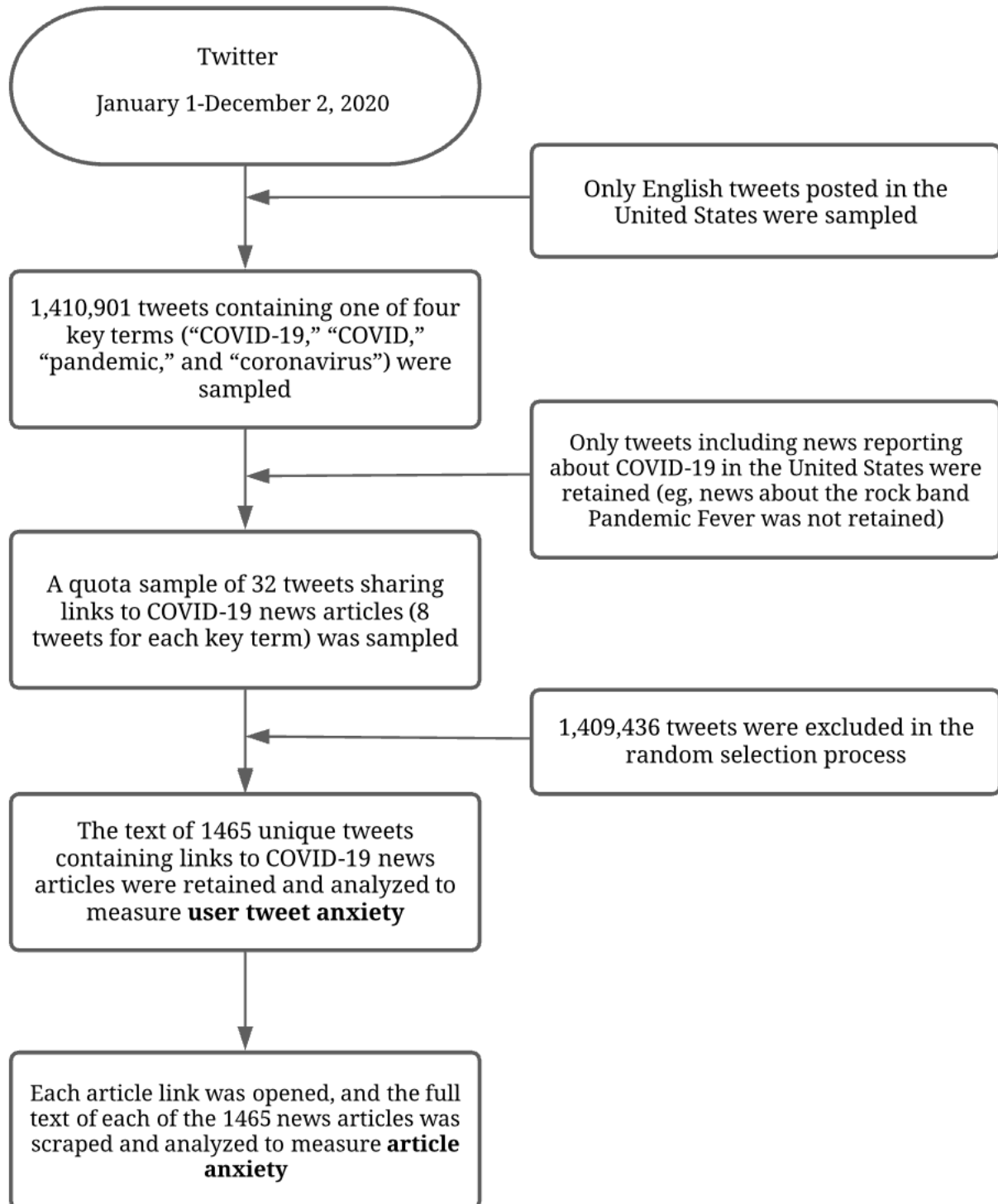
## Data Collection

The sample comprises content shared to Twitter, a popular social media platform used for sharing news [30]. The text of 1465 news articles and corresponding posts by users were collected from tweets containing the terms "COVID-19," "COVID," "pandemic," and "coronavirus" from January 1 to December 2, 2020. For an overview of the data collection process, see Figure 1.

The Python programming language was used to extract posts sharing news reports of COVID-19 health information. We collected a quota sample of 32,000 US tweets containing one of four key terms (ie, COVID, COVID-19, coronavirus, pandemic) each week from January 1 to December 2, 2020. The GetOldTweets3 Python3 library was used to scrape tweets for the months of January-July 2020 [33]. Twitter's application programming interface (version 2) was used to collect tweets from August-December 2020 [34].

Human coders then filtered through the sample of 1,410,901 tweets to randomly extract a quota of 8 original tweets per key term from each week sharing a news report about COVID-19. Data collection resulted in thousands of tweets containing links per week. To facilitate the representativeness of the news articles, 32 tweets were drawn from each week from a shuffled list of tweets containing hyperlinks. Since we aimed to assess users' reactions to the text of the article they read, without the confounding textual framing of other peoples' commentary about an article, retweets were excluded from the analysis. If a quota of 32 tweets each week (8 per key term) was not met, additional tweets were sampled for that week. Notably, the disease and pandemic were not commonly referred to as COVID-19 in early January; accordingly, three weeks did not have 8 tweets with the terms "COVID" and "COVID-19" per week.

The news articles were collected from links shared by Twitter users in general, regardless of who posted the tweet. We only included users sharing links to news articles regarding COVID-19 in the United States; all other content was excluded (eg, news about the rock band Pandemic Fever). If all posts for that week were excluded, another sample from that week was drawn. If a tweet linked to a news article that had been taken down, a replacement post was sampled from the same week. We then extracted the text from the news articles and their corresponding tweets. The final sample was comprised of n=1465 news-sharing tweets.

**Figure 1.** Flowchart of data collection process.



## Linguistic Inquiry and Word Count Sentiment Analysis

Once the final sample was collected (n=1465), we analyzed articles and tweets using the Linguistic Inquiry and Word Count (LIWC) program [35]. The body text of the news articles was analyzed to measure article anxiety, while the tweet text was analyzed to measure tweet anxiety. LIWC is a natural language processing text analysis program that classifies texts by counting the percentage of words in a given text that fall into prespecified categories, such as a linguistic category (eg, prepositions) or psychological processes (eg, anxiety, sadness). In this study, we focused on the percentage of LIWC anxiety lexicon words in news articles and tweets because this psychological process is germane to the efficacy of fear-based news messaging [14,36]. LIWC calculates the percentage of anxiety words relative to all words contained in a text to account for long versus short text classification. For example, we might discover that 15/745 (2.04%) words in a given article were anxiety lexicon words. The LIWC output would then assign that particular article an anxiety score of 2.04 (see Figure 2 for an example).

**Figure 2.** Sample text from a COVID-19 news article shared to Twitter [37]. The words highlighted in red are LIWC anxiety words. Since this article contains 15 anxiety words out of 745 words total (2.4%), this article is assigned a LIWC anxiety score of 2.4. LIWC: Linguistic Inquiry and Word Count.

## How Global Corruption Threatens the U.S. Pandemic Response

Abigail Bellows
Abigail Bellows is a nonresident scholar in the Democracy, Conflict, and Governance Program at the Carnegie Endowment for International Peace
@ABIGAIL_BELLOWS

As Ebola began to rage across the Democratic Republic of the Congo in 2018, the disease had a powerful accomplice: corruption. The country's health minister and his financial adviser embezzled $400,000 in relief funds—crimes for which they were recently sentenced to five years of forced labor. Yet the systemic **vulnerabilities** that enable this type of fraud persist around the world. How is the U.S. government assisting partners in Africa, Latin America, and South Asia in addressing corruption **vulnerabilities** before they are hit by full-blown outbreaks of the new coronavirus?

### IGNORING CORRUPTION IS RISKY

Corruption is like gasoline poured on the flames of a pandemic. Healthcare systems already debilitated by graft will **struggle** to address the most basic of needs during the crisis. Citizens who can't afford to pay bribes may be locked out of access to testing and treatment, a problem that would accelerate the virus's spread. Those who can bribe their way out of quarantines will probably do so, as has been reported already in Cameroon and Uganda. And government attempts to convey public health messages are likely to fall flat in places where decades of corruption have deeply undermined trust in the state.

At the elite level, the pandemic is setting off a flurry of public procurement spending, which faces serious **risks** for diversion, especially since traditional watchdog groups are also scrambling to adapt. Foreign assistance pouring in from the United States and other countries is also **vulnerable** to leakage. In normal times, various sources estimate that more than 10 percent of global healthcare spending is siphoned off by corruption, amounting to losses of more than $500 billion annually—and these **risks** are only heightened during a disaster. Meanwhile, oligarchs may be using the proceeds of corruption to buy up ventilators and arrange for private healthcare, as seen in Russia, a practice that drains resources from the public health system.

The pandemic's further spread around the globe, fueled by corruption, could cause serious harm to U.S. interests and foreign policy objectives. Public anger at government malfeasance could topple regimes, embolden antiestablishment populists, and provide openings for **terrorist** recruitment. Long-standing allies may turn away from the United States and toward China, **desperately** applying authoritarian measures in hopes of containing the virus. If corruption becomes more entrenched overseas, U.S. businesses will **struggle** to compete.

### CHARTING A DIFFERENT COURSE

The good news is that the U.S. government can take action to **avoid** this dark prognosis. The biggest near-term step would be inclusion of the Countering Russian and Other Overseas Kleptocracy (CROOK) Act (H.R. 3843/S. 3026) in Congress's planned fourth coronavirus-related spending package. The bill would form an Anti-Corruption Action Fund to surge support to countries eager to take rapid action against corruption, as the current crisis demands. The proposal, which is budget-neutral and enjoys bipartisan support, has already passed the House Foreign Affairs Committee and has been introduced in the Senate.

The next step would be for relevant committee leaders and the bill's bipartisan sponsors—Representatives Bill Keating and Brian Fitzpatrick alongside Senators Roger Wicker and Ben Cardin—to fold the measure into upcoming legislation. This would fill a glaring gap in congressional action to date on the coronavirus and signal that U.S. decisionmakers recognize the links between international corruption, public health, and U.S. national security.

Alongside congressional leadership, the State Department and the United States Agency for International Development (USAID) could take important steps to address the corruption-coronavirus nexus. Global health assistance should include strong anticorruption safeguards and support for emergency procurement mechanisms that are rigorous and transparent—in both U.S. and multilateral assistance. Diplomats should reinforce the need to maintain anticorruption law enforcement to deter crime during the pandemic. The United States can also celebrate local officials who act with integrity and urge civil society, media outlets, and whistleblowers to keep playing their vital roles in spite of rising **repression**.

The virus has yet to become a full-scale disaster in the most corruption-prone parts of the world—but time is running short. If the United States seeks to **avoid** a replay of the Ebola epidemic—on a far graver scale—it must act now to address global corruption **risks**.

## Statistical Analysis

We paired the final sample with the CDC's aggregate death toll on the day the tweet was posted. Contextualizing the articles and tweets allowed us to examine how fear elicited by COVID-19 health news manifests as individuals become desensitized to news of COVID-19–related deaths.

The outcome of interest was tweet anxiety. Note that the distribution of count data outcome variables (in our case, LIWC tweet anxiety) often contains excess zeros; this result is known as zero inflation. The positive values are skewed, and a considerable "clumping at zero" is trailed by a bump representing positive values [38]. In our specific distribution, the "clumping at zero" represents texts containing zero anxiety
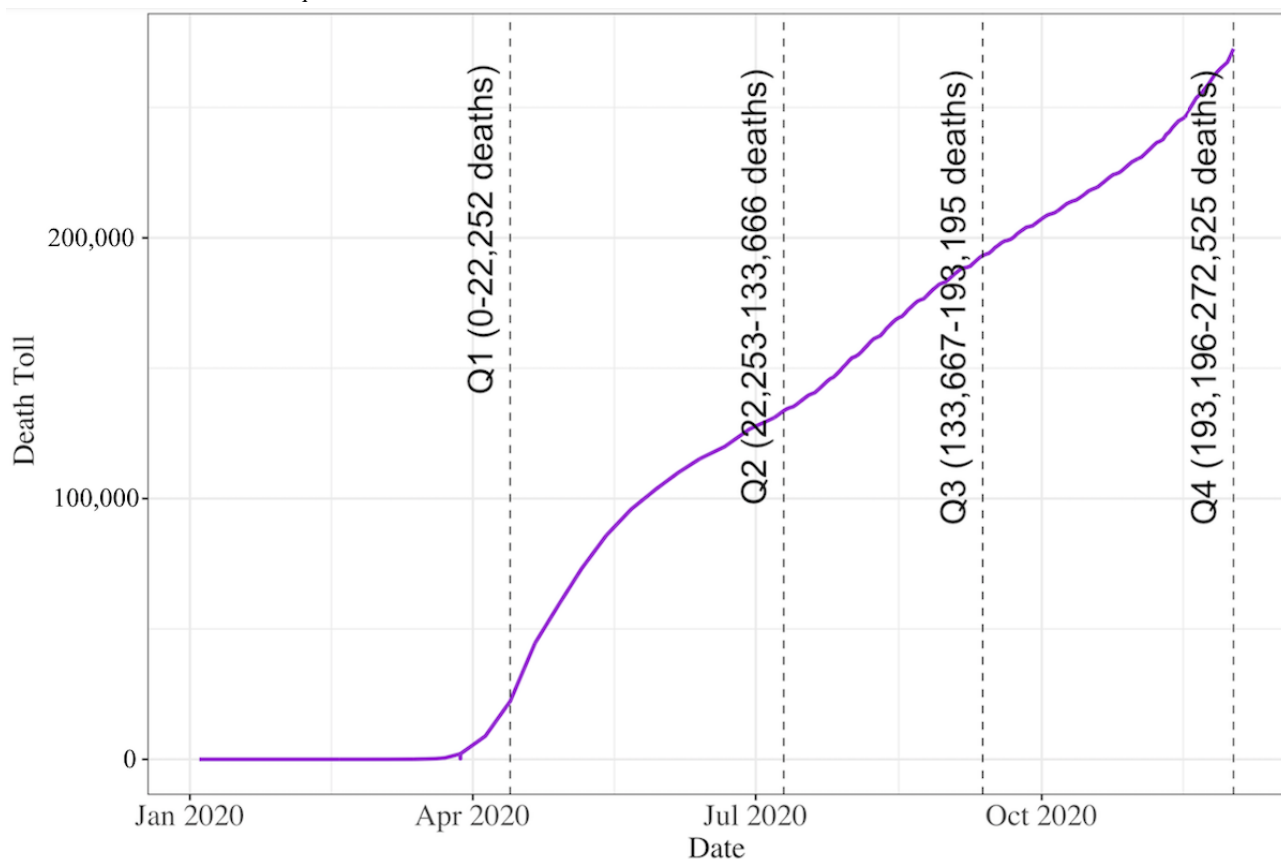
lexicon terms. Generalized linear models are not appropriate for zero-inflation data. As all observed zeros are unambiguous, they are best analyzed separately from the nonzeros.

Two distinct distributions generally characterize zero-inflation data; thus, a zero-inflated model, which separates the zero and nonzero counts, is appropriate [39,40]. In zero-inflated models, the distribution of positive count values depends on the probability of exceeding the hurdle and reaching the distribution of positive values. In other words, it considers the odds of having any anxiety in a tweet versus none at all. For tweets that clear the hurdle, it then considers how much anxiety will be in a tweet on a continuous distribution.

We employed a zero-inflated model using a gamma distribution with a log link to examine any association between article anxiety and death toll, along with their interaction with subsequent tweet anxiety for all values of tweet anxiety greater than zero. We paired that with a model that used a binomial distribution with a logit link to determine zero anxiety versus nonzero anxiety in tweets. We recoded the death toll into categories reflecting the death count at the second quartile, the third quartile, and the fourth quartile relative to the first quartile of the total death count (see Figure 3 for a breakdown). This was necessitated by the skewed and logarithmic character of the distribution. These values were then used in place of the continuous variable to model the interaction. We used R statistical software for data analysis (version 3.6.2; The R Foundation for Statistical Computing).

**Figure 3.** Distribution of death toll quartiles over time.



### Ethics Statement

This study only used information available in the public domain. No personally identifiable information was included in this study. Ethical review and approval was not required for this study because the institutional review board recognizes that the analysis of publicly available data does not constitute human subjects research.
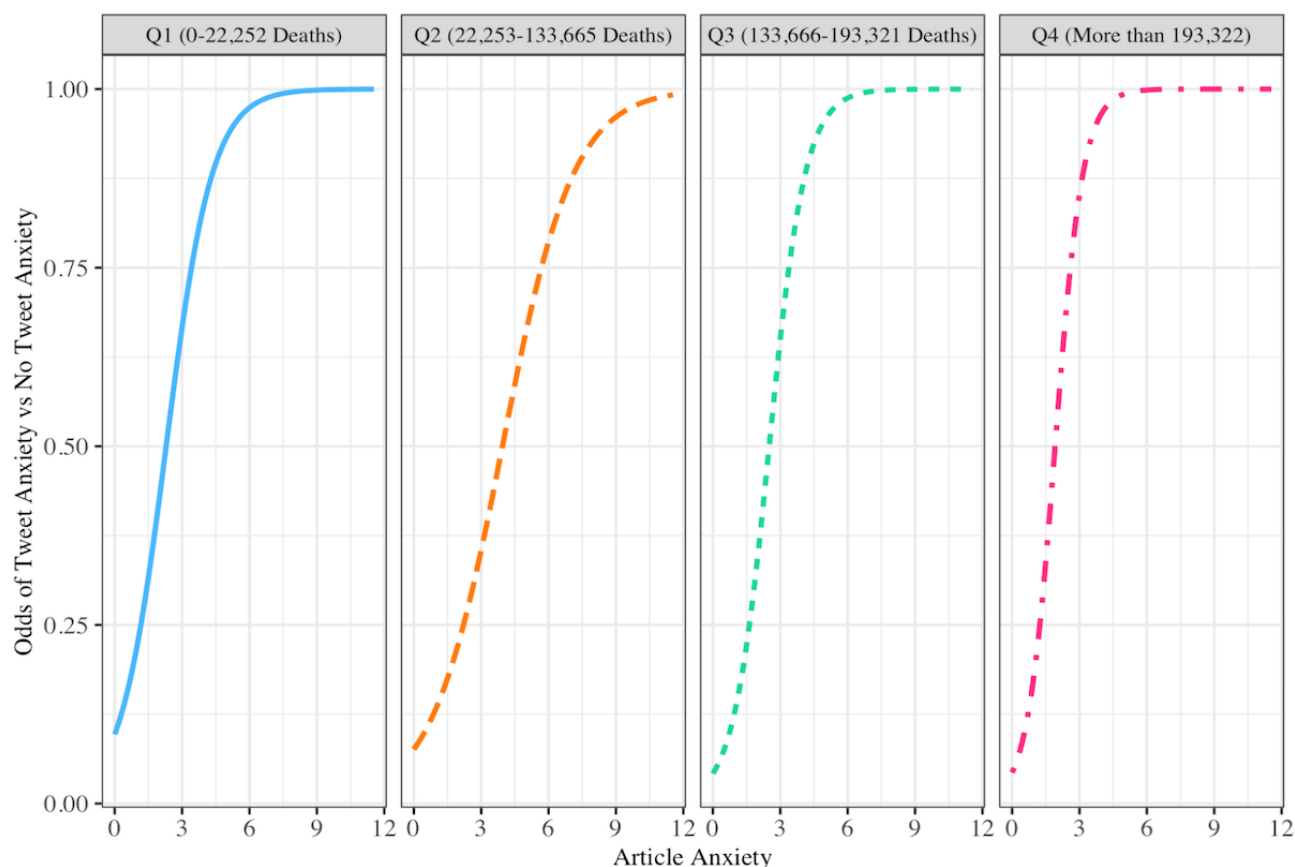
## Results

Results suggest that as the death toll increased over time, the baseline level of anxiety lexicon words in articles decreased; this was evidenced by our finding that when the pandemic's severity and threat increased, individuals shared less news coverage containing COVID-19 anxiety words (eg, "risk," "worried," "threatens"). When assessing the odds of a tweet having no anxiety versus anxiety, we found that the baseline odds of *not* having anxiety in a tweet were 0.11; the odds of having anxiety in a tweet increased (odds ratio [OR] 2.65, 95% CI 1.58-4.43, *P*<.001) with each unit increase in anxiety within an article. The odds of tweet anxiety decreased as paired with CDC total deaths in the third quartile (OR 0.41, 95% CI 0.2-0.83, *P*=.01) and fourth quartile (OR 0.42, 95% CI 0.21-0.85, *P*=.01), respectively (see Table 1 and Figure 4).

**Table 1.** The odds of a tweet containing anxiety language versus no anxiety language, as determined using a zero-inflated model with categorical death[a].

| Variable | Odds ratio (95% CI) | P value |
|---|---|---|
| Intercept | 0.11 (0.07-0.16) | <.001 |
| Anxiety in article | 2.65 (1.58-4.43) | <.001 |
| Second quartile (22,253-133,665 deaths) | 0.76 (0.41-1.41) | .39 |
| Third quartile (133,666-193,321 deaths) | 0.41 (0.2-0.83) | .01 |
| Fourth quartile (≥193,322 deaths) | 0.42 (0.21-0.85) | .02 |
| Interaction anxiety in article by second quartile deaths (22,253-133,665 deaths) | 0.71 (0.34-1.48) | .36 |
| Interaction anxiety in article by third quartile deaths (133,666-193,321 deaths) | 1.32 (0.54-3.24) | .55 |
| Interaction anxiety in article by fourth quartile deaths (≥193,322 deaths) | 1.9 (0.75-4.83) | .18 |

[a]This table reports the odds of no tweet anxiety versus tweet anxiety. Deaths were categorized based on the second, third, and fourth quartiles relative to the first quartile.

**Figure 4.** Article anxiety predicting the odds of tweet anxiety versus no tweet anxiety at the first, second, third, and fourth quartiles of the COVID-19 death toll.



We then examined the actual estimated linguistic anxiety of tweets, looking only at all of the values in a continuous distribution, excluding those values with zero anxiety (ie, the tweet did not contain any anxiety lexicon words). Although not statistically significant at $P<.05$, the results illuminate an emerging yet meaningful trend. The baseline level of anxiety in a tweet was 3.45. The tweet anxiety level trend increased (OR 1.25, 95% CI 0.99-1.59, $P=.068$) with each unit increase of article anxiety. Overall, tweets that shared links to more anxious articles expressed more anxious terms (eg, "avoid,"

"uncertain," "paranoid"). Notably, the interaction between article anxiety and deaths was not found to be a significant predictor of tweet anxiety level. Tweet anxiety rose rapidly with articles when the death toll was low and then decreased in the third quartile of deaths (OR 0.61, 95% CI 0.37-1.01, $P=.06$). As predicted, in addition to the increasing death toll being matched by a lower level of article anxiety, the extent to which article anxiety elicited tweet anxiety decreased when the death count reached the second quartile (see Table 2 and Figure 5).

**Table 2.** Actual anxiety expressed in tweets, as predicted by article anxiety and COVID-19 death toll: gamma regression model with categorical death[a].

| Variable | Coefficient (95% CI) | *P* value |
|---|---|---|
| Intercept | 3.45 (2.77-4.28) | <.001 |
| Anxiety in article | 1.25 (0.99-1.59) | .07 |
| Second quartile (22,253-133,665 deaths) | 1.53 (1.1-2.15) | .01 |
| Third quartile (133,666-193,321 deaths) | 1.47 (0.97-2.22) | .17 |
| Fourth quartile (≥193,322 deaths) | 1.21 (0.83-1.75) | .32 |
| Interaction anxiety in article by second quartile deaths (22,253-133,665 deaths) | 0.78 (0.56-1.08) | .14 |
| Interaction anxiety in article by third quartile deaths (133,666-193,321 deaths) | 0.61 (0.37-1.01) | .06 |
| Interaction anxiety in article by fourth quartile deaths (≥193,322 deaths) | 0.78 (0.5-1.22) | .28 |

[a]Deaths were categorized based on the second, third, and fourth quartiles relative to the first quartile. This table reports the actual estimated anxiety in the tweet, looking only at all of the values in a continuous distribution, excluding those with zero anxiety.

**Figure 5.** Article anxiety predicting nonzero tweet anxiety at the first, second, third, and fourth quartiles of the COVID-19 death toll.



## Discussion

### Principal Findings

This study reports exploratory findings on the effects of fear-inducing news messages during a pandemic. Most importantly, we demonstrated a link between the anxiety expressed in news articles and the odds of anxiety being expressed by those who shared the articles to Twitter. This likely reflects the ability of pandemic-related news messages to elicit a measure of fear in their readers, consonant with public health goals. However, likely as a function of the rising COVID-19 threat over time (as indicated by LIWC news article anxiety) and a low perceived ability to prevent the rapid spread of the virus, anxiety did not increase in response to climbing death tolls over time. Instead, anxiety in tweets increased sharply in response to article anxiety early on in the pandemic, but as the death toll climbed, it flattened out, and news articles seemingly lost their ability to elicit anxiety among readers.

Such findings from this study provide several insights and directions for future research. Our findings reveal that responses to COVID-19 news as well as the rising death toll are increasingly bland. Growing desensitization in the face of threatening pandemic information impedes public health experts' efforts to mitigate the COVID-19 crisis [41]. Therefore, future research should investigate how to "resensitize" the public and motivate them to take active roles in COVID-19–related responses (eg, wearing masks, washing hands, vaccination). Here, literature on behavioral theories may be helpful in implementing effective resensitization tactics. For instance, the transtheoretical model [42,43], which explains behavior change

through stages of change, suggests that to initiate and maintain health behaviors, it is important to have supportive relationships and motivate one another to share successes and experiences related to engaging in certain behaviors. In addition, it is suggested that reinforcement management—such as getting rewards from behavioral engagement—can be effective. In the context of COVID-19, health care providers can apply these tactics (ie, social support, reward) to motivate people to adhere to public health measures such as vaccination.

Second, since extant research shows that both statistics (eg, percentage of deaths) and cognitive dissonance can elicit desensitization [44,45], scholars should investigate the role of additional psychological processes in desensitization to the COVID-19 threat. Third, as self-disclosure varies by platform [46], more work is needed to explore how anxiety manifests on other platforms for discussing COVID-19 news. Finally, our findings suggest that health care practitioners should be prepared for public desensitization to future global pandemic scenarios. More specifically, it would be important to carefully monitor the public's level of desensitization to health news and implement appropriate resensitization strategies based on different stages in the pandemic.

## Limitations

Our findings illuminate desensitization to fear-inducing news messages during the pandemic; however, this study is not without limitations. By focusing on Twitter, we neglected to explore how anxiety manifests on other platforms for sharing news (eg, the comments section of digital news sites). As different platforms have different community norms [46], it is reasonable to expect manifestations of anxiety to vary by platform. Furthermore, Twitter users are younger, more democratic, and wealthier than the general population of Americans [47]. Acknowledging the biases associated with using computational social media data [48], our findings should be interpreted as representing a subset of the US population (ie,

Twitter users), not all US residents. Second, among 1.4 million tweets collected, only a small number of tweets were sampled in this study. Therefore, our study may lack generalizability. Additionally, the LIWC computerized coding tool does not allow for the nuanced coding that could be achieved with human coders. Although we have attempted to minimize this potential bias using a well-validated sentiment analysis procedure, LIWC [35], this study is limited in its use of anxiety in text as a measure of user anxiety.

## Conclusions

This work investigates how individuals' emotional reactions to news of the COVID-19 pandemic manifest as the death toll increases. Individuals become desensitized to an increased health threat and their emotional responses are blunted over time. Our results suggest desensitized public health reactions to threatening COVID-19 news, which could affect the propensity of individuals to adopt recommended health behaviors.

Public health agencies made recommendations to slow the pandemic's spread, including physically distancing from others when appropriate, wearing masks, engaging in frequent handwashing, and disinfecting frequently touched surfaces. The consequences of ignoring these guidelines initially incited widespread fear and anxiety around contracting the virus or having family and friends contract it and fall ill. Social scientists have tried to inform interventions aimed at promoting compliance with public health experts [49]. The results of this study suggest the increased threat conveyed in COVID-19 news has, however, diminished public anxiety, despite an increase in COVID-19–related deaths. Desensitization offers one way to explain why the increased threat is not eliciting widespread compliance with guidance from public health officials. This work sheds light on both the effectiveness and shortcomings of fear-based health messages during the pandemic, as well as the utility of natural language processing to gain an understanding of public responses to emerging health crises.

## Conflicts of Interest

None declared.

## References

1. Johns Hopkins Coronavirus Resource Center. URL: https://coronavirus.jhu.edu [accessed 2020-12-20]
2. World Health Organization. Managing the COVID-19 infodemic: Promoting healthy behaviours and mitigating the harm from misinformation and disinformation. URL: https://www.who.int/news/item/23-09-2020-managing-the-covid-19-infodemic-promoting-healthy-behaviours-and-mitigating-the-harm-from-misinformation-and-disinformation [accessed 2021-05-17]
3. Fitzpatrick KM, Harris C, Drawve G. Fear of COVID-19 and the mental health consequences in America. Psychol Trauma 2020 Aug;12(S1):S17-S21. [doi: 10.1037/tra0000924] [Medline: 32496100]
4. Marroquín B, Vine V, Morgan R. Mental health during the COVID-19 pandemic: Effects of stay-at-home policies, social distancing behavior, and social resources. Psychiatry Res 2020 Nov;293:113419 [FREE Full text] [doi: 10.1016/j.psychres.2020.113419] [Medline: 32861098]
5. Du H, Yang J, King RB, Yang L, Chi P. COVID-19 Increases Online Searches for Emotional and Health-Related Terms. Appl Psychol Health Well Being 2020 Dec;12(4):1039-1053 [FREE Full text] [doi: 10.1111/aphw.12237] [Medline: 33052612]
6. Arradondo B. Psychologists worry public is becoming desensitized to COVID-19 news. FOX 13 Tampa Bay. URL: https://www.fox13news.com/news/psychologists-worry-public-is-becoming-desensitized-to-covid-19-news [accessed 2020-12-31]

7.   Spector N. Why it's hard for us to fathom the COVID-19 death toll. TODAY. 2020. URL: https://www.today.com/health/why-it-s-hard-us-fathom-covid-19-death-toll-t191426 [accessed 2020-12-31]

8.   Young Americans Are Partying Hard and Spreading Covid-19 Quickly. Bloomberg. URL: https://www.bloomberg.com/news/articles/2020-07-01/young-americans-are-partying-hard-and-spreading-the-virus-fast [accessed 2020-12-31]

9.   TikTok stars charged over partying in LA during pandemic. BBC News. 2020 Aug 28. URL: https://www.bbc.com/news/world-us-canada-53954673 [accessed 2020-12-31]

10.  Hines M. 'Epicenter of the epicenter': Young people partying in Miami Beach despite COVID-19 threat. USA Today. 2020 Jul 16. URL: https://www.usatoday.com/story/travel/news/2020/07/16/vacationers-party-florida-even-covid-19-cases-surge/5449754002/ [accessed 2020-12-31]

11.  Chandrasekaran R, Mehta V, Valkunde T, Moustakas E. Topics, Trends, and Sentiments of Tweets About the COVID-19 Pandemic: Temporal Infoveillance Study. J Med Internet Res 2020 Oct 23;22(10):e22624 [FREE Full text] [doi: 10.2196/22624] [Medline: 33006937]

12.  Liu Q, Zheng Z, Zheng J, Chen Q, Liu G, Chen S, et al. Health Communication Through News Media During the Early Stage of the COVID-19 Outbreak in China: Digital Topic Modeling Approach. J Med Internet Res 2020 Apr 28;22(4):e19118 [FREE Full text] [doi: 10.2196/19118] [Medline: 32302966]

13.  Venkatesh A, Edirappuli S. Social distancing in covid-19: what are the mental health implications? BMJ 2020 Apr 06;369:m1379. [doi: 10.1136/bmj.m1379] [Medline: 32253182]

14.  So J. A further extension of the Extended Parallel Process Model (E-EPPM): implications of cognitive appraisal theory of emotion and dispositional coping style. Health Commun 2013;28(1):72-83. [doi: 10.1080/10410236.2012.708633] [Medline: 23330860]

15.  Witte K. Putting the fear back into fear appeals: The extended parallel process model. Communication Monographs 1992 Dec;59(4):329-349. [doi: 10.1080/03637759209376276]

16.  Witte K. The role of threat and efficacy in AIDS prevention. Int Q Community Health Educ 1991 Jan 01;12(3):225-249. [doi: 10.2190/U43P-9QLX-HJ5P-U2J5] [Medline: 20840971]

17.  Witte K. Fear control and danger control: A test of the extended parallel process model (EPPM). Communication Monographs 2009 Jun 02;61(2):113-134. [doi: 10.1080/03637759409376328]

18.  Davies H. Why your brain is desensitized to today's COVID-19 death toll. Medium. 2020. URL: https://medium.com/the-innovation/why-your-brain-is-desensitized-to-todays-covid-19-death-toll-ed79a76c0324 [accessed 2021-02-04]

19.  Funk JB, Baldacci HB, Pasold T, Baumgardner J. Violence exposure in real-life, video games, television, movies, and the internet: is there desensitization? J Adolesc 2004 Feb;27(1):23-39. [doi: 10.1016/j.adolescence.2003.10.005] [Medline: 15013258]

20.  Choudhury M, Monroy-Hernandez A, Mark G. 'Narco' emotions: affect and desensitization in social media during the Mexican drug war. 2014 Presented at: SIGCHI Conference on Human Factors in Computing Systems; April 2014; Toronto, ON. [doi: 10.1145/2556288.2557197]

21.  Cline VB, Croft RG, Courrier S. Desensitization of children to television violence. J Pers Soc Psychol 1973 Sep;27(3):360-365. [doi: 10.1037/h0034945] [Medline: 4741676]

22.  Krahé B, Möller I, Huesmann LR, Kirwil L, Felber J, Berger A. Desensitization to media violence: links with habitual media violence exposure, aggressive cognitions, and aggressive behavior. J Pers Soc Psychol 2011 Apr;100(4):630-646 [FREE Full text] [doi: 10.1037/a0021711] [Medline: 21186935]

23.  Scharrer E. Media Exposure and Sensitivity to Violence in News Reports: Evidence of Desensitization? Journalism & Mass Communication Quarterly 2008 Jun;85(2):291-310. [doi: 10.1177/107769900808500205]

24.  Soral W, Bilewicz M, Winiewski M. Exposure to hate speech increases prejudice through desensitization. Aggress Behav 2018 Mar;44(2):136-146. [doi: 10.1002/ab.21737] [Medline: 29094365]

25.  Daneback K, Ševčíková A, Ježek S. Exposure to online sexual materials in adolescence and desensitization to sexual content. Sexologies 2018 Jul;27(3):e71-e76. [doi: 10.1016/j.sexol.2018.04.001]

26.  Peter J, Valkenburg PM. Adolescents' Use of Sexually Explicit Internet Material and Sexual Uncertainty: The Role of Involvement and Gender. Communication Monographs 2010 Sep;77(3):357-375. [doi: 10.1080/03637751.2010.498791]

27.  Li J, Conathan D, Hughes C. Rethinking emotional desensitization to violence: Methodological and theoretical insights from social media data. In: Proceedings of the 8th International Conference on Social Media & Society. 2017 Presented at: 8th International Conference on Social Media & Society; July 2017; New York, NY. [doi: 10.1145/3097286.3097333]

28.  Isenberg N, Silbersweig D, Engelien A, Emmerich S, Malavade K, Beattie B, et al. Linguistic threat activates the human amygdala. Proc Natl Acad Sci USA 1999 Aug 31;96(18):10456-10459 [FREE Full text] [doi: 10.1073/pnas.96.18.10456] [Medline: 10468630]

29.  Abd-Alrazaq A, Alhuwail D, Househ M, Hamdi M, Shah Z. Top Concerns of Tweeters During the COVID-19 Pandemic: Infoveillance Study. J Med Internet Res 2020 Apr 21;22(4):e19016 [FREE Full text] [doi: 10.2196/19016] [Medline: 32287039]

30.  Kwak H, Lee C, Park H, Moon S. What is Twitter, a social network or a news media? In: Proceedings of the 19th International Conference on the World Wide Web. New York, USA: ACM Press; 2010 Presented at: WWW '10; April 2010; New York, NY. [doi: 10.1145/1772690.1772751]

31. Lwin MO, Lu J, Sheldenkar A, Schulz PJ. Strategic Uses of Facebook in Zika Outbreak Communication: Implications for the Crisis and Emergency Risk Communication Model. Int J Environ Res Public Health 2018 Sep 10;15(9):1 [FREE Full text] [doi: 10.3390/ijerph15091974] [Medline: 30201929]

32. Centers for Disease Control and Prevention. Provisional death counts for Coronavirus disease 2019 (COVID-19). URL: https://www.cdc.gov/nchs/nvss/vsrr/covid19/index.htm [accessed 2020-12-31]

33. GetOldTweets3. PyPI. URL: https://pypi.org/project/GetOldTweets3/ [accessed 2021-05-12]

34. Introducing a new and improved Twitter API. Twitter. URL: https://blog.twitter.com/developer/en_us/topics/tools/2020/introducing_new_twitter_api.html [accessed 2021-05-12]

35. Pennebaker J, Boyd R, Jordan K, Blackburn K. The development and psychometric properties of LIWC2015. The University of Texas at Austin. URL: https://repositories.lib.utexas.edu/bitstream/handle/2152/31333/LIWC2015_LanguageManual.pdf?Sequence=3 [accessed 2020-12-31]

36. Lewis I, Watson B, White KM. Extending the explanatory utility of the EPPM beyond fear-based persuasion. Health Commun 2013;28(1):84-98. [doi: 10.1080/10410236.2013.743430] [Medline: 23330861]

37. Bellows A. How global corruption threatens the US pandemic response. Carnegie Endowment for International Peace. URL: https://carnegieendowment.org/2020/04/13/how-global-corruption-threatens-u.s.-pandemic-response-pub-81545 [accessed 2021-05-12]

38. Tooze JA, Grunwald GK, Jones RH. Analysis of repeated measures data with clumping at zero. Stat Methods Med Res 2002 Aug;11(4):341-355. [doi: 10.1191/0962280202sm291ra] [Medline: 12197301]

39. Everitt BS, Johnson NL, Kotz S. Distributions in Statistics: Discrete Distributions. Journal of the Royal Statistical Society. Series A (General) 1970;133(3):482. [doi: 10.2307/2343567]

40. Lambert D. Zero-Inflated Poisson Regression, with an Application to Defects in Manufacturing. Technometrics 1992 Feb;34(1):1. [doi: 10.2307/1269547]

41. Fitzgerald M. Desensitization and the Coronavirus: How mass shootings have molded the American psyche. Medium. 2020. URL: https://medium.com/@mqfitzge/desensitization-and-the-coronavirus-how-mass-shootings-have-molded-the-american-psyche-be5217f52a6 [accessed 2021-02-04]

42. DiClemente C, Graydon M. Changing Behavior Using the Transtheoretical Model. In: Hagger MS, Cameron LD, Hamilton K, Hankonen N, Lintunen T, editors. The Handbook of Behavior Change. Cambridge: Cambridge University Press; Jul 2020:136-149.

43. Astroth DB, Cross-Poline GN, Stach DJ, Tilliss TSI, Annan SD. The transtheoretical model: an approach to behavioral change. J Dent Hyg 2002;76(4):286-295. [Medline: 12592920]

44. Slovic P, Västfjäll D. The More Who Die, the Less We Care: Psychic Numbing and Genocide. In: Kaul S, Kim D, editors. Imagining Human Rights. Berlin, München, Boston: De Gruyter; 2015:55.

45. Lawler M, Mackenzie S. What does cognitive dissonance mean? Everyday Health. URL: https://www.everydayhealth.com/neurology/cognitive-dissonance/what-does-cognitive-dissonance-mean-theory-definition/ [accessed 2021-02-04]

46. De Choudhury M, Sushovan D. Mental health discourse on reddit: Self-disclosure, social support, and anonymity. 2014 Presented at: Eighth International AAAI Conference on Weblogs and Social Media; June 2014; Ann Arbor, MI URL: https://ojs.aaai.org/index.php/ICWSM/article/view/14526

47. Hughes S, Adam W. 2018 Twitter Survey Dataset. Pew Research Center. 2019. URL: https://www.pewresearch.org/internet/dataset/2018-twitter-survey/ [accessed 2021-05-26]

48. Hargittai E. Potential Biases in Big Data: Omitted Voices on Social Media. Social Science Computer Review 2018 Jul 30;38(1):10-24. [doi: 10.1177/0894439318788322]

49. Bavel JJV, Baicker K, Boggio PS, Capraro V, Cichocka A, Cikara M, et al. Using social and behavioural science to support COVID-19 pandemic response. Nat Hum Behav 2020 May;4(5):460-471. [doi: 10.1038/s41562-020-0884-z] [Medline: 32355299]

## Abbreviations

**CDC:** Centers for Disease Control and Prevention
**LIWC:** Linguistic Inquiry and Word Count
**OR:** odds ratio

# The Role of the Canadian Media During the Initial Response to the COVID-19 Pandemic: A Topic Modelling Approach Using Canadian Broadcasting Corporation News Articles

Janhavi Patel[1*], BMSc; Harsheev Desai[2*], BS; Ali Okhowat[3], MD

[1]Michael G DeGroote School of Medicine, McMaster University, Hamilton, ON, Canada

[2]Faculty of Engineering and Architectural Science, Ryerson University, Toronto, ON, Canada

[3]UBC Global Health, Faculty of Medicine, University of British Columbia, Vancouver, BC, Canada

[*]these authors contributed equally

**Corresponding Author:**
Janhavi Patel, BMSc
Michael G DeGroote School of Medicine
McMaster University
1280 Main Street West
Hamilton, ON, L8S 4L8
Canada
Phone: 1 905 525 9140
Email: patelj10@mcmaster.ca

## Abstract

**Background:**   Beginning as a local epidemic, COVID-19 has since rapidly evolved into a pandemic. As countries around the world battle this outbreak, mass media has played an active role in disseminating public health information.

**Objective:**   The aim of this study was to get a better understanding of the role that the Canadian media played during the pandemic and to investigate the patterns of topics covered by media news reporting.

**Methods:**   We used a data set consisting of news articles published on the Canadian Broadcasting Corporation (CBC) website between December 2019 and May 2020. We then used Python software to analyze the data using Latent Dirichlet Allocation topic modelling. Subsequently, we used the pyLDAvis tool to plot these topics on a 2D plane through multidimensional scaling and divided these topics into different themes.

**Results:**   After removing articles that were published before the year 2019, we identified 6771 relevant news articles. According to the CV coherence value, we divided these articles into 15 topics, which were categorized into 6 themes. The three most popular themes were case reporting and testing (n=1738), Canadian response to the pandemic (n=1259), and changes to social life (n=1171), which accounted for 25.67%, 18.59%, and 17.29% of the total articles, respectively.

**Conclusions:**   Understanding the Canadian media's reporting on the COVID-19 pandemic shows that the Canadian pandemic response is a product of consistent government communication, as well as the public's understanding of and adherence to protocols.

## Introduction

COVID-19, which started as a local epidemic, evolved into a pandemic in a matter of months [1]. Countries around the world are battling the spread of this disease and the unfortunate consequences of COVID-19–related mortality and morbidity, resource limitations, and severe economic burden [1,2]. Canada

is no different and continues to observe a rising number of COVID-19 cases [3].

Due to the initial lack of vaccines and knowledge about the disease and its treatment, countries were forced to take unique approaches to combat the spread of the virus. Canada's response has been widely reported as being adequate, though much more could have been and remains to be done in tackling the spread

XSL•FO
**RenderX**

of COVID-19 [4]. The Canadian government website for its COVID-19 response highlights the measures Canada has taken, including the creation of the COVID Alert app, an ethics framework for policy makers, and economic support for Canadians. Support involves both financial measures and safety, including for Canadians abroad and vulnerable populations in Canada [3]. Additionally, there has been an emphasis on public education, collaboration, and guidance for researchers and frontline health workers [3].

With the uncertainty surrounding this novel coronavirus, the media—especially online news sources—have played a key role in informing the public about events related to the pandemic. Mass media has been successfully used for decades to increase public health awareness. News media outlets have been used across the globe for addressing public health issues like reducing tobacco use, participating in screening for cancer, and cardiovascular disease prevention [5]. The Eat Well campaign, which was advertised through a combination of news and commercial media outlets, increased awareness about meal prepping and healthy food choices in the Canadian population [6]. A postcampaign evaluation showed that low-resource communities had a greater uptake of information, thus highlighting the need to better understand the impact of different information dissemination campaigns to better cater to the target population [6].

The Canadian Broadcasting Corporation (CBC) is a daily source of local and national information for many Canadians [7]. The CBC's digital offering sees an average of 16.1 million new monthly visits [7] and continues to grow every month. Assessing the content of CBC articles can therefore provide insights into the information delivered to Canadians about the pandemic. Given that success in the fight against the pandemic greatly depends on the support of the public (eg, maintaining appropriate social distance and taking proper precautions), the information that media outlets report is important as it provides the public with up-to-date guidance.

The aim of this study was to better understand the role that news articles played in disseminating public health information, by specifically focusing on the topics reported and frequency of each topic reported regarding the COVID-19 pandemic. The methods used and results from this study could be relevant when reporting future events related to health care and national safety, which rely heavily on public support and awareness.

## Methods

### Data Collection

The data set was collected from the CBC website using a Python programming language script [8]. The script was used to extract information from over 6700 news articles, including the title, article summary, and main text for each article, using the term "coronavirus" as the search word. The extracted news articles were published between December 2012 and May 2020; however, only the articles published in 2019 and 2020 were included in this study.

We used Latent Dirichlet Allocation (LDA) to analyze these news articles. LDA is a three-level hierarchical Bayesian model, in which each item of a collection is modelled as a finite mixture over an underlying set of topics. The basic idea behind LDA is that documents can be represented as arbitrary mixtures over latent topics, which in turn are characterized by a distribution over words [9]. LDA has been extensively used and evaluated for its applicability in topic modelling research [10,11]. Moreover, Lancichinetti [12] showed that LDA has high reproducibility and accuracy for topic classification.
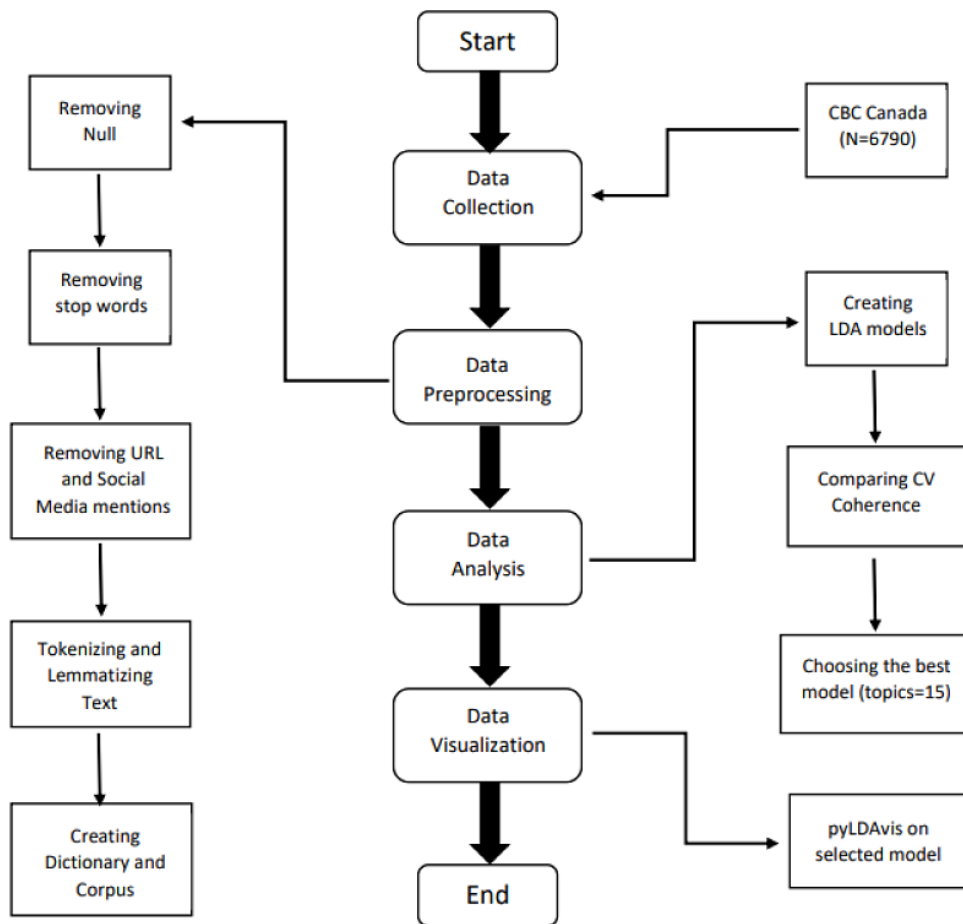
According to LDA, there are diverse topics in each news article, and the words in these articles can be allotted to one of these topics. However, LDA only groups inputs (ie, news articles in this case) based on the abovementioned distribution over words and it is subjective how these groups are interpreted as topics. To facilitate accurate representation, randomly selected articles from each topic were manually checked to make sure they were consistent with the interpreted topic.

### Data Processing

There were a total of 6771 news articles remaining after removing the articles published before 2019. These remaining articles were dated between December 22, 2019, and May 3, 2020.

Before moving forward with topic modelling, we used Python along with libraries, including the Pandas and Natural Language Toolkit (NLTK) libraries [13], to clean the data. The detailed process for this is displayed in Figure 1. We used the English language stop words provided by NLTK to remove common words such as "an," "all," "and," "for," and "from" as they hold no semantic value for our analysis. URLs and social media mentions consisting of "@" were also removed. The two primary inputs to the LDA model are the dictionary and the corpus, which were created using the Gensim library [14].

**Figure 1.** Workflow chart. CBC: Canadian Broadcasting Corporation; LDA: Latent Dirichlet Allocation.
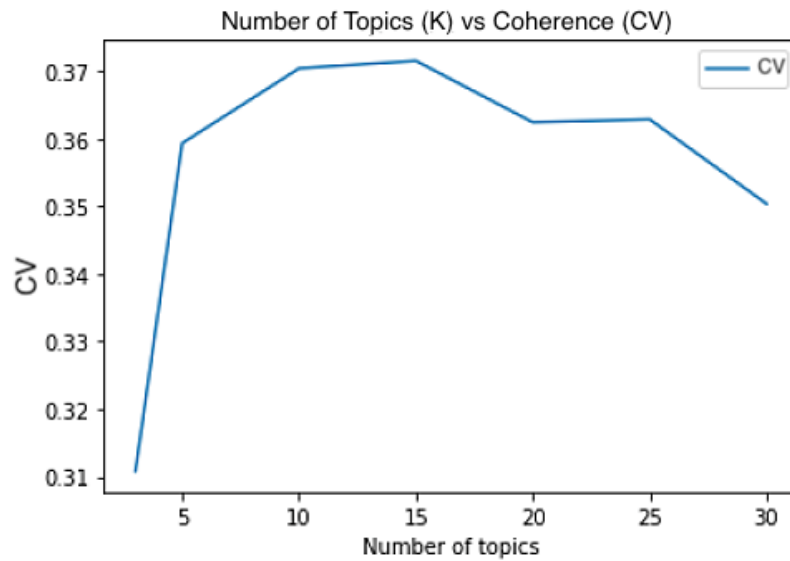


We used the CV coherence score to evaluate models with different numbers of topics and selected the one with the highest CV score. This approach mitigates one of LDA's limitations—the need to know the number of topics ahead of time. According to Röder et al [15], the CV coherence score is one of the fastest measures of coherence, and the most accurate. Henri Trenquier defines coherence as the human's semantic appreciation of a topic represented by its N top words [16]. We chose the top 15 (N) words in each topic to calculate the coherence.

As evident from Figure 2, the highest CV coherence score was achieved at 15 topics. This means that the top words in each topic were most closely related semantically when the news article data set was divided into 15 different topics using LDA.
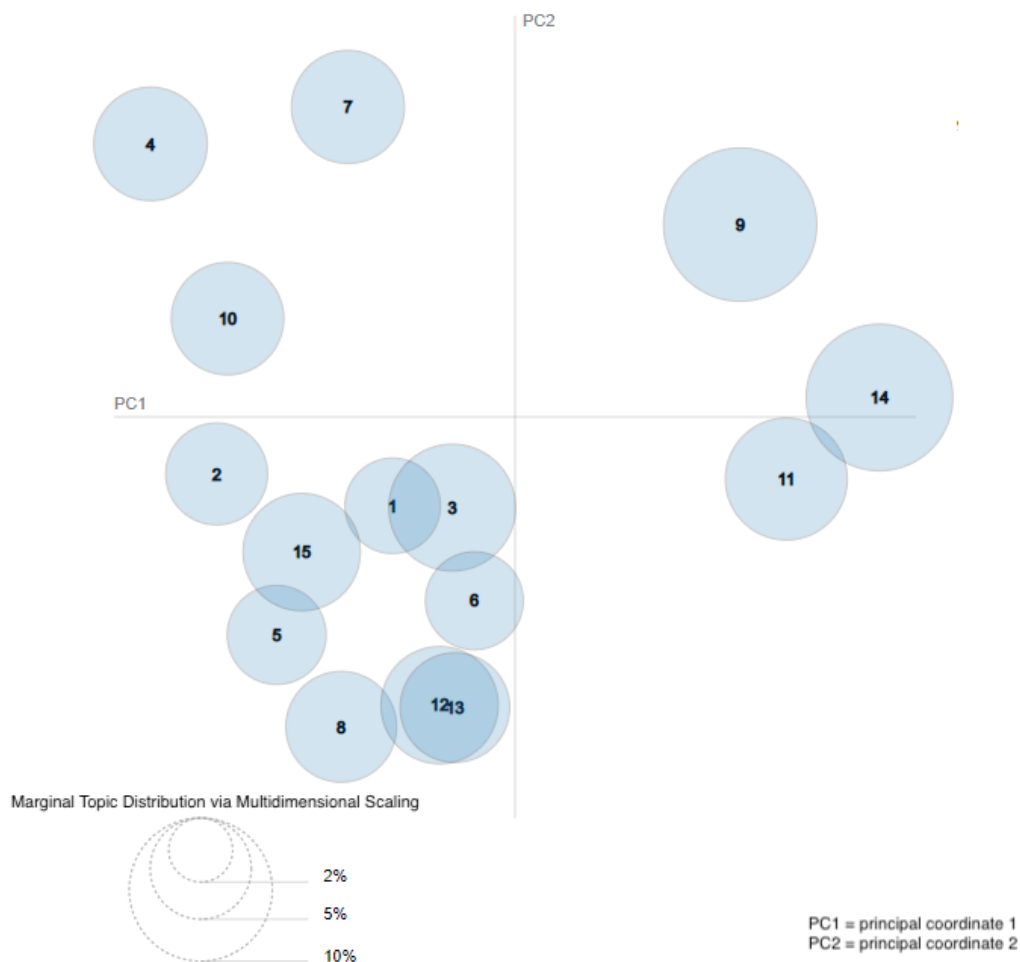
**Figure 2.** Coherence score for number of topics.



We then used the pyLDAvis tool [17] and Python to further analyze the 15 topics to extract valuable insights from the articles. The 15 topics were represented on an intertopic distance map, which is an interactive representation offered by the pyLDAvis tool (Figure 3). The topics are plotted as circles in a 2D plane whose centers are determined by computing the distance between topics [16].

**Figure 3.** Intertopic distance map.

The weight parameter λ was adjusted to find the theme for each topic based on the top words in the topic. Setting λ=1 ranks the words in a topic by frequency, while setting λ=0 ranks the words based on uniqueness to that topic [17]. We used the interactive bar provided by the pyLDAvis tool to adjust λ and understand the theme for each of the 15 topics. To use topic 9 as an example, Figure 4 shows the top 30 most frequently occurring words in topic 9. As multiple topics might have similar words

that occur frequently, we need to adjust the λ value to better gauge what topic an article might be about. For instance, when we set λ to 0.04, the terms most unique to topic 9 are captured, and presented in descending order in Figure 5. This analysis identifies words like "test," "positive," and "spread" as being unique to topic 9. Using the keywords, we identified that the general theme of articles in topic 9 is "testing." This process was repeated for each topic (Table 1).
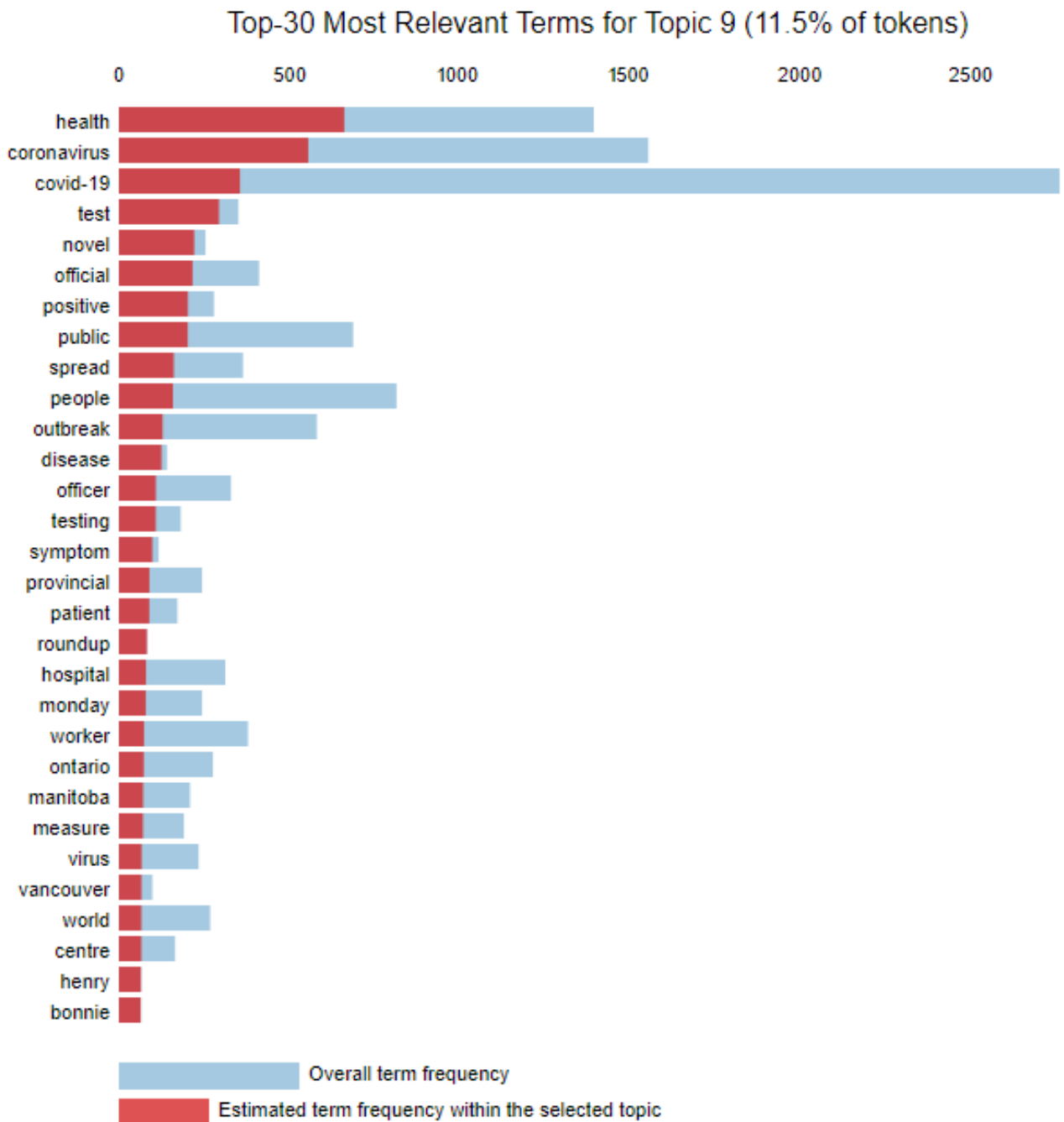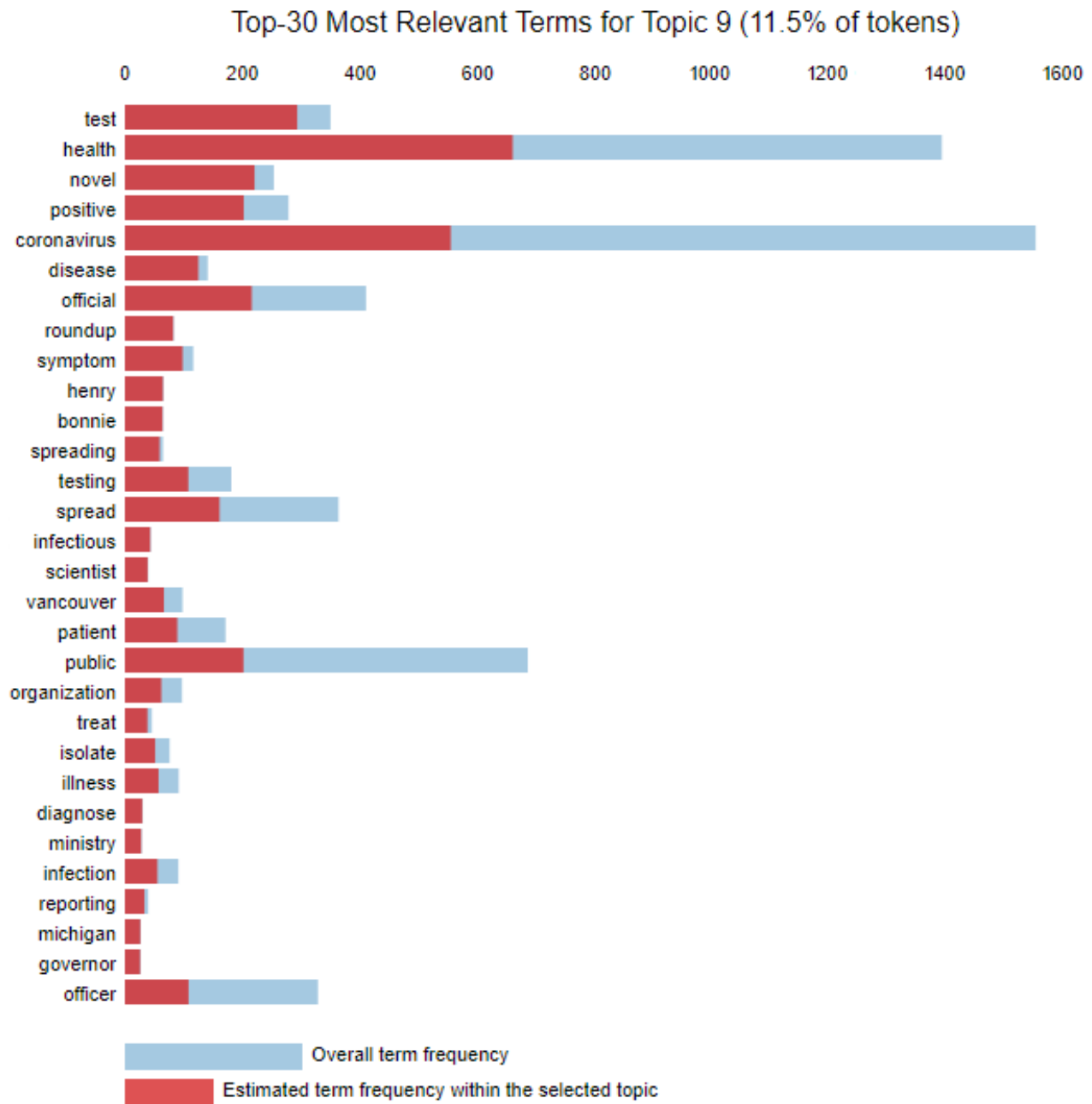
**Figure 4.** Top 30 most relevant terms (λ=1.0).



Top-30 Most Relevant Terms for Topic 9 (11.5% of tokens)

**Figure 5.** Top 30 most relevant terms (λ=0.4).



Top-30 Most Relevant Terms for Topic 9 (11.5% of tokens)

**Table 1.** Themes and topics (N=6771).

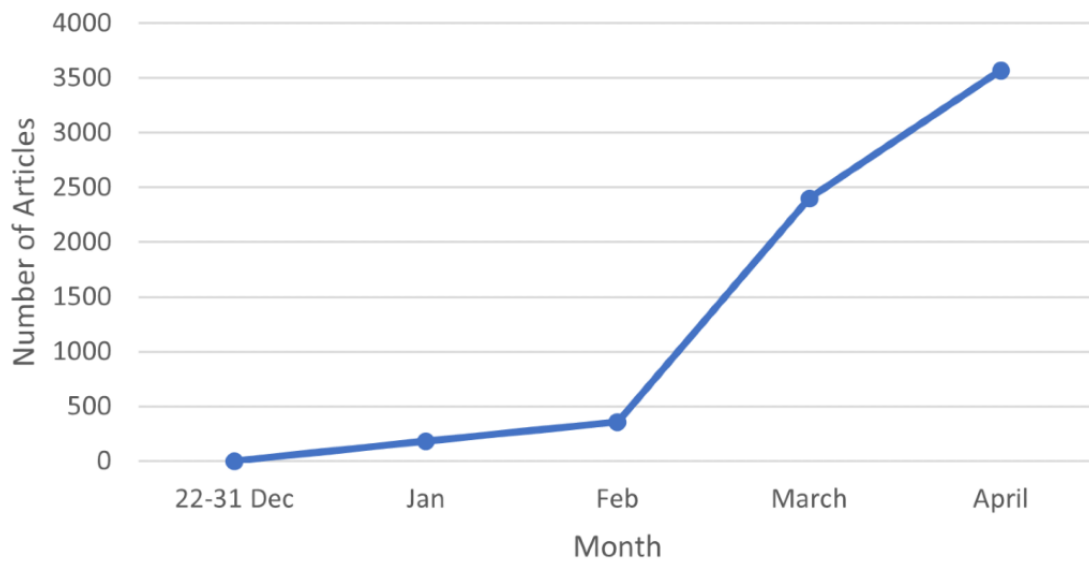| Themes and topics | Number of news articles, n (%)[a] | Keywords |
|---|---|---|
| **Theme 1: Case reporting and testing (n=1738)** | | |
| Topic 9: Testing | 876 (12.97) | Coronavirus, spread, public, positive, novel, official, people, health, test, covid |
| Topic 14: Case reporting | 862 (12.73) | Province, number, confirm, total, people, report, health, death, case, covid |
| **Theme 2: Canadian response to pandemic (n=1259)** | | |
| Topic 5: General response | 295 (4.37) | Pandemic, nation, member, Windsor, covid, community, family, first, local, want |
| Topic 6: Health care/hospital response | 267 (3.93) | Emergency, temporary, staff, state, hospital, Sudbury, worker, declare, general, covid |
| Topic 10: Vaccine research | 399 (5.88) | Ottawa, world, around, Canada, global, point, latest, covid, coronavirus, point |
| Topic 1: Medical supplies and resources | 298 (4.40) | Canadian, ventilator, doctor, could, Canada, happening, mask, province, available, covid |
| **Theme 3: Changes to everyday life (n=1171)** | | |
| Topic 2: Social gathering cancellations | 322 (4.76) | Summer, pandemic, coronavirus, cancel, festival, university, event, season, covid, plant |
| Topic 8: School closure/virtual learning | 397 (5.88) | Parent, school, family, child, learning, student, covid, equipment, worker, pandemic |
| Topic 12: General lifestyle changes | 452 (6.68) | People, avoid, coming, together, change, normal, covid, province, pandemic, government |
| **Theme 4: Communication from the government (n=1002)** | | |
| Topic 11: Public health announcements | 481 (7.04) | Medical, chief, public, officer, health, people, province, covid, Friday |
| Topic 3: Prime Minister's addresses | 521 (7.68) | Minister, prime, Justin, Trudeau, worker, pandemic, essential, coronavirus, health, covid |
| **Theme 5: International news (n=826)** | | |
| Topic 4: Articles related to news in the United States | 388 (5.73) | Trump, unite, outbreak, state, president, country, cruise, coronavirus, Canada, Canadian |
| Topic 7: Articles related to news in China | 438 (6.48) | Chinese, china, outbreak, answer, Canadian, flight, morning, expert, question, expert |
| **Theme 6: Government initiatives (n=775)** | | |
| Topic 13: Initiatives for vulnerable populations | 343 (5.08) | Shelter, homeless, social, distance, encourage, people, pandemic, covid, measure, physical |
| Topic 15: Economy and business | 432 (6.39) | Business, economy, government, federal, pandemic, support, covid, million, premier Canada |

[a]The percentages have been calculated using the N value (ie, 6771).

## Results

Using the pyLDAvis tool, we grouped the 15 topics into 6 themes as shown in Table 1. Theme 1 (case reporting and testing) had the greatest number of articles (n=1738), while theme 6 (government initiatives) represented the lowest number of news reports (n=775). The trend in the total frequency of articles related to COVID-19 over our study time period is shown in Figure 6.
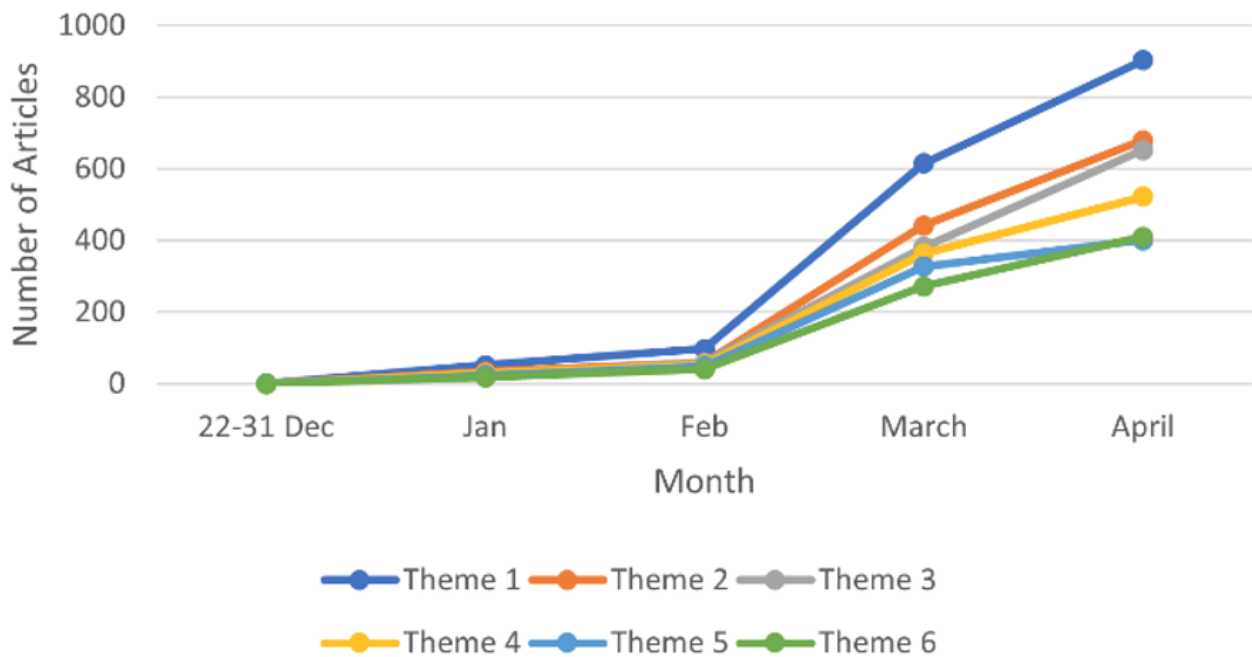
**Figure 6.** Time trend for number of articles.



The most frequent theme, theme 1 (case reporting and testing), consisted primarily of topics that covered articles related to information about testing (12.97%) and case reporting (12.73%). The information in the articles related to testing focused on information regarding the tests being conducted to assess the spread of the virus, whereas the articles related to case reporting primarily focused on reporting the number of confirmed cases

and deaths around the country, with words like "number," "report," "confirm," and "death" being frequently used. Similar to the trend in frequency of total articles about COVID-19, Theme 1 saw a sudden increase in the number of articles published, starting from the month of February and increasing throughout March and April (Figure 7).

**Figure 7.** Time trends for each theme.



The Canadian media's focus in relation to the outbreak and Canada's response is highlighted in themes 2 and 3. Theme 2 (Canadian response to pandemic) includes topics like general response (topic 5, n=295, 4.37%), health care/hospital response (topic 6, n=267, 3.93%), vaccine research (topic 10, n=399,

5.88%), and medical supplies and resources (topic 1, n=298, 4.40%). Theme 3 (changes to everyday life) discusses the changes resulting from the pandemic and includes topics like social gathering cancellations (topic 2, n=322, 4.76%), school closure/virtual learning (topic 8, n=397, 5.88%), and general

lifestyle changes (topic 12, n=452, 6.68%). Both these themes saw a steep increase in the number of reported articles after the month of February.

Theme 5 (international news) consisted primarily of topics related to the United States (n=388, 5.73%) and China (n=438, 6.48%). The information in articles related to the United States, Canada's geographic neighbor and largest trading partner, focused more on political relations, with frequently used words like "president," "Trump," and "state." On the contrary, the articles about China, the place of origin of the coronavirus, primarily focused on information that would enable a better understanding of the outbreak, with words like "outbreak," "question," and "answer" used more frequently.

Themes 4 and 6, although not high in frequency, focused on important themes like communication from the government and government initiatives. Communication from the government included both public health announcements as well as the Prime Minister's addresses to the public. The government initiatives theme included topics that discussed initiatives for vulnerable populations, specifically people experiencing homelessness (topic 13, n=343, 5.008%), as well as the economy and business (topic 15, n=432, 6.39%). The articles about government announcements had a steep increase leading to the declaration of pandemic in March 2020 but the slope reduced from March to April. Contrastingly, theme 6 saw a steep, consistent increase from February to April 2020.

## Discussion

### Principal Results

The COVID-19 pandemic has been a steep learning curve for all countries worldwide. Dissemination of information in a timely manner across communities and countries was crucial to limit the spread of COVID-19 and determine the efficacy of different treatment and management interventions. With the ensuing social isolation, online media took over as an important source of information available to the public; thus, understanding the role that the media played highlights key aspects of the challenges faced by Canada and its response to the pandemic. We used topic modelling using articles collected from the CBC's online platform and identified different themes reported through the articles.

Even though some reports about a fatal pneumonia of unknown cause had started coming out from China in early January, it was not until after the World Health Organization declared COVID-19 as a global health emergency that articles about the virus started increasing in Canada. The number of articles about COVID-19 showed a sharp increase starting February 2020 for most themes, after the World Health Organization declared it a global health emergency on January 30. Resource shortages and panic buying have been an issue in many countries battling COVID-19 [18]. Our study identified that there were about 300 articles focused on resources. It is postulated that anxiety around sudden lockdowns and uncertainty about the duration of the pandemic might have contributed to the response of preservation of self and family [18]. Retrospectively, it would be beneficial for health care professionals, the government, and the media to work closely together to provide better guidelines and policies for the public, both to reduce anxiety and ensure more equitable distribution of resources.

Additionally, our topic modelling showed that a considerable proportion of news articles in the study focused on the conditions of marginalized populations, such as people experiencing homelessness. Many people experiencing homelessness did not receive timely shelter and space to self-isolate, putting their lives and the lives of others at risk. Thus, our study results highlight the importance of creating an equitable response strategy during future pandemics.

Throughout the course of the pandemic, the most reported information was regarding testing and case reporting. This is consistent with any communicable disease, wherein proper testing, contact tracing, and case reporting are crucial to control the spread of the disease [19]. This information can contribute to increased anxiety, as witnessed by people's fear of acquiring the disease from health care facilities, and thus being reluctant to access care for other acute illnesses, including heart attacks and strokes [20]. On the other hand, having this information could make people feel more accountable for their actions and encourage them to be more socially responsible. Although the neglect of other conditions was an unintended, unfortunate consequence of pandemic-related public health measures, for future events, more holistic communication from health care professionals (ie, about considering other acute illnesses in times of crisis) and reporting from the media on this topic could aid in better management of people with acute and chronic illnesses.

### Limitations

This study only contains news articles published on CBC's online platform that were tagged with the term "coronavirus." There are several other sources of media available in Canada and future studies can focus on including multiple different sources of both digital and print media. The pandemic is ongoing and Canada's response and policies are constantly changing. Thus, doing a long-term study and constantly monitoring multiple media outlets' efforts will be helpful for future studies. This study nonetheless provides a glimpse of the Canadian media's role in the communication and dissemination of information. The LDA model has certain limitations; for example, the different topics need to be manually interpreted and are open to misinterpretation or overinterpretation. Some of its other limitations include the inability to capture correlations between topics and the use of a fixed number of topics, which must be known ahead of time.

### Comparison With Prior Work

Our study identified several similar and unique themes compared to the themes identified by another similar study on Chinese media reporting [21]. Topics like case reporting, disease spread, medical supplies and resources, and research and development were similarly observed in media in both studies. However, the Chinese study did not identify any themes related to communication from the government or the country's response regarding vulnerable populations. In contrast, although lower in frequency compared to other topics, Canada's media and response focused on ensuring proper communication from the

government and support for vulnerable populations. The government actively communicated with the public, not only through public health officials but also via regular addresses from the country's prime minister during the pandemic. Studies have shown that a leader's address to the public is very effective in reassuring people during times of crisis [22] and the media reports suggest that it played a big part in Canada's response to the pandemic.

Compared to initial communication during the H1N1 pandemic in 2009, which involved the dissemination of misinformation, leading to widespread panic, the slow dissemination of public health information by media outlets initially led to panic early in the COVID-19 pandemic. After the H1N1 pandemic, the Centers for Disease Control and Prevention conducted an audit on public health information dissemination and provided several guidelines for communications in future pandemics [23]. In line with the guidelines, our study topics found that the Canadian response had consistent messaging from federal government and public health officials; however, Canada's response still fell short with regard to prioritizing marginalized populations and reducing the public's initial stress stemming from widespread misinformation.

## Conclusions

Our study highlights that, based on the topical analysis of CBC news articles, the Canadian response to the COVID-19 pandemic was a joint effort guided by government policies and communications in conjunction with people's response and adherence to protocol.

One of the most important factors in preventing the spread of COVID-19 is to empower the public with accurate information [24].

The media plays an important bridging role by relaying information from the government to the public. Thus, by understanding and analyzing the extent to which certain events and policies affect public sentiment and response, policy makers can proactively improve communication for any similar future events, including pandemics, natural disasters, or issues related to national safety.

## Conflicts of Interest

None declared.

## References

1. Shereen MA, Khan S, Kazmi A, Bashir N, Siddique R. COVID-19 infection: Origin, transmission, and characteristics of human coronaviruses. J Adv Res 2020 Jul;24:91-98 [FREE Full text] [doi: 10.1016/j.jare.2020.03.005] [Medline: 32257431]
2. Clark A, Jit M, Warren-Gash C, Guthrie B, Wang HHX, Mercer SW, Centre for the Mathematical Modelling of Infectious Diseases COVID-19 working group. Global, regional, and national estimates of the population at increased risk of severe COVID-19 due to underlying health conditions in 2020: a modelling study. Lancet Glob Health 2020 Aug;8(8):e1003-e1017 [FREE Full text] [doi: 10.1016/S2214-109X(20)30264-3] [Medline: 32553130]
3. Government of Canada. Coronavirus disease (COVID-19): Outbreak update. 2020. URL: https://www.canada.ca/en/public-health/services/diseases/2019-novel-coronavirus-infection.html [accessed 2020-08-07]
4. Johns Hopkins University and Medicine. COVID-19 Dashboard by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (JHU). 2020. URL: https://coronavirus.jhu.edu/map.html [accessed 2020-08-07]
5. Institute of Medicine (US) Committee on Assuring the Health of the Public in the 21st Century. The Future of the Public's Health in the 21st Century. Washington, DC: National Academies Press (US); 2002.
6. Fernandez MA, Desroches S, Marquis M, Lebel A, Turcotte M, Provencher V. Promoting meal planning through mass media: awareness of a nutrition campaign among Canadian parents. Public Health Nutr 2019 Oct 30;22(18):3349-3359. [doi: 10.1017/s1368980019002957]
7. CBC/Radio-Canada. Our Performance - Media Lines. 2018. URL: https://site-cbc.radio-canada.ca/site/annual-reports/2017-2018/accountability-plan/our-performance-media-lines-english-services-highlights-en.html [accessed 2020-08-09]
8. Han R. COVID-19 News Articles Open Research Dataset, Version 3. URL: https://www.kaggle.com/ryanxjhan/cbc-news-coronavirus-articles-march-26 [accessed 2020-08-05]
9. Blei D, Ng A, Jordan M. Latent Dirichlet Allocation. The Journal of Machine Learning Research 2003;3:993-1022 [FREE Full text]
10. Blei D. Probabilistic topic models. Commun ACM 2012;55(4):77-84 [FREE Full text] [doi: 10.1145/2133806.2133826]
11. Jacobi C, van Atteveldt W, Welbers K. Quantitative analysis of large amounts of journalistic texts using topic modelling. Digital Journalism 2015 Oct 13;4(1):89-106 [FREE Full text] [doi: 10.1080/21670811.2015.1093271]
12. Lancichinetti A, Sirer M, Wang J, Acuna D, Körding K, Amaral L. High-Reproducibility and High-Accuracy Method for Automated Topic Classification. Phys Rev X 2015 Jan 29;5(1):1 [FREE Full text] [doi: 10.1103/physrevx.5.011007]
13. Bird S, Klein E, Loper E. Natural Language Processing with Python. Sebastopol, CA: O'Reilly Media; 2009.
14. Rehurek R, Sojka P. Software Framework for Topic Modelling with Large Corpora. In: Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks. 2019 Presented at: LREC 2010 Workshop on New Challenges for NLP Frameworks; May 22, 2010; Valletta, Malta p. 46-50 URL: https://is.muni.cz/publication/884893/en/Software-Framework-for-Topic-Modelling-with-Large-Corpora/Rehurek-Sojka

XSL•FO

RenderX

15.   Röder M, Both A, Hinneburg A. Exploring the Space of Topic Coherence Measures. In: Eighth ACM International Conference on Web Search and Data Mining. 2015 Feb 02 Presented at: Eighth ACM International Conference on Web Search and Data Mining; 2015; Shanghai, China p. 399-408 URL: https://doi.org/10.1145/2684822.2685324 [doi: 10.1145/2684822.2685324]

16.   Trenquier H. Improving Semantic Quality of Topic Models for Forensic Investigation. University of Amsterdam. 2018. URL: https://www.os3.nl/_media/2017-2018/courses/rp2/p76_report.pdf [accessed 2020-07-30]

17.   Sievert C, Shirley K. LDAvis: A method for visualizing and interpreting topics. In: Proceedings of the Workshop on Interactive Language Learning, Visualization, and Interfaces. 2014 Jun Presented at: Workshop on Interactive Language Learning, Visualization, and Interfaces; June 2014; Baltimore, MA. [doi: 10.3115/v1/w14-3110]

18.   Sim K, Chua HC, Vieta E, Fernandez G. The anatomy of panic buying related to the current COVID-19 pandemic. Psychiatry Res 2020 Jun;288:113015 [FREE Full text] [doi: 10.1016/j.psychres.2020.113015] [Medline: 32315887]

19.   Enanoria W, Liu F, Zipprich J, Harriman K, Ackley S, Blumberg S, et al. The Effect of Contact Investigations and Public Health Interventions in the Control and Prevention of Measles Transmission: A Simulation Study. PLoS One 2016;11(12):e0167160 [FREE Full text] [doi: 10.1371/journal.pone.0167160] [Medline: 27941976]

20.   New data confirms significant drop in heart attack patients presenting at hospital. Heart and Stroke Foundation. URL: https://www.heartandstroke.ca/what-we-do/media-centre/news-releases/news-release-new-data-confirms-significant-drop-in-heart-attack-patients-presenting-at-hospital [accessed 2020-08-09]

21.   Liu Q, Zheng Z, Zheng J, Chen Q, Liu G, Chen S, et al. Health Communication Through News Media During the Early Stage of the COVID-19 Outbreak in China: Digital Topic Modeling Approach. J Med Internet Res 2020 Apr 28;22(4):e19118 [FREE Full text] [doi: 10.2196/19118] [Medline: 32302966]

22.   Deitchman S. Enhancing crisis leadership in public health emergencies. Disaster Med Public Health Prep 2013 Oct;7(5):534-540. [doi: 10.1017/dmp.2013.81] [Medline: 24274133]

23.   Institute of Medicine (US) Forum on Medical and Public Health Preparedness for Catastrophic Events. The 2009 H1N1 Influenza Vaccination Campaign: Summary of a Workshop Series. Washington, DC: National Academies Press (US); 2010.

24.   Reddy B, Gupta A. Importance of effective communication during COVID-19 infodemic. J Family Med Prim Care 2020 Aug;9(8):3793-3796 [FREE Full text] [doi: 10.4103/jfmpc.jfmpc_719_20] [Medline: 33110769]

## Abbreviations

**CBC:** Canadian Broadcasting Corporation
**LDA:** Latent Dirichlet Allocation
**NLTK:** Natural Language Toolkit

XSL·FO
RenderX

Original Paper

# A Public Health Research Agenda for Managing Infodemics: Methods and Results of the First WHO Infodemiology Conference

Neville Calleja[1*], MD, PhD; AbdelHalim AbdAllah[2*], BA; Neetu Abad[3*], PhD; Naglaa Ahmed[4*], PhD; Dolores Albarracin[5*], PhD; Elena Altieri[6*], MA, MBA; Julienne N Anoko[7*], PhD, MPH; Ruben Arcos[8*], PhD; Arina Anis Azlan[9*], PhD; Judit Bayer[10,11*], PhD, Habil; Anja Bechmann[12*], PhD; Supriya Bezbaruah[13*], MPH, PhD; Sylvie C Briand[14*], MD, MPH, PhD; Ian Brooks[15*], PhD; Lucie M Bucci[16*], MA; Stefano Burzo[17*], MA; Christine Czerniak[14*], PhD; Manlio De Domenico[18*], PhD; Adam G Dunn[19*], PhD; Ullrich K H Ecker[20*], PhD; Laura Espinosa[21*], MPH, DVM; Camille Francois[22*], PhD; Kacper Gradon[23*], PhD, Hab; Anatoliy Gruzd[24*], PhD; Beste Sultan Gülgün[25*], MSc; Rustam Haydarov[26*], MSc; Cherstyn Hurley[27*], MA; Santi Indra Astuti[28*], MSi; Atsuyoshi Ishizumi[3,29*], MPH, MSc; Neil Johnson[30*], PhD; Dylan Johnson Restrepo[30*]; Masato Kajimoto[31*], PhD; Aybüke Koyuncu[3*], MPH; Shibani Kulkarni[3,29*], MPH, PhD; Jaya Lamichhane[14*], MA, MBA; Rosamund Lewis[32*], MDCM, MSc; Avichal Mahajan[14*], PhD; Ahmed Mandil[4*], MBChB, DrPH; Erin McAweeney[22*], PhD; Melanie Messer[33*], MA, RN, DrPH; Wesley Moy[34*], MBA, MSSI, MSS, PhD; Patricia Ndumbi Ngamala[35*], PhD; Tim Nguyen[14*], MSc; Mark Nunn[*], MA; Saad B Omer[36*], MD, PhD; Claudia Pagliari[37*], PhD; Palak Patel[3,29*], MBBS, MPH; Lynette Phuong[14*], MPH; Dimitri Prybylski[3*], MPH, PhD; Arash Rashidian[4*], PhD; Emily Rempel[38*], MSc, PhD; Sara Rubinelli[39,40*], PhD; PierLuigi Sacco[41,42*], PhD; Anton Schneider[43*], BA; Kai Shu[44*], PhD; Melanie Smith[22*], PhD; Harry Sufehmi[45*], MSc; Viroj Tangcharoensathien[46*], PhD; Robert Terry[47*], MPhil, PhD; Naveen Thacker[48*], MD; Tom Trewinnard[49*], MBA; Shannon Turner[50,51*], MSc, PhD; Heidi Tworek[52*], PhD; Saad Uakkas[53*], MD; Emily Vraga[54*], PhD; Claire Wardle[55*], PhD; Herman Wasserman[56*], PhD; Elisabeth Wilhelm[3*], MSc; Andrea Würz[21*], MA; Brian Yau[57*], BSc; Lei Zhou[58*], MD; Tina D Purnat[35*], MSc

[1]Directorate for Health Information & Research, Ministry for Health, Valetta, Malta

[2]WHO Regional Office for Africa, Brazzaville, Congo

[3]US Centers for Disease Control and Prevention, Atlanta, GA, United States

[4]WHO Regional Office for Eastern Mediterranean, Cairo, Egypt

[5]Department of Psychology, College of Liberal Arts & Sciences, University of Illinois Urbana-Champaign, Champaign, IL, United States

[6]Department of Communications, World Health Organization, Geneva, Switzerland

[7]WHO Regional Office for Africa, Dakar, Senegal

[8]Department of Communication Sciences and Sociology, Communication Sciences Faculty, University Rey Juan Carlos, Madrid, Spain

[9]Faculty of Social Sciences and Humanities, Universiti Kebangsaan Malaysia, Bangi, Malaysia

[10]Department of Communication, Budapest Economics University (BGE), Budapest, Hungary

[11]Institute for Information, Telecommunications and Media Law, University of Münster (WWU), Münster, Germany

[12]DATALAB - Center for Digital Social Research, School of Communication and Culture, Aarhus University, Aarhus, Denmark

[13]WHO Regional Office for South East Asia, New Delhi, India

[14]Department of Infectious Hazards Management, Emergency Preparedness Division, World Health Organization, Geneva, Switzerland

[15]Center for Health Informatics, School of Information Sciences, University of Illinois at Urbana-Champaign, Champaign, IL, United States

[16]Immunize Canada, Canadian Public Health Association, Ottawa, ON, Canada

[17]Department of Political Science, University of British Columbia, Vancouver, BC, Canada

[18]CoMuNe Lab, Fondazione Bruno Kessler, Povo, Italy

[19]Biomedical Informatics and Digital Health, School of Medical Sciences, The University of Sydney, Sydney, Australia

[20]School of Psychological Science, The University of Western Australia, Perth, Australia

[21]European Centre for Disease Prevention and Control, Stockholm, Sweden

[22]Graphika, New York, NY, United States

[23]Department of Security and Crime Science, University College London, London, United Kingdom

[24]Ted Rogers School of Management, Ryerson University, Toronto, ON, Canada

[25]Ministry of Health, Ankara, Turkey

[26]UNICEF Headquarters, New York, NY, United States

[27]Immunisation and Countermeasures Department, Public Health England, London, United Kingdom

[28]The Faculty of Communication Science, Bandung Islamic University (UNISBA), Bandung, Indonesia

[29]Oak Ridge Institute for Science and Education, Oak Ridge, TN, United States

[30]Department of Physics, George Washington University, Washington, DC, United States

[31]Journalism and Media Studies Centre, The University of Hong Kong, Hong Kong, China

[32]Emergency Preaparedness Division, World Health Organization, Geneva, Switzerland

[33]Faculty I, Department of Nursing Science II, Trier University, Trier, Germany

[34]Advanced Academic Programs, Johns Hopkins University, Washington, DC, United States

[35]Department of Digital Health and Innovation, Science Division, World Health Organization, Geneva, Switzerland

[36]Yale Institute for Global Health, Yale University, New Haven, CT, United States

[37]Usher Institute, Edinburgh Medical School, University of Edinburgh, Edinburgh, United Kingdom

[38]British Columbia Centre for Disease Control, Vancouver, BC, Canada

[39]Department of Health Sciences and Medicine, University of Lucerne, Lucerne, Switzerland

[40]Swiss Paraplegic Research, Lucerne, Switzerland

[41]Department of Humanities Studies, Free University of Languages and Communication IULM, Milan, Italy

[42]metaLAB (at) Harvard, Harvard University, Cambridge, MA, United States

[43]Office of Infectious Disease, Global Health Bureau, United States Agency for International Development (USAID), Washington, DC, United States

[44]Computer Science Department, Illinois Institute of Technology, Chicago, IL, United States

[45]Masyarakat Anti Fitnah Indonesia (MAFINDO), Jakarta, Indonesia

[46]International Health Policy Programme, Ministry of Public Health, Bangkok, Thailand

[47]Science Division, World Health Organization, Geneva, Switzerland

[48]Deep Children Hospital and Research Centre, Gandhidham, India

[49]Fathm, London, United Kingdom

[50]Public Health Association of British Columbia, Victoria, BC, Canada

[51]Vaccine Safety Net (VSN), Geneva, Switzerland

[52]Department of History, University of British Columbia, Vancouver, BC, Canada

[53]Faculty of Medicine, Mohamed V University in Rabat, Rabat, Morocco

[54]Hubbard School of Journalism and Mass Communication, University of Minnesota, Minneapolis, MN, United States

[55]First Draft News, New York, NY, United States

[56]Centre for Film and Media Studies, University of Cape Town, Cape Town, South Africa

[57]Department of Regulation and Prequalification, Access to Medicines and Health Products Division, World Health Organization, Geneva, Switzerland

[58]Public Health Emergency Center, Chinese Center for Disease Control and Prevention, Beijing, China

[*]all authors contributed equally

**Corresponding Author:**
Tim Nguyen, MSc
Department of Infectious Hazards Management
Emergency Preparedness Division
World Health Organization
Avenue Appia 20
Geneva, 1211
Switzerland
Phone: 41 22 791 21 11
Email: nguyent@who.int

## Abstract

**Background:** An infodemic is an overflow of information of varying quality that surges across digital and physical environments during an acute public health event. It leads to confusion, risk-taking, and behaviors that can harm health and lead to erosion of trust in health authorities and public health responses. Owing to the global scale and high stakes of the health emergency, responding to the infodemic related to the pandemic is particularly urgent. Building on diverse research disciplines and expanding the discipline of infodemiology, more evidence-based interventions are needed to design infodemic management interventions and tools and implement them by health emergency responders.

XSL•FO

**RenderX**

**Objective:** The World Health Organization organized the first global infodemiology conference, entirely online, during June and July 2020, with a follow-up process from August to October 2020, to review current multidisciplinary evidence, interventions, and practices that can be applied to the COVID-19 infodemic response. This resulted in the creation of a public health research agenda for managing infodemics.

**Methods:** As part of the conference, a structured expert judgment synthesis method was used to formulate a public health research agenda. A total of 110 participants represented diverse scientific disciplines from over 35 countries and global public health implementing partners. The conference used a laddered discussion sprint methodology by rotating participant teams, and a managed follow-up process was used to assemble a research agenda based on the discussion and structured expert feedback. This resulted in a five-workstream frame of the research agenda for infodemic management and 166 suggested research questions. The participants then ranked the questions for feasibility and expected public health impact. The expert consensus was summarized in a public health research agenda that included a list of priority research questions.

**Results:** The public health research agenda for infodemic management has five workstreams: (1) measuring and continuously monitoring the impact of infodemics during health emergencies; (2) detecting signals and understanding the spread and risk of infodemics; (3) responding and deploying interventions that mitigate and protect against infodemics and their harmful effects; (4) evaluating infodemic interventions and strengthening the resilience of individuals and communities to infodemics; and (5) promoting the development, adaptation, and application of interventions and toolkits for infodemic management. Each workstream identifies research questions and highlights 49 high priority research questions.

**Conclusions:** Public health authorities need to develop, validate, implement, and adapt tools and interventions for managing infodemics in acute public health events in ways that are appropriate for their countries and contexts. Infodemiology provides a scientific foundation to make this possible. This research agenda proposes a structured framework for targeted investment for the scientific community, policy makers, implementing organizations, and other stakeholders to consider.

## Introduction

A pneumonia of unknown cause detected in Wuhan, China, was first reported to the World Health Organization (WHO) Country Office in China on December 31, 2019. The disease, caused by a novel coronavirus (SARS-CoV-2), was subsequently named COVID-19, and it was declared a Public Health Emergency of International Concern on January 30, 2020. On March 11, 2020, the WHO characterized the outbreak as a pandemic. Globally, as of August 23, 2021, 211,373,303 confirmed cases of COVID-19, including 4,424,341 deaths, had been reported to the WHO [1].

On February 15, 2020, WHO Director-General Tedros Adhanom Ghebreyesus raised the concern that the epidemic was accompanied by an infodemic [2]. An infodemic is an overflow of information of varying quality that surges across digital and physical environments during an acute public health event and makes it difficult for people to find information to better protect themselves and their communities [3]. An infodemic can lead to confusion, misunderstanding of health information, risk-taking, and behaviors that can harm health, hinder the public health response, and lead to mistrust in health authorities. [4]. Therefore, people need timely, accurate, and accessible information in the right format and amount during epidemics to adopt health-promoting behavior to protect themselves, their families, and their communities against the infection.

The International Health Regulations (2005) list risk communication as one of eight core capacities that WHO

Member States need to build and sustain as part of a global agreement to strengthen national and global systems to detect and respond to public health threats [5]. Risk communication and community engagement (RCCE) is an important approach for developing and disseminating accurate information, and it has been associated with more successful empowerment of affected local communities in disease outbreaks [6]. Experiences from the HIV, Ebola, Zika, and polio epidemics have demonstrated the cost to public health and health systems when rumors and misinformation are amplified in an environment where there is already a high level of distrust, which is aggravated by a poor public health communications response [7]. In a public health emergency or outbreak, existing service delivery may be disrupted and health authorities may not yet know the facts and have adequate evidence; this can lead to an information void, causing confusion and anxiety in the affected population [8]. If information voids are not responded to with high-quality health information, they can quickly be filled with misinformation and disinformation. Pieces of information of unknown validity can be benign and transient, or they can be false, causing damage if they affect individual and community decision-making. Rumors can be detrimental to health, especially in emergencies and crisis situations [4]. Rumors, unlike misinformation or disinformation, may be found to be true, and they can be either persistent and long-standing or evolve quickly after an acute event [4].

Overall, health emergencies give rise to information overload, which has been shown to influence people's behavior, risk perception, and protective actions during health emergencies

[9] and subsequently give rise to information avoidance. In emergencies, affected individuals and populations may have difficulty processing complex information and may retain only some of the early information they receive. In such circumstances, rumors can propagate quickly, challenging emergency responses that rely on the affected population to follow accurate health advice and enacting behaviors to protect individual and community health [8].

Although rumors and health misinformation have been around for as long as diseases, today's environment is different. The COVID-19 infodemic has been an unprecedented challenge because we are experiencing an epidemic in a digitized globalized society. Digital tools and technologies have not only changed the way we communicate but have also changed our lives, altering the way we live, work, interact, and build our social identities and sense of community. For example, rumors and information have travelled across borders very quickly and influenced traditional media news cycles and coverage, emotive misinformation travels much more quickly across the digital media than fact-based health information, and epidemic control decisions or controversy in one country can cause debate and comparison with responses in other countries [9].

This infodemic has placed strain not just on how to communicate the evolving scientific knowledge but also on how public health authorities can implement a nimbler pandemic response that addresses the needs and concerns of local communities. During the COVID-19 response, health authorities have faced full-on the changed information and communication ecosystem [10] and its challenges, such as:

- Computational amplification of polarizing messages over factual ones, and use of bots and cyborgs to manipulate the outcome of online petitions, change search engine results, and boost certain messages on social media;
- Widespread microtargeting of social media users that is enabled by the social media and search engine platform business models, putting individuals into their own personalized "information bubbles";
- Changed practices in TV and radio newsrooms that enable dissemination and amplification of poor-quality information that originates online;
- Weakened local media and collapse of local journalism, which have enabled mis- and disinformation to take hold.

In response to the infodemic, health authorities have needed to build partnerships beyond their usual networks—with fact-checkers; broader groups of media and journalists; social media, search engines, and digital interaction platforms; community organizations; civil society; and others. However, there is still room for improvement based on experience from the COVID-19 response. For example, although fact-checking organizations are relatively mature worldwide, half of them do not work with health professionals when fact-checking and debunking health-related claims, leaving room for better collaboration with health authorities and medical associations [11]. Moreover, whereas communication campaigns can raise the visibility of a set of messages, they are often not effective at debunking false claims, which require more quantitative and qualitative pretesting of messages; also, they must respond to questions, concerns, and narratives that are currently capturing people's attention in a specific geographical area or a vulnerable community [12]. Mis-, dis-, and malinformation (also referred to as *information disorder*) are major and growing challenges, not only for emergency response but also for other societal actions [10].

Because of these challenges, the COVID-19 infodemic is not only a communication challenge but a challenge for the whole information ecosystem. Already at the beginning of the COVID-19 pandemic in April 2020, the WHO had crowdsourced a framework for managing infodemics that calls for whole-of-society involvement and response [3]. This framework recognized that in digitized society, the harmful effects of the infodemic cannot be managed through the prevailing approaches to communication, community engagement, and messaging alone. Infodemic response must take into account the information ecosystem, the ways we interact within the information ecosystem, and how information affects our health behavior. Consequently, this dynamic environment requires interventions across multiple levels, such as individual, community, medium, platform, policy, and others. The WHO infodemic management framework called for a multidisciplinary research agenda that informs the use of evidence-based interventions and surveillance across all phases of an epidemic [13], which led to the convening of this technical conference.

Between June and October 2020, the WHO Information Network for Epidemics (EPI-WIN) organized a global online technical conference followed by an asynchronous expert review exercise to develop a public health research agenda for infodemic management [3,13-17]. This transdisciplinary scientific consultation and review gathered infodemic insights and approaches from a wide range of relevant fields to inform and expand frameworks in infodemiology. Along with strengthening the foundations of an expanding infodemiology discipline [18] and creating the research agenda to direct focus and investment toward this emerging field, other aims of the conference were to improve understanding of the multidisciplinary nature of infodemic management; identify current examples and tools to understand, measure, and control infodemics; and establish a community of practice and research, preparing the ground for sustainable, long-term practices for responding to infodemics. The full conference report is available on the WHO website [17]. This paper summarizes the methods and results of the research agenda and the development of the research questions.

## Methods

### Overview

The research question prioritization exercise was designed in line with the WHO research agenda development guide for staff [19]. Held in the context of the COVID-19 pandemic and with travel restrictions in place, the consultation necessarily took place online via videoconference. The virtual discussions took place over 8 meeting days during 4 weeks in June and July 2020, and they resulted in a research agenda frame and a list of priority research questions. This was followed by asynchronous email communication from August to October 2020, during which

participants were led through a structured expert opinion exercise to review and prioritize research questions within the set research agenda frame. Institutional Review Board review was not sought because the work described in this paper was based on observation of discussions at the conference, and it focused on the synthesis of expert opinion following the Chatham House Rule [20]. No personal information was collected from the participating experts.

## Format of the Virtual Conference

The 110 invited participants represented over 35 countries across 19 time zones, with a 56% to 44% gender split in favor of women (62 female, 48 male). They were academics selected by the organizers for the relevance of their publication record in the past two years for the purpose of this consultation, or practitioners who were working in pandemic response. A total of 60 additional invited academics were not available to participate. The conference participants represented 20 different academic and professional fields, such as digital health, computer science, communications and graphic design, media studies and journalism, history, applied mathematics, information science, data science and computational social sciences, complexity science, social and behavioral sciences, ethics, governance, marketing, and user experience and design; they were joined by colleagues from the fields of risk communication and community engagement, epidemiology, and public health, as well as by global public health implementing partners. Conflicts of interest were reviewed in accordance with WHO procedures for the management of declaration of interest for expert consultations [21]. The conference and follow-up communication were supported by a team of 49 organizers.

The meetings took the format of plenary sessions at the beginning and end of the conference and an in-between working session with four discussion sprints. Each participant was engaged in the meeting process for 18 hours (10 hours in plenary and 8 hours in topic discussions). Participants were split into four teams, grouping by similar time zone location but ensuring academic and practitioner diversity of the teams. Each team met four times for 2-hour "sprint sessions" of intense discussion on one of four topics, led by dedicated "topic masters" (scientific facilitators). The topic masters were scientists established in their scientific disciplines; 7 were academics employed by universities, and 1 was a WHO staff member with an academic affiliation. As the teams rotated from topic to topic, the topic masters facilitated discussions to collect insights from the discussion and validate expert opinion they had collected from discussion with preceding teams. By the end of the process, each team had discussed each topic, and each topic was discussed with four teams in an additive fashion—a total of 32 sprint hours of expert discussion.

The discussion sprints were oriented around four topics that mirror the epidemiological method for outbreak detection and management across the phases of the epidemic curve, enabling the actions of "preparing, monitoring, detecting, intervening, strengthening, and enabling" infodemic management. The topics were (1) how to measure and monitor digital and physical information environments; (2) how information originates and spreads; (3) how information affects individuals and populations; and (4) what interventions work to protect and mitigate against mis- and disinformation. By the end of the working session, a frame for a research agenda emerged based on the feedback from all the team discussions, seeded with draft research questions that were identified by the discussion facilitators.

In addition, the facilitator leaders of each of the four discussion streams at the conference summarized the discussions they had with all four teams of participants. Their reports summarized discussions about the main suggested research questions for the research agenda as well as enablers and challenges to researching them. This initial collected set of research questions became the basis for the follow-up process after the conference.

## Asynchronous Expert Ideation and Prioritization Exercise

After the virtual conference, the same participants were led through a 3-month asynchronous structured exercise that aimed to collect and rank research questions and to guide the participants toward a refined research agenda. In the exercise, structured expert judgment was collected through an adapted Delphi consultation using the Investigate Discuss Estimate Aggregate (IDEA) protocol [22]. The method involved asking the participants to devise and submit research questions that were relevant to the topic, answerable in the short or medium term, ideally capable of producing knowledge that could be put to use in the short or medium term, and focused on scope (ie, an answer to the research question should be provided in a single academic paper). They were also asked to focus on what would be scientifically feasible to answer and what had an expected public health benefit. To improve the reach beyond the pool of conference participants, each expert could invite up to two additional experts, based on their expertise and the value of their potential contributions. In total, 38 experts submitted additional research questions to the pool of candidate research questions and the following ranking survey. To maximize transparency in the categorization, experts could themselves choose which category or subcategory to submit a research question to. To identify potential gaps in the overall research agenda, the survey included open-ended questions.

A list of candidate research questions was built by combining the questions that were proposed by the topic facilitators based on the discussions at the conference and those that were collected through the survey round after the conference. The collected candidate research questions were assessed for topic overlap and scope, and they were edited and merged for clarity by three reviewers. The three experts were present in the discussions during the conference and are coauthors of this paper. Two are staff members of health authorities, and one is an academic. This reduced the questions to a consolidated list that was used in the research question ranking exercise.

The questions in the consolidated list were then anonymously scored and ranked through another exercise using the LimeSurvey platform [23]. There, participating experts were asked to rank the research questions based on two dimensions: public health impact and feasibility. These two ranking indicators were selected to point the agenda to evidence that can inform COVID-19 infodemic response quickly or with high
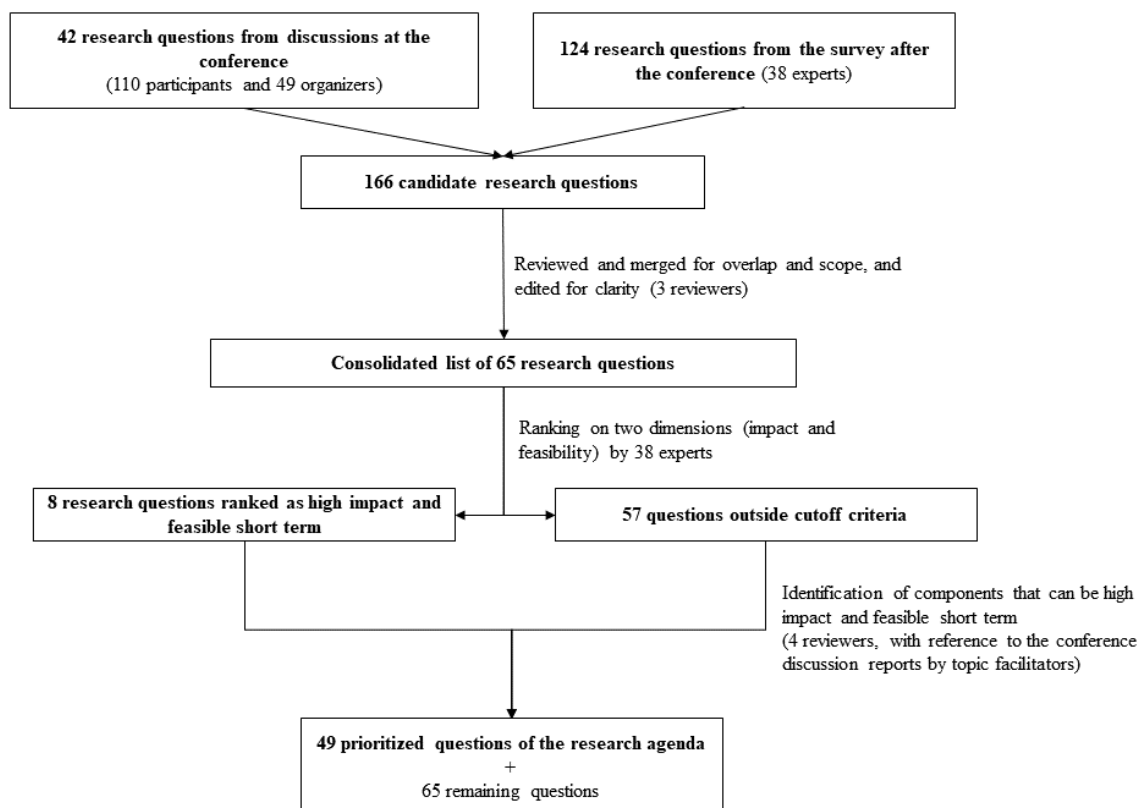
impact, anticipating its importance in light of pandemic fatigue and the protracted use of public health and social measures to manage the pandemic, as well as ahead of the eventual introduction of COVID-19 vaccines. Public health impact was assessed through the question: Can this research lead individuals or communities to take healthy actions or help them understand why and how they do not take healthy actions? Research questions that could lead to meaningful change or adaptation of behaviors would be considered more impactful. Experts were asked to rate each question on a 5-point Likert scale (1, very low impact; 2, minor impact; 3, moderate impact; 4, high impact; 5, very high impact). Feasibility was assessed through answering the question "Can you think of a research project that would answer this specific question in a set number of months?" The faster the research project could be initiated and deliver results, the higher its feasibility and usefulness for the COVID-19 pandemic response. Experts were asked to rate this question on a 5-point Likert scale (1, 3 months; 2, 6 months; 3, 12 months; 4, 18 months; 5, 24+ months, based on emergency response planning time periods). A research question was considered high priority when it scored above 3 on impact and below 3 on feasibility.

To reduce potential survey fatigue and to avoid systematic missingness in the rankings (ie, due to respondents ranking only the first few questions within each stream), the order of the research questions to rank was randomized within each research stream. The randomization was conducted via the LimeSurvey platform. Following the ranking exercise, four experts reviewed the questions that fell outside the prioritization area—below the 3.5 consensus impact rank and with feasibility of >1 year. The four experts were three researchers that had previously reviewed the submitted research questions, with an additional staff member of a health authority who was not a coauthor on this paper. The fourth health authority staff member was added because the research agenda questions were meant to be feasible in a short time frame or highly important to the health authority response to the infodemic. The experts reviewed the questions outside the cutoff and, on consensus, they identified prerequisites or parts of these research questions that could be delivered with quicker feasibility and high public health impact. These research components were added to the research agenda as research questions.

This exercise reduced the number of questions to a shortlist of top priority and second tier priority questions per work steam, totaling 49 priority research questions. The remaining questions that were part of the exercise and did not make the prioritization cutoff were retained for reference, and they can be used for future reviews. The results of the recursive refinement of research questions through structured expert judgment exercise are summarized in Figure 1.

**Figure 1.** Refinement of the research questions through the structured expert judgment process.

## Results

### Themes That Emerged During the Discussion Sprints

The discussion at the virtual conference reflected the complexity of the information ecosystem and the way it influences the strategies for managing the COVID-19 infodemic and other infodemics to support health behaviors and the management of epidemic risk. Several themes surfaced in the topic discussion sprints, as follows:

A common theme across discussions was that it is necessary to identify reproducible patterns and crossdisciplinary metrics for the science of infodemiology. Because access to full data sets from social media is rare and they do not represent the engagement of all populations, and because metrics vary from platform to platform, it is difficult to produce generalizable or comparable results. Mathematical modeling, such as epidemiological modeling, does not necessarily take human behavior into account, which can limit its efficacy to predict future human behavior and the impact thereof on an outbreak; however, modeling can aid the development of hypotheses for how information/infection flows, how networks might respond, and how interventions should be designed to test them. There are also limits to applying the epidemiological framework as a way to monitor and measure spread, especially if we assume that the unit we are working with is information instead of a virus, because viruses do not have an agenda and they infect opportunistically. Detangling the differences between rumors, misinformation, and disinformation requires a common taxonomy of information classification, some of which may be labelled as more harmful or less harmful. This could inform identification of the "tipping points" or when action needs to be taken to address more harmful misinformation by offering a more tailored and effective response.

Although it is important to describe the flow of health information, there needs to be a balance between a system-level understanding that "washes over" details and a case-study understanding that captures details but may miss the "bigger picture." Substantial amounts of social and behavioral and health data are available; however, determining which data sources and types of analyses would improve a response needs a clearer definition. The degree of detail is needed to understand the infodemic while balancing privacy and ethical concerns, and managing limited analytic capacity in short time frames should be discussed. Amid a pandemic, speed is of the essence, and balancing rapid data collection and analysis methods with the desire for rigor may mean prioritizing specific kinds of data for short-term operational use versus longer-term, longitudinal trend analysis and use. Understanding the diffusion of information through certain networks may require other data collection approaches and discussion of how closed messaging apps and offline networks challenge this.

One area of research that needs further study is the extent to which offline behavior is being influenced by online behavior (and vice versa). There is limited research on how exposure to information or misinformation affects behavior because behavioral processes can be quite complex. Amid a crisis, people might use cognitive shortcuts and rely on the first information

they hear; also, they may be less adept at processing more complex information. At the same time, there is little known about the longitudinal effects of the exposure to false claims that may not seem harmful at any one point in time but could have a cumulative harmful effect over time. In addition, when misinformation is easy to spread, this can create a harmful mixture. Anecdotal evidence suggests that people can exhibit negative health behaviors because of misinformation they heard during the COVID-19 outbreak; however, we need better measures of how knowledge connects to intent and behavior, both online and offline. For example, does increased exposure to misinformation make it more likely that someone will exhibit a behavior that is detrimental to their health? Further research is needed to develop better monitoring metrics, in addition to consolidated and validated indicators that predict behaviors or serve as proxies for specific behaviors.

The participants also emphasized that there is an interplay between information ecosystem actors and the resilience of communities and individuals. It was agreed that trust is a key element of building resilient communities. This leads to the need to establish and maintain trustworthy information sources. Some work must be done to identify these sources of information and to ensure easy and equal access. The discussions also highlighted the urgent need to empower communities to manage infodemics and build resilient communities through co-designed interventions. This would be made possible by understanding the context in which infodemics occur and spray. Community engagement goes along with building self-efficacy and self-capability through practice. It should focus on the "middle ground," as in, the majority of "silent lurkers"—those who have not yet formed strong opinions. Besides individuals, communities, and states, members of the private sector should be regarded as actors. Internet platforms can be active vectors or targets of campaigns and can also be influential members of communities.

When considering long-term interventions, critical thinking and literacy (eg, health, information, digital, and media literacies) play important roles as a basis for interventions to address infodemics. Health literacy is a major topic in health communication research and practice. It includes critical literacy as the ability to evaluate and apply health information, and it is considered a major asset in managing an infodemic. Similarly, information, news, digital, and media literacies contribute to individuals' ability to distinguish high- from low-quality information, especially online, and to the ability to improve their offline lives through digital technology use. Research into each type of literacy has developed in isolation, and questions remain on how to empower populations to think critically; what normative models of thinking are most appropriate for an infodemic; who is responsible for building literacy; and how literacy efforts can be integrated into existing societal systems (eg, school education) and be adapted to reach populations outside of the traditional educational settings.

To help prioritize interventions and actions, it is also necessary to identify priority populations based on key vulnerabilities. Population studies need to be conducted to identify specific individuals and groups of individuals who are at the greatest risk of not being able to critically assess misinformation and of

spreading it. This approach should include studying people's perceptions, beliefs, and knowledge, as well as the barriers and facilitators that can affect the access to and evaluation of credible health information as well as its use in offline life. Additionally, the alignment of information vulnerabilities with disease vulnerabilities should be considered.
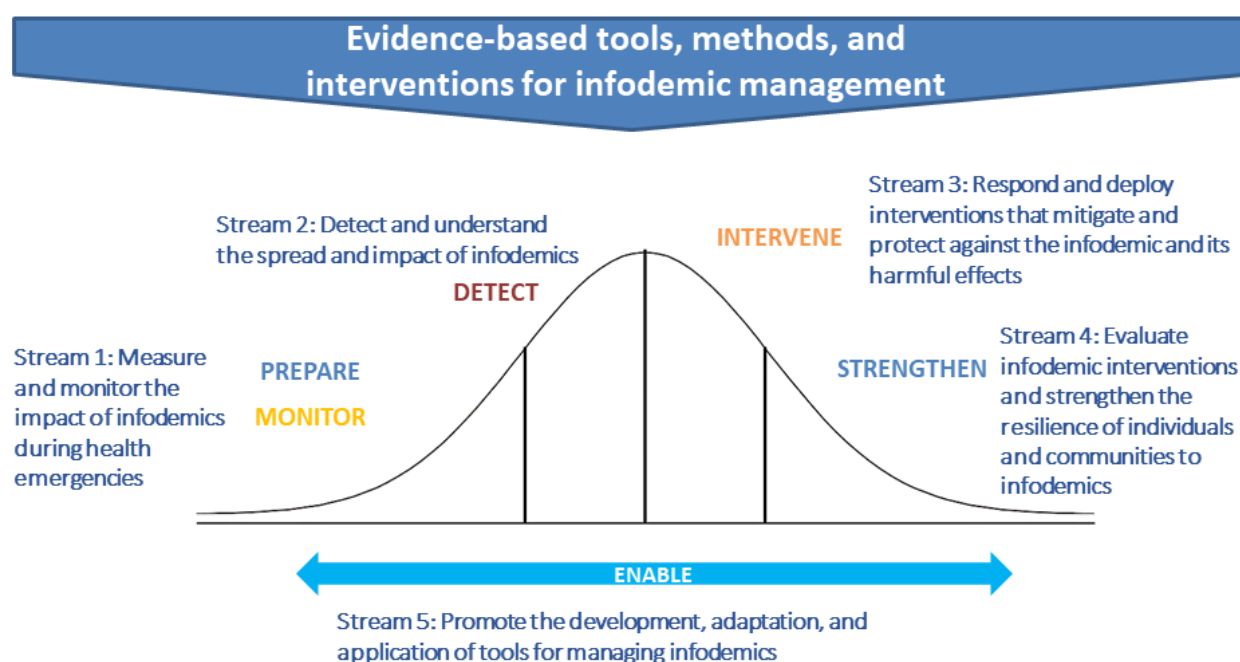
## Public Health Research Agenda for Managing Infodemics

In addition to reviewing the current evidence and research gaps across different scientific fields, the conference participants sought to identify a research frame that could structure a public health research agenda for infodemic management. The aforementioned themes that emerged converged to broader landscaping of research gaps (eg, the need for better monitoring and metrics; localized and system-level characterization of infodemics; and understanding the components of the information ecosystem, individuals, communities, states, and private social media platforms). The themes also focused on some specific areas of knowledge gaps or promising interventions (eg, understanding the linkage between online and offline behavior; the roles of critical thinking and health literacy; and identification of priority populations). An agenda for future research should not only aim to fill gaps in the existing evidence base but, at high priority, also to link research evidence to public health practice. Therefore, the conference participants agreed to establish the research agenda frame analogous to the lens of epidemic management and to fit the identified themes, issues, and gaps to this epidemiological frame (see Figure 2). The framework's streams were built on the activities of a health authority that supports outbreak response along the phases of an epidemic curve [24]—prepare and monitor, detect, intervene, strengthen, and enable infodemic response, as outlined below and in Table 1:

- Stream 1 supports the preparedness and monitoring of infodemics through measurement and monitoring of impact of infodemics. Standardized metrics and measurement tools can help characterize infodemics online and offline, identify absence of information where misinformation can gain

more traction, and help recognize tipping points when detailed investigations need to take place. Last, evaluation of infodemic management interventions needs more elaboration.

- Stream 2 addresses the need to detect and understand the spread and impact of infodemics. In the context of infodemics, communities and vulnerable groups are no longer defined only in terms of geographies but can also be formed through shared values, goals, or motivations. Development of interventions therefore needs localized contextualized understanding of the infodemic, how misinformation affects behaviors in vulnerable groups, and understanding of the ethical and regulatory approaches needed to mitigate the spread of misinformation.

- Stream 3 addresses the response and deployment of interventions that protect and mitigate the infodemic and its harmful effects. Thinking about implementation of interventions needs to be built into the infodemic management activities and research so that the research is linked to what health authorities need to respond. To achieve this, behavior/change models relevant to infodemic management need to be developed, and interventions need to be designed.

- Stream 4 aims at research that strengthens infodemic management by development of common frames to improve intervention development and programmatic response to infodemics. Using the continuum of community engagement, local cultural context, and building resilience to infodemics and misinformation at individual, community, platform, and societal levels are addressed.

- Stream 5 supports the overarching aim to strengthen infodemic management practice by enhancing transferability of lessons learned and evidence-based interventions between contexts, countries, and infodemics. The information ecosystem and socioeconomic determinants of access and use of health information differ across countries; we therefore need to understand how interventions can be successfully transferred across countries and what impact they will have in other settings.

**Figure 2.** The frame of the research agenda mapped onto the phases of epidemic preparedness and response.



**Table 1.** Framework of the public health research agenda for managing infodemics.

| Stream | Subtopics |
|---|---|
| Stream 1: Measure and monitor the impact of infodemics during health emergencies | 1.1. Standardize taxonomies and classifications |
| | 1.2. Develop new metrics to measure and quantify infodemics |
| | 1.3. Analyze and triangulate data from multiple sources |
| | 1.4. Improve evaluation approaches for infodemic interventions |
| Stream 2: Detect and understand the spread and impact of infodemics | 2.1. Understand how information originates, evolves, and spreads on different platforms and channels |
| | 2.2. Assess the role of actors, influencers, platforms, and channels |
| | 2.3. Understand how misinformation affects behavior in different populations |
| | 2.4. Develop regulatory and ethical principles to mitigate the spread and propagation of harmful health information |
| Stream 3: Respond and deploy interventions that mitigate and protect against the infodemic and its harmful effects | 3.1 Design a behavioral/change model applicable to infodemic management |
| | 3.2. Design interventions for different levels of action to mitigate the infodemics |
| Stream 4: Evaluate infodemic interventions and strengthen the resilience of individuals and communities to infodemics | 4.1. Develop interventions that address individual, community, cultural and societal-level factors affecting trust and resilience to misinformation |
| | 4.2. Understand and learn from how misinformation has affected behavior among different populations and in different contexts for specific infodemics |
| | 4.3. Identify factors associated with successful infodemic management by health authorities, the media, civil society, the private sector, and other stakeholders |
| Stream 5: Promote the development, adaptation, and application of tools for managing infodemics | 5.1. Use implementation research evidence in program improvement and policy development |
| | 5.2. Promote evidence-based interventions and approaches among countries |
| | 5.3. Improve effectiveness and response times to the infodemic during acute health events |

At the conclusion of the conference, 42 research questions were collected from the topic discussions, as curated by the scientific topic facilitators. During the follow-up research question generation exercise, 38 experts submitted an additional 124 research questions across 5 research streams and 16 subcategories. This added up to 166 candidate research questions. These research questions were reviewed and merged for repetition, overlap, and scope; they were then edited for clarity by three reviewers. Suggestions that were not formulated as research questions were excluded. This review identified a high degree of overlap and repetition, pointing to a saturation of topics submitted for the ranking exercise. It resulted in a consolidated list of 65 questions that were subjected to the ranking exercise.

The research questions to be ranked were evenly distributed, with at least 10 questions included for ranking in each of the five research streams (18 in stream 1, 16 in stream 2, 10 each in streams 3 and 4, and 11 in stream 5). The ranking exercise results for these questions are depicted in Figure 3.

Following the ranking exercise, four experts reviewed the results. Based on the ranking exercise, only 8 research questions covering streams 1, 2, and 3 were prioritized with a consensus rank greater than 3.5 and feasibility of <1 year. Therefore, the experts reviewed the 8 ranked questions and the remaining 57 questions that fell outside the cutoff limit. Based on their expert judgment and guidance from the reports of topic discussions during the conference, they identified precursor questions or components of these research questions that could be delivered with quicker feasibility or higher public health impact. The experts worked with the goal to use the questions and the

feedback collected in the ranking exercise and used them as a guide to formulate research questions that could be their precursors. They worked on consensus and formulated the final list of 49 research questions, and they retained the additional 65 questions for future reference.

Expert review of the results of the ranking exercise identified 3 top research questions per work stream, resulting in a list of 15 top priority research questions for the public health research agenda for infodemic management (Table 2). Further, a second tier of important research questions was set for each subtopic, totaling 34 questions. Multimedia Appendix 1 shows the prioritized research questions and agenda in more detail, as well as the additional 65 research questions that were not identified as a high priority in the short term. These can be used to map further evidence gaps on the topics and for reference and guidance in subsequent research agenda reviews [17].

**Figure 3.** Ranking of the surveyed research questions across two indicators: public health impact and feasibility. Research questions that were within the cutoff limit of minimum 3.5 impact and less than 12 months feasibility are marked in yellow. Questions that were ranked outside the cutoff limits were reviewed and broken into additional smaller component questions that were of high value.

**Table 2.** Top 15 research questions across five streams of the research agenda.

| Stream | Top 3 questions per stream |
|---|---|
| Stream 1: Measure and monitor the impact of infodemics during health emergencies | • What are ways to score health-related misinformation according to its potential for harm (to people's health and behaviors; social cohesion; trust in health service delivery, government, communities, media; etc)?<br>• How do the infodemic curve and measures of spread and impact change over time during the phases of a disease outbreak?<br>• What are the potential indicators or their proxies for measuring trust, resilience, behavior change, exposure to misinformation, susceptibility to misinformation, social cohesion, depth of community engagement, etc? |
| Stream 2: Detect and understand the spread and impact of infodemics | • How does misinformation mutate, adapt, or become remixed between infodemics and within infodemics?<br>• What are the strategies used to reduce misinformation's potential harmfulness in closed networks (online and offline)?<br>• How do different types of health misinformation affect online and offline behavior, and what are some measures that can help forecast the impact of the health misinformation types on behavior? |
| Stream 3: Respond and deploy interventions that mitigate and protect against the infodemic and its harmful effects | • What behavioral or process models can inform the development of an infodemic strategy and measure its impact at the individual, community, platform, or societal level?<br>• What are the promising interventions at the societal/community/individual/health system levels to address and mitigate health misinformation?<br>• What types of participatory or human-centered design approaches can be used to produce more tailored and effective infodemic management interventions? |
| Stream 4: Evaluate infodemic interventions and strengthen the resilience of individuals and communities to infodemics | • How might we define and measure the gradient of community engagement, trust, and empowerment at the individual and community levels as they relate to infodemic management and reduction of harm from health misinformation?<br>• What are the sociobehavioral, mental heuristics, and design hierarchies that need to be considered when developing an intervention at the individual and community level?<br>• What are the "best buy interventions" to be used by different types of actors in society to maximize the impact on the infodemic at a lower marginal cost? |
| Stream 5: Promote the development, adaptation, and application of tools for managing infodemics | • What considerations should be included in the assessment of risk, harms, and opportunities during the design and implementation of research and infodemic management interventions?<br>• What would a readiness assessment look like for infodemic preparedness for a new COVID-19 health intervention?<br>• What recommendations can be made to update the International Health Regulations to incorporate infodemic management more strongly as a core capacity of Member States? |

## *Discussion*

### Principal Findings

Throughout the consultation, the discussions built progressively; participants shared a wealth of experience, discussed the challenges and benefits of various approaches, clarified the initial topics, and ultimately achieved a high degree of consensus about the needs that the research agenda would have to meet. The overarching conclusion was the need to complete and implement the research agenda along with the framework for action [3,13]. The takeaway action points from the conference are as follows.

Information, misinformation, and public health are intertwined by nature: the WHO has dealt with issues at the intersection of misinformation, trust, and demand for health services since it was founded. Lessons from this experience have led to evolved epidemic response methods, tools, and the global response community over time. The WHO and other partners who work in the fields of public health communication, risk communication, and community engagement have been challenged by the scale of the COVID-19 infodemic, which has been amplified by the global digitized information ecosystem.

In a new century, addressing new types of outbreaks requires innovative and precise public health tools [25]. Different populations have different information needs, channels, and barriers. Evidence-based interventions are needed at all levels—for individuals, communities, platforms, health systems, and societies—to reduce the transmission and impact of the disease. Coordination, connection, and integration across disciplines and sectors must be central to expanding the scientific discipline of infodemiology. Rapid application of the science during the COVID-19 pandemic needs (1) sustained integration across the various disciplines of research; (2) integration between research, practice, and lived experience; and (3) inclusion of representation and voices from different sociocultural contexts in practice and lived experience.

At the same time, infodemic management must broaden its tools beyond communication and consider all components of the information ecosystem [7,26,27]. Because the information ecosystem spans both online and offline environments, it is more difficult to detect and respond to the infodemic in communities as well as to work proactively to build resilience and a healthier information ecosystem overall. Media, policy makers, and the private sector influence the information ecosystem where individuals, interest groups, civil society,

academia, fact-checkers, and others also interact. Partnerships between health authorities, fact-checkers, media organizations, and other global public health partners, such as the Africa Infodemic Response Alliance [28] are critical to effectively promoting high quality health information and countering health misinformation at local level. The RCCE collective service [29] was started in June 2020 to concentrate the RCCE capacities across global RCCE partners. Strengthened partnerships at local levels are also needed to focus on community engagement in offline communities. On the other hand, regulatory interventions could help standardize access to social platform data, ensuring that the data we do have access to is comprehensive and regular. Access to regular and better data/metadata would increase accountability for how a healthier online information ecosystem is built. Based on data availability, this access could also facilitate the design of research and interventions to give us a better understanding of which interventions work online and the conducting of independent analyses of information provided by the platforms.

Addressing the harms of infodemics is important because they impact health behaviors and are barriers to healthy life and well-being. It is important to better understand proactive strategies that apply social inoculation theory or literacies theory in building resilience. Developing health literacy is critical and includes access to health services literacy, and it is dependent on digital nativity/technological skills, access to information, and media literacy/interrogative skills. The reasons why mis- and disinformation spread are complex; therefore, it is important not to reduce that complexity by framing infodemic management as simply a battle against misinformation [7,30-34]. It is equally important to reinforce and accelerate health-enhancing behaviors and generate information to help people develop resilience to information disorder. In the long term, this will help people build trust, make informed decisions, and access essential health services, and it will have impact far beyond the COVID-19 pandemic.

Given the urgency of pandemic response, the new transdisciplinary practice will have to learn from practice and iteration even as it develops, reporting experience gained through implementation to provide more evidence on what works and what does not [35]. Ultimately, health authorities need to identify and allocate the necessary capacity to manage infodemics. This is a programmatic and process issue. Once that capacity is in place, decision-makers and the private sector need to develop, validate, implement, and adapt tools for infodemic management during acute public health events in culturally and contextually appropriate ways. The issue of connecting this evolving practice and research is not trivial: the community that implements the research agenda must be, and must remain, a community of practice *and* research that prioritizes questions to inform operations and improve contemporary practice, foregrounding the pragmatic needs of people in the field and on the ground.

We also need to think about how to build systems for social listening, signal detection, and the analysis of infodemics and misinformation. For example, investment is needed to develop a shared, open reference database for characterizing misinformation (including examples) to identify appropriate

interventions and when and how to deploy them [36]. Effectively, the content would be re-contextualized to enable characterization and use in the analysis. This database could include different types and sources of misinformation, the intent of those creating or sharing misinformation, the degree of inaccuracy (based on the level of expert consensus and scientific evidence that exists), its impact on attitudes or behaviors, the likely audience, its virality, or its alignment with politics. This reference database could be populated with specific examples of misinformation that fall into each domain, which could then be aligned with interventions based on best practices shown to be effective for that type of misinformation. Such a reference database would help to answer questions about the differences between the content people will merely share online and the content that will affect their decision-making and offline behavior. It would also aid the investigation of whether we can use content characteristics to predict the likelihood of spreading a piece of content in different ways.

The community of research and practice could also develop and use a shared "living systematic review" for interventions measured in terms of effectiveness on a set range of criteria, strength of evidence, generalizability, and likely contexts for application. Interventions across disciplines could be collected, with a rubric describing the outcomes against which the intervention has been tested, its generalizability or application to specific populations, the contexts in which it has been tested, its feasibility and costs, and the confidence in the findings. This could include determining the consistent metrics appropriate to evaluating the success of an intervention to prioritize efforts. However, it is unlikely that only one intervention will be successful; a toolkit of different approaches will likely be appropriate. This living systematic review could be aligned with the misinformation reference database to identify gaps where good evidence-based research does not exist to address certain types of misinformation.

## Conclusions

The resulting public health research agenda for infodemic management will be maintained on the WHO website as a living document; its implementation and priorities will be reviewed and adjusted regularly.

Infodemics impact people, including health professionals, globally. Although infodemics are not new, addressing them in the new digitized society is a different and centrally important challenge in responding to the COVID-19 pandemic as well as future pandemics. The research agenda that emerged from this consultation crystallizes themes that can inform initiatives to build the foundations of effective infodemic management in all countries. The main target audience for these research questions are researchers and practitioners. They will also be of interest to public health experts, nongovernmental organizations, the media, and other stakeholders.

There is a large gap between infodemiology research and evidence that has been generated by the academic disciplines and the response to the infodemics. Tools and interventions that are grounded in this evidence are sorely needed by health authorities worldwide. This is partially because scientific disciplines have worked in a mostly disconnected fashion on

addressing the challenge of information overload, communication, design, media studies, sociobehavioral factors, misinformation, and the ethics and regulation of the information ecosystem. The WHO and its Member States and partners must close this gap by developing and adopting evidence-based tools that are appropriate for their local contexts. This consultation and the previous infodemic management meeting [3,13] may have been among the first opportunities for many people working toward this goal to hear about the expertise and activities of others, and to frame the entirety of this activity within the problems of disease control and public health.

Following the conference, the WHO partnered with five scientific journals in a joint call for papers for special issues on infodemiology [37], two of which have already been published [38,39]. The WHO EPI-WIN team has used the outcomes of this conference as the input in the third and fourth WHO infodemic management conferences [40,41] and the upcoming fifth WHO infodemic management conference, which will focus on the development of measurements and metrics for infodemic management. The research gaps that were identified have also guided the WHO in the review of the COVID-19 research

blueprint [42] and in the development of partnerships that foster filling of research gaps and for translation of evidence into use by health authorities and other partners [28,43]. The WHO has also applied evidence and infoveillance methods to inform its own work and contribute to the development of metrics for health authorities [18,44,45].

The challenge of a novel pandemic pathogen intertwined with an infodemic is a double burden that demands action-oriented research to inform public health response. The new research agenda will strengthen the scientific understanding of how infodemics impact populations and their health, but it will also serve as a basis for action and learning for future preparedness, strengthened through cross-sectoral pilot projects and continuous after-action reviews to build capacity. After the acute phase of the COVID-19 pandemic, we need to shift the focus to strengthening longer-term capacities and advocating for the inclusion of new tools and indicators. When applied to acute health events, the evolving research discipline of infodemiology can provide crucial evidence and facilitate multidisciplinary expertise and coordination.

## Conflicts of Interest

AA, EA, JNA, SB, SCB, CC, JL, RL, AM, PNN, TN, TDP, AR, and BY are staff of the World Health Organization; NA, AI, AK, SK, PP, DP, and EW are staff members of the US Centers for Disease Control and Prevention; LE and AW are staff members of the European Centre for Disease Prevention and Control. RH is staff of UNICEF. These authors alone are responsible for the views expressed in this paper, and they do not represent the views of their organizations. BS participated in the conference while being staff of WHO. SB consults for WHO in data analysis for emergency preparedness and supported the work described in the paper during his consultancy contract. IB is Director of WHO Collaborating Center on information systems for health, which supports WHO with broader digital health analytics and policy analysis. The center has supported Pan American Health Organization (PAHO)/WHO with infodemic analytics during COVID-19. LMB works for Immunize Canada/Canadian Public Health Association, which has received educational grants/funding from Merck Canada, Pfizer Canada, Pfizer Global, Moderna Canada, Seqirus, Sanofi Canada, GSK Canada and the Public Health Agency of Canada (PHAC). These funds are not related to the paper. AB has received funding from Carlsberg Foundation for a research project about online hostility for which she is project co-primary investigator. This project is not related to the deliberation described in this publication.

Multimedia Appendix 1
Public health research agenda for managing infodemics (all prioritized and collected research questions).
[DOCX File , 37 KB - infodemiology_v1i1e30979_app1.docx ]

## References

1.  COVID-19 dashboard. World Health Organization. URL: https://covid19.who.int/ [accessed 2021-04-19]
2.  WHO Director-General's speech at the Munich Security Conference, 15 February 2020. World Health Organization. URL: https://www.who.int/dg/speeches/detail/munich-security-conference [accessed 2021-04-29]
3.  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]

XSL•FO
RenderX

4.   Kou Y, Gui X, Chen Y, Pine K. Conspiracy talk on social media: collective sensemaking during a public health crisis. Proc ACM Hum-Comput Interact 2017 Dec 06;1(CSCW):1-21. [doi: 10.1145/3134696]

5.   International Health Regulations (2005) Second Edition. World Health Organization. 2008 Jan 01. URL: https://www.who.int/publications/i/item/9789241580410 [accessed 2021-09-07]

6.   Walker B, Adukwu E. The 2013-2016 Ebola epidemic: evaluating communication strategies between two affected countries in West Africa. Eur J Public Health 2020 Feb 01;30(1):118-124. [doi: 10.1093/eurpub/ckz104] [Medline: 31177274]

7.   Vicol D, Tannous N, Belesiotis P, Tchakerian N, Stewart R. Health Misinformation in Africa, Latin America and the UK: Impacts and Possible Solutions. Full Fact. URL: https://fullfact.org/media/uploads/en-tackling-health-misinfo.pdf [accessed 2021-09-07]

8.   Reynolds B, W Seeger M. Crisis and emergency risk communication as an integrative model. J Health Commun 2005;10(1):43-55. [doi: 10.1080/10810730590904571] [Medline: 15764443]

9.   Brennen J, Simon F, Howard P, Nielsen R. Types, sources, and claims of COVID-19 misinformation. Reuters Institute at University of Oxford. 2020. URL: https://reutersinstitute.politics.ox.ac.uk/types-sources-and-claims-covid-19-misinformation [accessed 2021-09-07]

10.  Wardle C, Derakhshan H. Information Disorder: toward an interdisciplinary framework for research and policy making. Council of Europe. 2017. URL: https://rm.coe.int/information-disorder-report-version-august-2018/16808c9c77 [accessed 2021-09-07]

11.  Health authorities and innovative collaborations across society to combat the infodemic. World Health Organization. 2020 Dec 04. URL: https://www.who.int/news-room/events/detail/2020/12/04/default-calendar/health-authorities-and-innovative-collaborations-across-society-to-combat-the-infodemic [accessed 2021-09-07]

12.  Vaccine Misinformation Management Field Guide. UNICEF. 2020 Dec. URL: https://www.unicef.org/mena/reports/vaccine-misinformation-management-field-guide [accessed 2021-09-07]

13.  An ad hoc WHO technical consultation managing the COVID-19 infodemic: call for action, 7-8 April 2020. World Health Organization. 2020. URL: https://apps.who.int/iris/handle/10665/334287 [accessed 2021-09-07]

14.  1st WHO Infodemiology Conference. World Health Organization. 2020. URL: https://www.who.int/teams/risk-communication/infodemic-management/1st-who-infodemiology-conference [accessed 2021-08-19]

15.  Pre-conference: WHO Infodemiology Conference. World Health Organization. 2020. URL: https://www.who.int/teams/risk-communication/infodemic-management/pre-conference-1st-who-infodemiology-conference [accessed 2021-09-07]

16.  Post-conference: 1st WHO Infodemiology Conference. World Health Organization. 2020. URL: https://www.who.int/teams/risk-communication/infodemic-management/post-conference-1st-who-infodemiology-conference [accessed 2021-09-07]

17.  WHO public health research agenda for managing infodemics. World Health Organization. 2021 Feb 03. URL: https://www.who.int/publications/i/item/9789240019508 [accessed 2021-07-09]

18.  Purnat T, Vacca P, Burzo S, Zecchin T, Wright A, Briand S, et al. WHO digital intelligence analysis for tracking narratives and information voids in the COVID-19 infodemic. In: Studies in Health Technology and Informatics Volume 281: Public Health and Informatics. Amsterdam, the Netherlands: IOS Press; 2021:989-993.

19.  A systematic approach for undertaking a research priority-setting exercise: guidance for WHO staff. World Health Organization. 2020. URL: https://apps.who.int/iris/handle/10665/334408 [accessed 2020-09-08]

20.  Chatham House Rule. Chatham House. URL: https://www.chathamhouse.org/about-us/chatham-house-rule#:~:text=When%20a%20meeting%2C%20or%20part,other%20participant%2C%20may%20be%20revealed [accessed 2021-08-19]

21.  Declarations of interest. World Health Organization. URL: https://www.who.int/about/ethics/declarations-of-interest [accessed 2021-09-08]

22.  Hemming V, Burgman MA, Hanea AM, McBride MF, Wintle BC. A practical guide to structured expert elicitation using the IDEA protocol. Methods Ecol Evol 2017 Sep 05;9(1):169-180. [doi: 10.1111/2041-210X.12857]

23.  LimeSurvey. URL: https://www.limesurvey.org/ [accessed 2021-05-01]

24.  Emergency response framework (ERF), 2nd edition. World Health Organization. 2017 Jun 03. URL: https://www.who.int/publications/i/item/emergency-response-framework-(-erf)-2nd-ed [accessed 2021-08-23]

25.  Dunn AG, Mandl KD, Coiera E. Social media interventions for precision public health: promises and risks. NPJ Digit Med 2018 Sep 19;1(1) [FREE Full text] [doi: 10.1038/s41746-018-0054-0] [Medline: 30854472]

26.  Shane T, Noel P. Data deficits: why we need to monitor the demand and supply of information in real time. 28 Sep 2020. 2020 Sep 28. URL: https://firstdraftnews.org/long-form-article/data-deficits/ [accessed 2021-09-08]

27.  Balog-Way DHP, McComas KA. COVID-19: Reflections on trust, tradeoffs, and preparedness. J Risk Res 2020 Apr 27;23(7-8):838-848. [doi: 10.1080/13669877.2020.1758192]

28.  The Africa Infodemic Response Alliance. World Health Organization. URL: https://www.afro.who.int/aira [accessed 2021-09-08]

29.  The Collective Service: a new partnership for strengthening risk communication and community engagement in public health and humanitarian emergencies. World Health Organization. URL: https://www.who.int/teams/risk-communication/the-collective-service [accessed 2021-09-08]

30.    Tripodi F. Searching for Alternative Facts: Analyzing Scriptural Inference in Conservative News Practices. Data & Society. URL: https://datasociety.net/wp-content/uploads/2018/05/Data_Society_Searching-for-Alternative-Facts.pdf [accessed 2021-08-19]

31.    Pool J, Fatehi F, Akhlaghpour S. Infodemic, misinformation and disinformation in pandemics: scientific landscape and the road ahead for public health informatics research. Stud Health Technol Inform 2021 May 27;281:764-768. [doi: 10.3233/SHTI210278] [Medline: 34042681]

32.    Soroya SH, Farooq A, Mahmood K, Isoaho J, Zara S. From information seeking to information avoidance: understanding the health information behavior during a global health crisis. Inf Process Manag 2021 Mar;58(2):102440 [FREE Full text] [doi: 10.1016/j.ipm.2020.102440] [Medline: 33281273]

33.    Tentolouris A, Ntanasis-Stathopoulos I, Vlachakis PK, Tsilimigras DI, Gavriatopoulou M, Dimopoulos MA. COVID-19: time to flatten the infodemic curve. Clin Exp Med 2021 May;21(2):161-165 [FREE Full text] [doi: 10.1007/s10238-020-00680-x] [Medline: 33417084]

34.    Lewandowsky S, van der Linden S. Countering misinformation and fake news through inoculation and prebunking. Eur Rev Soc Psychol 2021 Feb 22:1-38. [doi: 10.1080/10463283.2021.1876983]

35.    Purnat TD. Building systems for respond to infodemics and build resilience to misinformation. LinkedIn. 2020 Dec 02. URL: https://www.linkedin.com/pulse/building-systems-respond-infodemics-build-resilience-tina-d-purnat/ [accessed 2021-09-08]

36.    Dunn A, Steffens M, Dyda A, Mandl K. Knowing when to act: a call for an open misinformation library to guide actionable surveillance. Big Data Soc 2021 May 21;8(1):205395172110187 [FREE Full text] [doi: 10.1177/20539517211018788]

37.    Joint call for papers - special issues on Infodemiology. World Health Organization. 2020 Aug 18. URL: https://www.who.int/news-room/articles-detail/joint-call-for-papers-special-issues-on-infodemiology [accessed 2021-09-08]

38.    Gruzd A, De Domenico M, Sacco P, Briand S. Studying the COVID-19 infodemic at scale. Big Data Soc 2021 Jun 10;8(1):205395172110211 [FREE Full text] [doi: 10.1177/20539517211021115]

39.    Special feature: infodemics and health security. Health Security. 2021. URL: https://www.liebertpub.com/toc/hs/19/1 [accessed 2021-09-08]

40.    3rd virtual global WHO Infodemic Management conference. World Health Organization. 2020. URL: https://www.who.int/teams/risk-communication/infodemic-management/3rd-virtual-global-who-infodemic-management-conference [accessed 2021-09-08]

41.    4th virtual WHO Infodemic Management conference: advances in social listening for public health. World Health Organization. 2021 May 04. URL: https://www.who.int/news-room/events/detail/2021/05/04/default-calendar/4th-virtual-who-infodemic-management-conference-advances-in-social-listening-for-public-health [accessed 2021-09-08]

42.    World Health Organization. URL: https://www.who.int/teams/blueprint/covid-19/covid-19-global-research-innovation-forum [accessed 2021-09-08]

43.    Gesualdo F, Bucci LM, Rizzo C, Tozzi AE. Digital tools, multidisciplinarity and innovation for communicating vaccine safety in the COVID-19 era. Hum Vaccin Immunother 2021 Mar 25:1-4. [doi: 10.1080/21645515.2020.1865048] [Medline: 33764272]

44.    Purnat TD, Vacca P, Czerniak C, Ball S, Burzo S, Zecchin T, et al. Infodemic signal detection during the COVID-19 pandemic: development of a methodology for identifying potential information voids in online conversations. JMIR Infodemiology 2021 Jul 28;1(1):e30971 [FREE Full text] [doi: 10.2196/30971] [Medline: 34447926]

45.    Purnat TD, Wilson H, Nguyen T, Briand S. EARS - a WHO platform for AI-supported real-time online social listening of COVID-19 conversations. Stud Health Technol Inform 2021 May 27;281:1009-1010. [doi: 10.3233/SHTI210330] [Medline: 34042825]

## Abbreviations

**EPI-WIN:** World Health Organization Information Network for Epidemics
**IDEA:** Investigate Discuss Estimate Aggregate protocol
**RCCE:** risk communication and community engagement
**WHO:** World Health Organization

XSL•FO

RenderX

<u>Original Paper</u>

# Characterization of Vaccine Tweets During the Early Stage of the COVID-19 Outbreak in the United States: Topic Modeling Analysis

Li Crystal Jiang[1], PhD; Tsz Hang Chu[1], MPhil; Mengru Sun[2], BA

[1]Department of Media and Communication, City University of Hong Kong, Hong Kong, Hong Kong
[2]College of Media and International Culture, Zhejiang University, Hangzhou, China

**Corresponding Author:**
Li Crystal Jiang, PhD
Department of Media and Communication
City University of Hong Kong
M5082, Run Run Shaw Creative Media Centre
18 Tat Hong Avenue, Kowloon
Hong Kong
China (Hong Kong)
Phone: 852 034429332
Email: crystal.jiang@cityu.edu.hk

## *Abstract*

**Background:** During the early stages of the COVID-19 pandemic, developing safe and effective coronavirus vaccines was considered critical to arresting the spread of the disease. News and social media discussions have extensively covered the issue of coronavirus vaccines, with a mixture of vaccine advocacies, concerns, and oppositions.

**Objective:** This study aimed to uncover the emerging themes in Twitter users' perceptions and attitudes toward vaccines during the early stages of the COVID-19 outbreak.

**Methods:** This study employed topic modeling to analyze tweets related to coronavirus vaccines at the start of the COVID-19 outbreak in the United States (February 21 to March 20, 2020). We created a predefined query (eg, "COVID" AND "vaccine") to extract the tweet text and metadata (number of followers of the Twitter account and engagement metrics based on likes, comments, and retweeting) from the Meltwater database. After preprocessing the data, we tested Latent Dirichlet Allocation models to identify topics associated with these tweets. The model specifying 20 topics provided the best overall coherence, and each topic was interpreted based on its top associated terms.

**Results:** In total, we analyzed 100,209 tweets containing keywords related to coronavirus and vaccines. The 20 topics were further collapsed based on shared similarities, thereby generating 7 major themes. Our analysis characterized 26.3% (26,234/100,209) of the tweets as *News Related to Coronavirus and Vaccine Development*, 25.4% (25,425/100,209) as *General Discussion and Seeking of Information on Coronavirus*, 12.9% (12,882/100,209) as *Financial Concerns*, 12.7% (12,696/100,209) as *Venting Negative Emotions*, 9.9% (9908/100,209) as *Prayers and Calls for Positivity*, 8.1% (8155/100,209) as *Efficacy of Vaccine and Treatment*, and 4.9% (4909/100,209) as *Conspiracies about Coronavirus and Its Vaccines*. Different themes demonstrated some changes over time, mostly in close association with news or events related to vaccine developments. Twitter users who discussed conspiracy theories, the efficacy of vaccines and treatments, and financial concerns had more followers than those focused on other vaccine themes. The engagement level—the extent to which a tweet being retweeted, quoted, liked, or replied by other users—was similar among different themes, but tweets venting negative emotions yielded the lowest engagement.

**Conclusions:** This study enriches our understanding of public concerns over new vaccines or vaccine development at early stages of the outbreak, bearing implications for influencing vaccine attitudes and guiding public health efforts to cope with infectious disease outbreaks in the future. This study concluded that public concerns centered on general policy issues related to coronavirus vaccines and that the discussions were considerably mixed with political views when vaccines were not made available. Only a small proportion of tweets focused on conspiracy theories, but these tweets demonstrated high engagement levels and were often contributed by Twitter users with more influence.

XSL·FO
**RenderX**

## Introduction

### Background

The COVID-19 pandemic has affected more than 200 countries and territories, killed more than 1.2 million people, devastated the global economy, and disrupted the daily life of billions of people [1]. Owing to the lack of effective containment measures during the early stages of the COVID-19 outbreak, many of those heavily affected placed their hope on the development of coronavirus vaccines. Ever since the early stages of the outbreak, extensive news coverage followed the progress of vaccine developments, while web users engaged in heated discussions about coronavirus vaccines or vaccines in general on various social media platforms such as Facebook, Twitter, and Instagram [2-4]. It is crucial to understand media portrayals and public discussions of coronavirus vaccines during the early stages of the outbreak because they influenced policy-making in public health and public perceptions of and attitudes toward vaccination in the later stage [5-11]. A comprehensive understanding of the public opinion during the initial phase of infectious outbreaks will inform how public health professionals and policymakers make decisions in addressing public concerns in future outbreaks of infectious diseases [12].

### Infodemic and Early Stages of Outbreaks

Frequent infectious outbreaks are an ongoing reality for globalized societies, and the early stage of an outbreak is always challenging. The beginning of an outbreak is typically characterized by a lack of accuracy, widespread misinformation, as well as heightened uncertainty and fear among the general public [13,14]. In the first couple of months of the COVID-19 pandemic, policymakers had limited knowledge about coronavirus and largely relied on data modeling for predictions and decisions. Similarly, owing to the lack of knowledge, there was little consensus among media professionals, public health professionals, and politicians over containment measures [15]. Instead, geopolitical discourses, conspiracy theories, and racial bigotry created significant amounts of noise for officials trying to manage the pandemic [16-19]. All of these issues brought intensified fear and anxiety to the public.

Social media platforms shape public experience and opinions, while also serving as platforms for public health. During the initial phase of the pandemic, social media became the hotspot of all sorts of issues for the pandemic. Previous studies have shown that social media content about COVID-19 is mixed with a deluge of stigmas, rumors, and misinformation [16-18] and is highly biased by political and social ideologies [19-21]. On February 15, 2020, the World Health Organization officially coined a phenomenon "infodemic," which refers to the rapid spread of misinformation through social media platforms and other outlets on a global scale [22-25]. An infodemic is a serious threat to public health as it greatly advocates hostile attitudes toward preventive measures and complicates our fight with the COVID-19 pandemic [26].

### COVID-19, Vaccines, and Social Media

Despite the scientific consensus that vaccination is a safe and effective approach to prevent infectious diseases, there is more controversy over the use of vaccines than over other preventive measures (eg, hand hygiene, social distancing). These concerns include fear of side effects, uncertainty about vaccine efficacy, and general mistrust of the sciences and the government. These contentions have resulted in vaccine hesitancy, declines in immunization, and even small outbreaks of vaccine-preventable diseases [27-29]. Controversies over vaccination have often manifested in social media communities, leading to increasing research investigating the spread of information and opinions about vaccines on various social media platforms. This inquiry mainly focuses on the intensified competition between provaccination and antivaccination views on social media in recent years [30]. Both manual coding and computational methods have identified similar proportions of provaccination and antivaccination content on YouTube and Twitter [31-33], but antivaccination content—produced by closely connected communities and employing sophisticated antivaccination advocacy strategies—often outweigh the provaccination content [34,35]. There is some variation across specific types of vaccines. For example, influenza-related videos contain more anti-immunization content compared to videos on measles, presumably because influenza vaccination is normally perceived as new and less efficacious [32].

Scholars propose several strategies for tackling the vaccine controversy and addressing antivaccination information on social media [23,24], such as infoveillance. Infoveillance is an emerging approach that tracks what people do and write on the internet to reflect public opinions, behaviors, knowledge, and attitudes related to health issues [36]. Major applications of infoveillance include but are not limited to monitoring health-relevant messages on the internet (eg, antivaccination sites), outlining web-based health information availability (eg, vaccine advocacies), and analyzing search engine queries to predict disease outbreaks (eg, syndromic inquiry). By analyzing social media posts related to public health issues, previous studies have successfully performed surveillance on public opinions and public sentiments [36-40], predicted prevalence and mortality across time and space [41,42], and explained how intended or unintended behavioral responses are shaped by social networks and other information features [43-45]. In the case of analyzing vaccine-related social media messages, infoveillance can provide key stakeholders (eg, health organizations, governments) the benefits of revealing public concerns over vaccines and monitoring public sentiments in real-time. It also helps identify influencers and advocates, directly engages with the vaccine targets (ie, people who are at high risk of infection), and manage misinformation and hostile messages efficiently. Infoveillance is particularly powered by big data and computational techniques as they offer very useful tools for understanding social media content in an unstructured, bottom-up manner. Previous studies have successfully used computational methods to examine public perceptions on

influenza vaccine [46], human papilloma virus vaccines [47], and childhood vaccinations [48,49].

This study aims to investigate the discussion related to coronavirus vaccines on Twitter during the early stage of the COVID-19 outbreak in the United States (February 20, 2020 to March 31, 2020). This study will contribute to our understanding of coronavirus vaccines and connections to attitudes related to vaccines by tracking back to the initial public concerns. The findings of this study will elucidate the public discussions on new vaccines or vaccines under development, and the concerns and issues revealed in this study can show the implications on public health efforts in coping with infectious disease outbreaks in the future. Provided that the coronavirus vaccines show plenty of uncertainty in efficacy and effectiveness, using unsupervised learning methods, we aim to explore the main themes that emerged from the tweets related to coronavirus vaccines during the initial stage of the pandemic in the United States (RQ1). We also seek to examine how these themes evolved over time (RQ2).

Out of the different types of misinformation, conspiracy theories have merged as a significant concern in the "social media infodemic." Since the COVID-19 pandemic, several studies have analyzed certain types of conspiracy theories such as the coronavirus as a bioweapon [49], the 5G coronavirus [50], or "Film Your Hospital" [51]. However, few studies have looked at the overall spread of conspiracy theories related to COVID-19. A recent analysis of German tweets indicated that less than 1% of the tweets analyzed were related to conspiracy theories, although partisanship boosted the spread of conspiracy theory tweets [52]. This study examined how conspiracy theories related to coronavirus vaccines were represented in the American tweets at the early stage of the outbreak (RQ3).

Previous studies also indicate that the spread of conspiracy theories and antivaccination messages follow a different pattern compared to that of provaccination messages. On social media, antivaccination content, in general, attracts more likes and engages more discussion because content producers are inclined to use a variety of persuasive strategies (eg, health narratives) and present antivaccination in the form of public criticism aggressively [30,53]. Antivaccination messages are normally produced by a small proportion of powerful influencers, but antivaccination supporters perpetuate echo chambers by actively spreading conspiracy theories and misinformation through a more decentralized network [34]. This study also expected some differences in the influences (ie, number of followers) and engagement levels when comparing different themes in Twitter vaccine discussions. Specifically, compared to the tweets discussing other vaccine-related themes, tweets discussing conspiracy theories were likely contributed by Twitter users with more followers (H1a) and produced more engagement than tweets that discuss other themes (H1b).

## Methods

### Data Source

The study period was set from February 20 to March 31, 2020. We marked this period as the early stage because it corresponded

to a sharp increase in the coronavirus case count and death toll in the United States (eg, over 181,000 cases and 3606 deaths by March 31, 2020). At the end of March 2020, the United States became the country with the most number of confirmed cases in the world. Moreover, in March 2020, most state and local governments declared COVID-19 as a public health emergency, issued stay-at-home orders, and mandated closures of schools and public meeting places [54]. We purposely chose this time frame to capture the tweets during the first phase of the COVID-19 outbreak in the United States. Meltwater [55], a commercial web-based media monitoring service, was used for data collection. Meltwater has access to the full Twitter pipelining data hosting service, providing customized reporting options with the last 15 months of Twitter history. Meltwater geotagged each tweet using the user's Twitter bio-related or other geo-related information, thus ensuring that all tweets included in the sample were posted by American Twitter users.

Using the social media monitoring and data collection platform provided by Meltwater, we collected tweets originating from the United States and written in English that were related to the coronavirus vaccine by using the following Boolean query: (covid OR coronavirus) AND (vaccine OR vaccines OR vaccination OR vaccinations OR vaccinate OR vax OR vaxine OR vaxx OR vaccinated). Using this strategy, we identified 117,718 tweets (including original tweets and quote tweets but not replies and retweets). The text of the tweet and relevant metadata, including username, date of the post, and follower count, were stored. We also stored the engagement metric provided by Meltwater, which was a composite score representing how many times a tweet was retweeted, quoted, liked, or prompted a reply by other users. A higher engagement value indicated that the tweet received more attention by other Twitter users.

### Topic Modeling

To analyze the obtained data set, we applied topic modeling—an unsupervised machine learning algorithm that allows researchers to uncover hidden thematic structures in a sizable collection of documents [56]. A topic model can "produce a set of interpretable topics (groups of words that are associated under a single theme) and assess the strength with which each document exhibits those topics" [57]. In this study, we used Latent Dirichlet Allocation (LDA), one of the widely used topic models that groups words that frequently co-occur in documents into various topics. By providing the text input and setting the desired number of topics, LDA automatically produces a set of topics, words are allocated to the topics, and the topic proportions are attributed for each document [58]. We decided to use LDA, as findings yielded by prior studies indicate that it performs well with both long and short texts. In addition, it has been previously used to examine COVID-19–related discussions on Twitter [59].
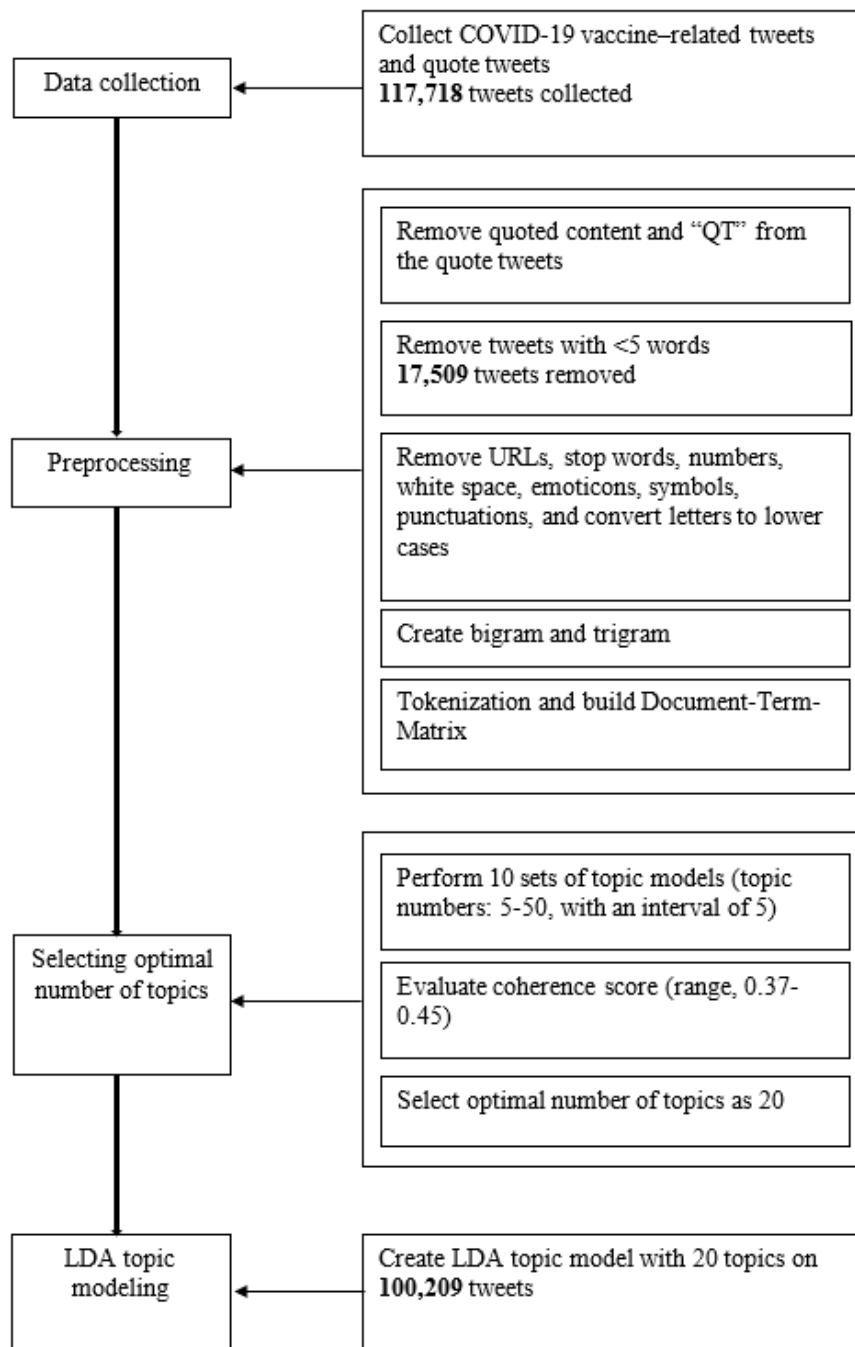
### Data Preprocessing

To prepare the corpus for LDA topic modeling, we first removed the quoted content within the quote tweets and the "QT" (meaning a quote tweet) to retain only the original content of the tweet. As the length of document plays a significant role in the topic modeling method [60], tweets with fewer than 5 words

were removed, leaving a total sample of 100,209 tweets. Following this, all the URLs within the tweets were removed. Next, the tweets were preprocessed using standard natural language processing practice [61]. We converted all the letters to lower case, removed all the stop words (eg, the, it, that), lemmatized the words, and removed numbers, white space, emoticons, symbols, and punctuation, with the use of Python packages such as NLTK (Apache) [62] and spaCy (Explosion AI) [63]. Bigram and trigram were also created and added. After tokenization, Document-Term-Matrix was built and used for the LDA topic modeling.

## Number of Topics

To determine the optimal number of topics for this tweet set, we performed 10 sets of topic models with topic numbers ranging from 5 to 50 (with intervals of 5) by implementing the LDA model from the Python package MALLET. The topic coherence —a metric focusing on the interpretability—of the 10 topic models were then calculated and evaluated for selecting the appropriate number of topics [64,65]. We decided to use the topic model with 20 topics in this study because it presented the highest topic coherence as compared with the other candidate models. Figure 1 presents the steps of data processing and creating topic models.

**Figure 1.** Data processing and analysis flowchart. LDA: Latent Dirichlet Allocation; QT: quote tweet.

## Topic Interpretation and Further Analyses

The output of the LDA topic model based on 20 topics was reviewed. Although the LDA model presumed that each document contained a mixture of topics and the model produced a probability topic distribution for each topic, we considered only the dominant topic, that is, the topic with the highest probability in that document, and categorized each tweet subject to its dominant topic [66,67]. We then reviewed the 20 top-associated terms, together with the top 5 tweets with the highest topic percentage contribution of each topic, before labeling each topic. These labels were based on the authors' background knowledge regarding vaccine hesitancy as well as the observation of coronavirus vaccine–related news and user-generated opinions on Twitter during the data collection and analysis [68]. The 3 authors involved in this study independently labeled the topics, and the resulting 3 sets of topic labels were compared. The diverse topic labels were discussed and 100% agreement between the authors was reached. The labeled topics were further grouped into distinct themes deductively following discussion. Lastly, differences between the themes on the number of followers and levels of engagement were examined. Nonparametric tests were used, as the outcome variables were not normally distributed within the current data.

## Results

### Topic Modeling

We analyzed 100,209 tweets in this study. The average number of followers of each tweet was 19,300.62 (SD 431,794.41), and the average engagement value of the tweets was 29.41 (SD 624.91). To examine themes that have emerged in coronavirus vaccine discourses on social media (RQ1), LDA modeling with 20 topics was performed. During the labeling process, it was noticed that 4 of the topics were related to news concerning human trials and testing of coronavirus vaccines. As the 4 topics were similar and closely related, the 3 authors agreed to merge these discussions into 1 overall topic, that is, *News of Vaccine Development*. Next, the remaining 17 topics were organized into 7 themes. The themes, topic labels, and associated words for each topic are presented in Table 1. The majority of the tweets were labeled as *News Related to the Coronavirus and Vaccine Development* (26,234/100,209, 26.2%) and *General Discussion and Seeking of Information on the Coronavirus* (25,425/100,209, 25.4%), followed by *Financial Concerns* (12,882/100,209, 12.9%), *Venting Negative Emotions* (12,696/100,209, 12.7%), *Prayers and Calls for Positivity* (9908/100,209, 9.9%), *Efficacy of Vaccines and Treatments* (8155/100,209, 8.1%), and *Conspiracies about Coronavirus and Its Vaccines* (4909/100,209, 4.9%).

**Table 1.** Themes and topics from coronavirus vaccine discussions on Twitter.

| Themes/topics of discussion | Associated words | Frequency of discussion (N=100,209) | Percentage of discussion | Examples of tweets |
| --- | --- | --- | --- | --- |
| **News related to coronavirus and vaccine developments** | | | | |
| News of vaccine developments | human, trial, begin, volunteer, receive | 17,435 | 17.4% | As said by the US authorities, the first clinical trial of COVID-19 vaccine on humans has been planned to begin today. The first human subject is going to get the dose today. |
| News of US government research funding/plans for the pandemic | research, fund, system, medical, government | 4012 | 4.0% | The White House approved the emergency fund to deal with COVID-19 in the United States and abroad. The fund will support the development of the COVID-19 vaccine by providing money for new equipment as well as supplies. |
| News of research plans for vaccines | pandemic, develop, effort, outbreak, step | 4787 | 4.8% | COVID-19: Mainland China has taken a new step for developing the vaccine. A team from around the world will investigate the initial results on youngsters. |
| **General discussion and seeking of information on coronavirus** | | | | |
| Seeking of information on vaccines | question, understand, information, real, cure | 3460 | 3.5% | I would like to know how the coronavirus vaccine interacts with the flu shot. Although I am not that clever to tell if we should be concerned about it, I want to raise this question out of curiosity. |
| Discussion about coronavirus trend | case, number, low, current, increase | 3268 | 3.3% | This will never work. For example, there is a rising number of confirmed cases in the Republic of Korea and Taiwan after the ease of restrictions. There is going to be nonstop waves of infection if there is no vaccine. The main purpose of isolation is reducing the load on the health care system. |
| Discussion about coronavirus and its vaccines | virus, spread, vaccine, fast, mutate | 4860 | 4.8% | Although the mutation of coronavirus is much slower than that of the flu viruses, it is an RNA virus, which normally mutates nearly 100 times faster than viruses based on DNA. It will be much difficult to control or vaccinate in the future if millions of people are infected by it as it will provide more chance for the coronavirus to mutate. |
| Comparisons with influenza | flu, kill, deadly season, thousand | 10,145 | 10.1% | First, we have the flu vaccine already. Second, compared with the influenza and the Spanish flu that have caused over 50 million deaths, the coronavirus seems more infectious. Third, compared with that with the Spanish flu, the death rate with the coronavirus is higher. Fourth, while the influenza virus has more impact on individuals older than 65 years, the coronavirus does not discriminate individuals according to age. |
| Preventive measures | protect, hand, safe, force, home | 3692 | 3.7% | The following steps can help in defeating COVID-19: stay calm and keep washing your hands with water and soap or use hand sanitizer. Keep social distancing, open doors with your elbow, and do not rub your nose, face, or shake hands with others. |
| **Financial concerns** | | | | |
| Disparity in income | American, rich, poor, capitalism, afford | 5653 | 5.6% | We, the taxpayers, are going to pay for the research on COVID-19 vaccines, which we deliver to the select few without any compensation. Rich people can acquire billions from tax cuts and chief executive officers can acquire millions from compensation. The capitalism of the Republican Party is socialism for the rich. We all are the targets. |

| Themes/topics of discussion | Associated words | Frequency of discussion (N=100,209) | Percentage of discussion | Examples of tweets |
|---|---|---|---|---|
| Price of vaccine | free, affordable, cost, charge, insurance | 7229 | 7.2% | The COVID-19 vaccine should be free of charge for people who do not have enough money for copayment for insurance or those who do not have medical insurance. The fee of my vaccination will be covered by my insurance, and I am able to pay for the difference. We have to ensure that the health insurance companies pay their part first. |
| **Efficacy of vaccines and treatments** | | | | |
| Efficacy of vaccines | prevent, cancer, infect, immunity, antibody | 4096 | 4.1% | A lot of people don't know about the COVID-19 vaccines. They are not injecting your body with the dead virus but harmless spikes. Immunity will be built to the spikes after injecting the vaccine. This could ease the worries of those opposing the vaccine. |
| Efficacy of treatments/preventions | test, treatment, effective, hospital, prove | 4059 | 4.1% | Lately, many physicians from the United States and France have asserted the effectiveness of antimalarial medication in treating COVID-19. Does it mean that the malaria vaccine would work against COVID-19 also? |
| **Conspiracies about coronavirus and its vaccines** | | | | |
| Conspiracies related to companies/stock/government | profit, market, stock, government, attempt | 4909 | 4.9% | Is it possible that the Republican Party and Trump manipulated the stock market and profited through insider trading of Moderna's stock? This biotechnology company, which invented the new vaccine, had its stock increased by 15%. |
| **Venting negative emotions** | | | | |
| Negative emotions (toward Trump and big pharmacies) | wrong, damn, business, stupid, idiot | 7019 | 7.0% | The vaccine makers could create whatever they want. Even if someone got injured or died, we cannot sue them. If someone dies, that's just bad luck. If some child dies, that's just bad luck. If someone becomes paralyzed, that's just bad luck. The profits of the pharmacies grow because we never fight back. We are just the slaves of the big pharmacies. |
| Trump-related | trump, lie, truth, blame, reality | 5677 | 5.7% | Agreed. What Trump and his incompetent administration do is to lie about everything: the seriousness of the disease, keeping the disease on a tight rein already, getting a vaccine soon. |
| **Prayers and calls for positivity** | | | | |
| Emotions/prayers | good, hope, happen, pretty, remember | 5228 | 5.2% | That is some good news! All of us need some optimism. |
| Calls for positivity | great, love, call, idea, good | 4680 | 4.7% | I enjoy seeing the positivity in the current state. |

## Themes and Topics From Coronavirus Vaccine Discussions on Twitter

### News Related to the Coronavirus and Vaccine Development

During the COVID-19 pandemic, social media users frequently shared news related to coronavirus as well as the development of coronavirus vaccines. There were 3 topics under this theme: *News of Vaccine Development*, *News of US government Research Funding/Plans for the Pandemic*, and *News of Research Plans for Vaccines*. Tweets categorized in this theme included general news on the progress of human trials and vaccine development across different countries (eg, Germany), announcements of US government funding for scientists and companies conducting research, and upcoming prevention plans for the pandemic released by official bodies as well as coronavirus vaccine research across the globe (eg, "The White House has permitted an emergency funding of US $1 billion overall in order to fight the COVID-19 outbreak. The emergency fund will offer resources as well as financial support for COVID-19 vaccine development for the states").

### General Discussion and Seeking of Information on the Coronavirus

A total of 5 topics were grouped under this theme: *Seeking of Information on Vaccines, Discussion of the Coronavirus and Its Vaccines, Discussion of Coronavirus Spread and Infection Trends, Comparisons with Influenza,* and *Preventive Measures.* The coronavirus was often compared with the influenza virus in terms of death rate, speed of transmission, and so on (eg, "Up till now, there is no cure for COVID-19 but only treatment for the symptoms. The long-term plan is to invent a new vaccine; yet, there would be no vaccine available in the next couple of months"). The importance of preventive measures, including handwashing and social distancing, was also stressed because there is currently no vaccine nor effective treatment for the coronavirus infection (eg, "We all need to get rid of bad habits. Stop touching your face when you are in public space. Scratch your nose only after washing your hands or scratch it with your sleeve. And remember to wash your hands once you get home").

### Financial Concerns

There were 2 topics under this theme: *Disparity over Income* and *Price of the Vaccines.* In the topic *Disparity over Income*, conversations were related to the gap between the rich and the poor during the pandemic as well as the differences in access to future coronavirus vaccines (eg, "All this is turning into a class war now. Only rich people can get the COVID-19 vaccine as none of us can be sure that the vaccine will be affordable for everyone"). Worries of inequality brought about by capitalism in obtaining vaccination were also expressed (eg, "Capitalism should never get closed to health care systems. The operating costs of the traditional Medicare and the administrative costs of the US health spending is extremely high. The new vaccine should be free for everyone"). As for the price of vaccination, "free" instead of "affordable" coronavirus vaccines for all Americans were urged (eg, "Citizens who were not able to pay for the COVID-19 vaccine will just keep spreading the disease. The COVID-19 vaccine should be affordable or free for everyone!").

### Venting Negative Emotions

This theme had 2 topics: *Negative Emotions (toward Trump and big pharmaceutical companies)* and *Trump-related frustrations.* Negative emotions, including anger and disappointment, toward Donald Trump or big pharmaceutical companies, were presented, as those Twitter users believed that Trump and Big Pharma were trying to profit from the pandemic. Additionally, negative emotions were expressed toward Trump explicitly owing to claims he made that are believed to have been mistaken, such as the claim that receiving the influenza vaccine would prevent COVID-19 (eg, "He [Trump] actually believed that a flu shot could fight COVID-19. I do not understand how people with brains elected this guy").

### Prayers and Calls for Positivity

The 2 topics under this theme were *Emotional Expressions/Prayers* and *Calls for Positivity.* Tweets allocated within these 2 topics included messages that aimed to encourage others during the pandemic, expressed hopes for and needs for effective coronavirus vaccines, and hopes for an end to the pandemic (eg, "Let us hope that the COVID-19 situation will be resolved when we have a vaccine/cure for it!").

### Efficacy of the Vaccine and Treatment

The 2 topics under this theme were *Efficacy of Vaccine* and *Efficacy of Treatment/Prevention.* These topics stressed the uncertainties of how well the vaccines for coronavirus work as well as the effectiveness of the current treatment and prevention strategies (eg, "I have learnt from some journals that medicines such as chloroquine, hydroxychloroquine, and azithromycin could be used as treatment or prophylaxis of COVID-19. I hope such treatments can help buying us time while getting a vaccine for COVID-19").

### Conspiracies About Coronavirus and Its Vaccines

There were different conspiracies about coronavirus and its vaccines on social media (RQ3). Many of these were related to the companies developing coronavirus vaccines, stock markets, as well as the government. For example, some tweets were claiming that the coronavirus vaccine would contain a microchip that would allow the government or company to track the vaccine receivers (eg, "Once the COVID-19 vaccines are launched, people will be motived by fear to receive the vaccines that have microchips in it."). There were also claims that the US government spread the coronavirus deliberately and withheld the coronavirus vaccines (eg, "I think the US has the cure already because it invented this bioweapon. It does not want other parties to wreck its cautiously crafted plans for devastation and racketeering")
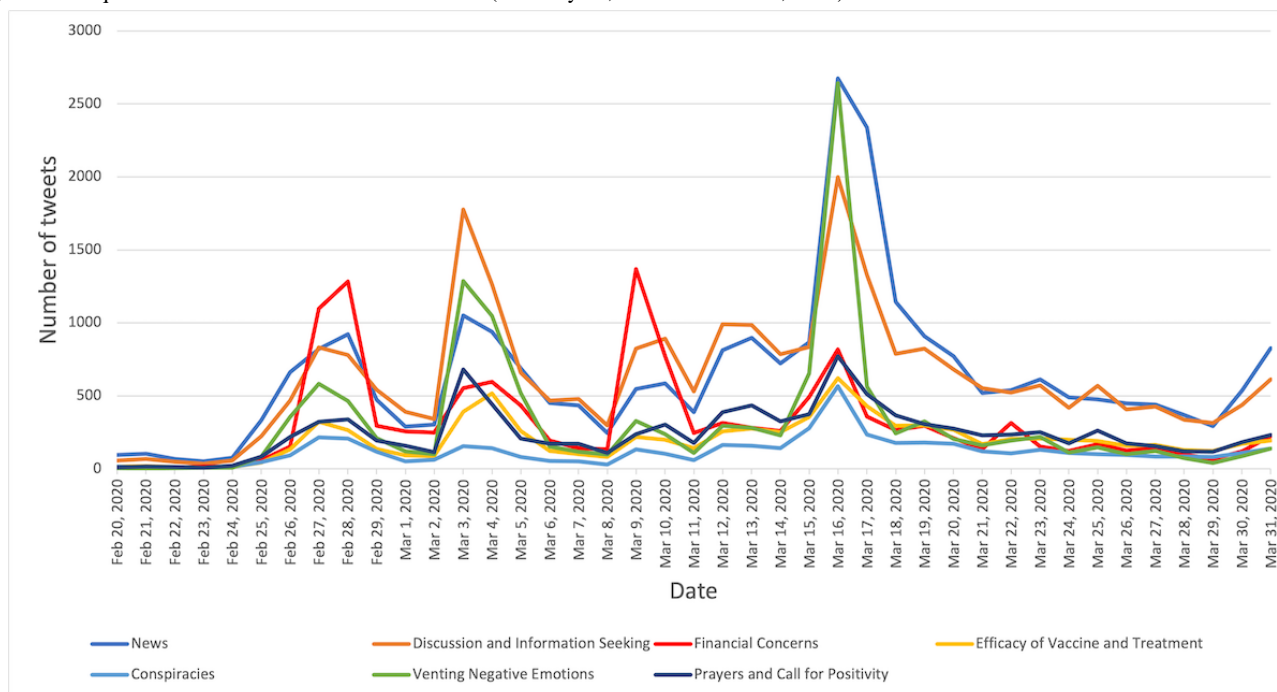
### Themes Across Time

Figure 2 shows the changes in the coronavirus vaccine–related discussions on Twitter based on the themes identified across the data collection period, that is, February 20, 2020 to March 31, 2020 (RQ2). Among the 7 themes identified, *News Related to Coronavirus and Vaccine Development* and *General Discussion and Seeking of Information on Coronavirus* were the most frequently presented overall. Coronavirus vaccine–related discourses on Twitter were promoted by breaking news or announcements and speeches made by the governments and political elites. As shown in Figure 2, there were several peaks in the coronavirus vaccine–related discussions at the early stage of the outbreak. The discussions were elevated on March 16, 2020, which corresponded to the trending news that the Trump administration was attempting to offer large sums of money to a German company in exchange for exclusive access to a possible coronavirus vaccine on March 15, 2020. Coronavirus vaccine discourses regarding *Financial Concerns* reached the peak and exceeded other themes on February 28, 2020 and March 9, 2020. The rising discussions about the prices and affordability of the coronavirus vaccines were related to Health and Human Services Secretary Alex Azar's refusal of promising affordable coronavirus vaccines for all US citizens on February 27, 2020 and Bernie Sanders' promises of free coronavirus vaccine for all Americans on March 9, 2020, respectively. We observed some co-occurring patterns across the themes during the same peaks and time periods. General discussions and information on coronavirus highly mirrored the themes of news related to coronavirus and vaccine

developments. The expression of negative emotions also increased when the discussions of these 2 themes reached a spike. There were also some observed differences in the themes across time. The efficacy of the vaccine and treatment, conspiracies about the coronavirus and vaccines, and prayers and calls for positivity appeared more periodically, while other themes (ie, news related to coronavirus and vaccine developments, general discussion and information on coronavirus, financial concerns, venting negative emotions) were more episodic, featured with several peaks instigated by breaking news or events related to vaccine developments.

**Figure 2.** Frequencies of themes of the tweets over time (February 20, 2020 to March 31, 2020).



### Differences in the Follower Numbers and Engagement Level

#### Analyses of Follower Numbers and Engagement Level

H1a and H1b hypothesized that conspiracy tweets' contributors had more followers, and conspiracy tweets received higher levels of engagement than tweets with other themes. To examine the differences between themes in the number of followers and levels of engagement, further analyses were performed. First, the results of classification of LDA topic modeling were attached in the original data (which contained the metadata, including the number of followers and engagement metric of each tweet).

The data were then entered into the SPSS software (IBM Corp). Next, we created a categorical variable according to the themes and used it as the independent variable to examine the differences in the number of followers and levels of engagement across themes by using the Kruskal-Wallis $H$ test. Post-hoc analysis was also performed using the Bonferroni-corrected Dunn test. As the hypotheses focus on the difference between conspiracy tweets and the other tweets with themes that presented attitudes and concerns toward the coronavirus vaccines, tweets labeled as news or discussion/information seeking were excluded from the analyses. Table 2 presents the median and mean rank of numbers of followers and levels of engagement.

**Table 2.** Median and mean ranks of the followers and engagement among themes.

| Themes | Followers | | Engagement | |
|---|---|---|---|---|
| | Median | Mean rank | Median | Mean rank |
| Financial concerns | 685 | 24,411.87[a] | 1 | 20,836.52[a] |
| Efficacy of vaccines and treatments | 740 | 24,957.68[a,b] | 1 | 20,784.34[a] |
| Conspiracies about coronavirus and vaccines | 770 | 25,095.54[b] | 1 | 20,631.67[a] |
| Venting negative emotions | 616 | 23,615.43[c] | 0 | 19,275.31[b] |
| Prayers/calls for positivity | 620 | 23,776.11[c] | 1 | 20,807.51[a] |

[a-c]Same superscripts in the same column indicate no significant statistical differences ($P>.05$); different superscripts in the same column indicate significant statistical differences ($P<.05$).

### Followers

The results of our study suggested that there were significant differences in the number of followers between different themes ($\chi^2_4$=77.8, $P$<.001). The post-hoc test further suggested that conspiracy tweets were more likely to be posted by users with a large number of followers than the tweets classified as *Venting Negative Emotions* ($P$<.001), *Prayers and Calls for Positivity* ($P$<.001), and *Financial Concerns* ($P$=.04). Tweets that discussed the efficacy of vaccines and treatments were also more likely to be posted by users with a large number of followers than the tweets classified as *Venting Negative Emotions* ($P$<.001) and *Prayers and Call for Positivity* ($P$<.001). Similarly, tweets expressing financial concerns were more likely to be posted by users with a large number of followers than the tweets classified as *Venting Negative Emotions* ($P$<.001) and *Prayers and Calls for Positivity* ($P$=.007). As conspiracy tweets were more likely to be posted by users with more followers than the tweets identified as *Financial Concerns, Venting Negative Emotions,* and *Prayers and Calls for Positivity*, H1a was partially supported.

### Engagement Levels

The results of our study suggested that there were significant differences in the levels of engagement between different themes ($\chi^2_4$=155.8, $P$<.001). The post-hoc test further suggested that tweets classified as *Venting Negative Emotions* significantly received lower levels of engagement than tweets classified as *Conspiracies about Coronavirus and Its Vaccines* ($P$<.001), *Efficacy of Vaccines and Treatments* ($P$<.001), *Prayers and Calls for Positivity* ($P$<.001), and *Financial Concerns* ($P$<.001). As conspiracy tweets only received higher levels of engagement than the tweets classified as *Venting Negative Emotions*, H1b was partially supported.

## Discussion

### Principal Findings

This study examined how American Twitter users discussed coronavirus vaccines during the initial stage of the COVID-19 pandemic. Using the technique of topic modeling, this study identified 7 themes in Twitter discussions. While approximately one-fourth of the tweets were about news updates related to coronavirus and vaccine developments, the remaining tweets consisted of general discussion and information seeking on coronavirus, expressions of financial concerns, disclosures of negative emotions, prayers and calling for positivity, discussions of vaccine and treatment efficacy, and conspiracy theories. In a close association with news or events related to vaccine developments, some themes demonstrated episodic changes and high degrees of co-occurrences. However, the themes of conspiracies about coronavirus and vaccines, prayers and calling for positivity, and efficacy of vaccines and treatments appeared in more periodic patterns. This study enriches our understanding of the public concerns related to vaccines during the early stage of the outbreak, and these shared concerns can inform public health organizations and professionals for more tailored health messages and vaccination policies.

Our results suggest that during the early stage of the pandemic, Twitter discussions related to coronavirus vaccines were centered on general policy issues and were largely mixed with political discussions. Two contextual factors presumably contributed to such characteristics. First, because key stakeholders did not quickly achieve a consensus on containment measures in the initial phase of the pandemic, vaccines were often staged in the public discourses as a potential remedy [12]. It is also understandable that when there was no specific vaccine available, individuals and communities addressed the vaccine issues from a policy-related perspective by discussing the investment and cost aspects of vaccination. Second, the discussions on coronavirus vaccines were situated in the political discourses during the presidential election. A topic revealed from this study was negative emotions toward Donald Trump and explicitly for his claim of using influenza shots to prevent coronavirus infections. Other COVID-19 studies also similarly demonstrated that Donald Trump and other politicians deeply influenced the vaccine discussions and even contributed to the spread of misinformation [34]. This was not surprising as vaccination is one of the politicized health controversies [52,69,70], and it was a strategic effort to feature the vaccine in political discourses. However, political disagreement over vaccines could be detrimental because they were often associated with vaccine hesitancy, reduced confidence in scientific and health facts [71], and decreased policy support for immunizations [72]. Recent research suggested that as different vaccines passed phase trials and were made available to the public, the discussions over vaccine efficacy and safety sharply increased in the United States [73].

Consistent with other studies that examined coronavirus vaccine sentiments and attitudes on social media over different periods of the pandemic [34,74,75], our study indicates that the public had mixed opinions and emotions over coronavirus vaccines, which may create significant barriers to reaching the vaccine-induced herd immunity. Antivaccination arguments and conspiracy theories were one of the major sources for vaccination opposition, although they did not constitute a large part of social media discussions. However, this small proportion of tweets was contributed by Twitter users with more influence. They also demonstrated higher engagement levels, thus resulting in echo chamber effects among small-size subnetworks [52]. It is observed that most themes demonstrated peaks and troughs over time but some themes (eg, *Conspiracies* and *Efficacy of Vaccines and Treatments*) were more periodic and some themes (eg, *Venting Negative Emotions*) were more episodic. We speculate that conspiracies and efficacy concerns were largely about unconfirmed but expected issues (eg, pharmacy conspiracies apply for all the vaccines); thus, such discussions were likely to merge periodically. However, unexpected events (eg, Trump's claim of using influenza shots to prevent coronavirus) will stimulate heated discussions, leading to a peak in the data. When these events were later addressed by the authority's responses, the discussions gradually vanished.

### Practical Implications

This study offers several practical implications for addressing the infodemic at the early stage of outbreaks or health crises. First, public health professionals should timely and appropriately

address the public needs for vaccine-related information. Our analysis revealed that many Twitter discussions were by people seeking more information or expressing concerns on coronavirus and vaccines. Such surges in information demand should be addressed by supplying with appropriate information that is easy to follow.

Second, health communication may differentiate communication strategies for episodic and periodic themes. As indicated by the results, episodic themes (eg, financial concerns, venting negative emotions) tended to emerge when breaking news or unexpected events occurred. Quick and appropriate responses to these events would effectively reassure the public and eliminate "epidemics of fear" [76]. For periodic themes such as conspiracy theories and efficacy concerns, regular surveillance and tailored responses can counterbalance the negative effects of these themes.

Health organizations and health professionals should make more systematic and organized efforts to address antivaccination content and other vaccine-related misinformation. Together with other studies [47-49], this study indicated that antivaccination content and misinformation about vaccines were contributed by closely connected communities and followed several clear and predictable patterns. When coronavirus vaccines were still under development, antivaccination content had been spreading on the internet along with these recurring conspiracy themes, which indicates that the battle with conspiracies and antivaccination messages is a long fight. A prebunking approach could effectively reduce the negative outcomes of conspiracy theories and misinformation about vaccines [77]. For example, recent research shows that attitudinal inoculation (eg, prewarning the audiences with common vaccine-related conspiracy theories) can develop resistance to the influence of vaccine conspiracy theories at a later stage [78].

Last but not the least, social media influencers (ie, accounts with many followers) play an important role in the spread of vaccine-related opinions. Fact-checking the content published by social media influencers may effectively limit the spread of conspiracy theories, which requires efforts from both social media platforms and the influencers themselves [79,80]. Twitter recently made some initial moves by introducing a labeling and striking system to identify and remove COVID-19 misinformation [81]. In a related vein, social media influencers are also encouraged by social media and health organizations to enhance their health literacy and their capacities for fact-checking before they self-proclaim as vaccine activists or public health activists on social media platforms.

## Limitations

This study has the following limitations. First, the study findings are limited to the Twitter discussions during the first phase of the COVID-19 pandemic in the United States. The public's

concerns might have changed over time as the development of vaccines progressed. The recent suspension of a coronavirus vaccine owing to adverse effects has brought a lot of discussions on vaccine safety [9]. The peculiar political environment (eg, presidential election year) may also have contributed to the patterns of the results. Further research is encouraged to look at discussions related to coronavirus vaccines on different social media forms and in different countries. Longitudinal studies and comparisons across countries or regions are particularly preferred to examine the dynamics and heterogeneity in the spread of information and opinions. This study only captured the Twitter discussions during the first stage of the COVID-19 pandemic. Future research could employ larger data sets from Twitter or other social media platforms, especially the latest data sets, to reveal the bigger picture of public concerns over coronavirus vaccines.

There are also some limitations in this analysis. For example, we relied on keyword inquiry to extract vaccine tweets from a database, but we cannot guarantee that all posts were related to coronavirus vaccine conversations. Some outliers might have been included in the data. When interpreting the topic themes, although the 3 authors independently coded the 20 topics, the intercoder reliability was not calculated owing to the small number of topics revealed from the LDA results. The analysis also did not distinguish the nature of Twitter accounts, which may be a mixture of personal, organizational, and bot accounts. Bot accounts may have contributed to a certain portion of the Twitter discussions, but we did not estimate the potential bot traffic. Because Twitter data did not account for users' demographics, while we had a limited understanding of the types of users engaged in the discussion (ie, the number of followers), we do not know more details about the users who were contributing to the discourse. Provided the difficulty of manual classification, future studies should seek to apply more sophisticated machine learning techniques to identify the types of Twitter accounts—ideally, the characteristics of personal accounts (eg, political ideology). Such knowledge will allow us to go beyond the aggregated data to look at individual users.

## Conclusion

Overall, the spread of information and opinions on social media platforms during the early stage of the outbreak has profoundly affected individuals' beliefs and attitudes toward vaccines and, ultimately, their vaccination decisions. During the early stage of the COVID-19 pandemic in the United States, Twitter discussions related to coronavirus vaccines were centered on general policy issues and were largely mixed with political discussions. The public discussions demonstrated mixed concerns for coronavirus vaccines even before the vaccines were available, and some concerns appeared periodically. These issues call for more preparatory work to cope with the infodemic challenge and to handle infectious breaks in the future.

XSL•FO

**RenderX**

## Authors' Contributions

LCJ and THC conceptualized and designed the study. THC and MS collected and analyzed the data. All authors drafted the manuscript and revised the final manuscript.

## Conflicts of Interest

None declared.

## References

1.  COVID-19 live updates. Worldometer. URL: https://www.worldometers.info/coronavirus/ [accessed 2020-12-30]
2.  Merchant RM, Lurie N. Social Media and Emergency Preparedness in Response to Novel Coronavirus. JAMA 2020 May 26;323(20):2011-2012. [doi: 10.1001/jama.2020.4469] [Medline: 32202611]
3.  Chan A, Nickson C, Rudolph J, Lee A, Joynt G. Social media for rapid knowledge dissemination: early experience from the COVID-19 pandemic. Anaesthesia 2020 Dec;75(12):1579-1582 [FREE Full text] [doi: 10.1111/anae.15057] [Medline: 32227594]
4.  Callaway E. Coronavirus vaccines: five key questions as trials begin. Nature 2020 Mar;579(7800):481. [doi: 10.1038/d41586-020-00798-8] [Medline: 32203367]
5.  Blasi P, King D, Henrikson N. HPV Vaccine Public Awareness Campaigns: An Environmental Scan. Health Promot Pract 2015 Nov;16(6):897-905. [doi: 10.1177/1524839915596133] [Medline: 26220277]
6.  Manning ML, Davis J. Journal Club: Twitter as a source of vaccination information: content drivers and what they're saying. Am J Infect Control 2013 Jun;41(6):571-572. [doi: 10.1016/j.ajic.2013.02.003] [Medline: 23726549]
7.  Chen N, Murphy S. Examining the role of media coverage and trust in public health agencies in H1N1 influenza prevention. International Public Health Journal. 2011. URL: https://www.proquest.com/docview/1727485229/abstract/B20E36AA192647FCPQ/1?accountid=10134 [accessed 2021-09-07]
8.  Chen W, Stoecker C. Mass media coverage and influenza vaccine uptake. Vaccine 2020 Jan 10;38(2):271-277. [doi: 10.1016/j.vaccine.2019.10.019] [Medline: 31699506]
9.  Moran MB, Chatterjee JS, Frank LB, Murphy ST, Zhao N, Chen N, et al. Individual, Cultural and Structural Predictors of Vaccine Safety Confidence and Influenza Vaccination Among Hispanic Female Subgroups. J Immigr Minor Health 2017 Aug;19(4):790-800 [FREE Full text] [doi: 10.1007/s10903-016-0428-9] [Medline: 27154236]
10.  Shropshire AM, Brent-Hotchkiss R, Andrews UK. Mass media campaign impacts influenza vaccine obtainment of university students. J Am Coll Health 2013;61(8):435-443. [doi: 10.1080/07448481.2013.830619] [Medline: 24152021]
11.  Xu Z, Ellis L, Laffidy M. News Frames and News Exposure Predicting Flu Vaccination Uptake: Evidence from U.S. Newspapers, 2011-2018 Using Computational Methods. Health Commun 2020 Sep 14:1-9. [doi: 10.1080/10410236.2020.1818958] [Medline: 32927970]
12.  Cuello-Garcia C, Pérez-Gaxiola G, van Amelsvoort L. Social media can have an impact on how we manage and investigate the COVID-19 pandemic. J Clin Epidemiol 2020 Nov;127:198-201 [FREE Full text] [doi: 10.1016/j.jclinepi.2020.06.028] [Medline: 32603686]
13.  Zhang Y, Tambo E, Djuikoue IC, Tazemda GK, Fotsing MF, Zhou XN. Early stage risk communication and community engagement (RCCE) strategies and measures against the coronavirus disease 2019 (COVID-19) pandemic crisis. Glob Health J 2021 Mar;5(1):44-50 [FREE Full text] [doi: 10.1016/j.glohj.2021.02.009] [Medline: 33850632]
14.  Lal A, Ashworth HC, Dada S, Hoemeke L, Tambo E. Optimizing Pandemic Preparedness and Response Through Health Information Systems: Lessons Learned From Ebola to COVID-19. Disaster Med Public Health Prep 2020 Oct 02:1-8 [FREE Full text] [doi: 10.1017/dmp.2020.361] [Medline: 33004102]
15.  Mocatta G, Hawley E. The coronavirus crisis as tipping point: communicating the environment in a time of pandemic. Media International Australia 2020 Aug 17;177(1):119-124. [doi: 10.1177/1329878x20950030]
16.  Apuke O, Omar B. Fake news and COVID-19: modelling the predictors of fake news sharing among social media users. Telematics and Informatics 2021 Jan;56:101475 [FREE Full text] [doi: 10.1016/j.tele.2020.101475]
17.  Freckelton Qc I. COVID-19: Fear, quackery, false representations and the law. Int J Law Psychiatry 2020;72:101611 [FREE Full text] [doi: 10.1016/j.ijlp.2020.101611] [Medline: 32911444]
18.  Pennycook G, McPhetres J, Zhang Y, Lu JG, Rand DG. Fighting COVID-19 Misinformation on Social Media: Experimental Evidence for a Scalable Accuracy-Nudge Intervention. Psychol Sci 2020 Jul;31(7):770-780 [FREE Full text] [doi: 10.1177/0956797620939054] [Medline: 32603243]
19.  Calvillo DP, Ross BJ, Garcia RJB, Smelter TJ, Rutchick AM. Political Ideology Predicts Perceptions of the Threat of COVID-19 (and Susceptibility to Fake News About It). Social Psychological and Personality Science 2020 Jul 22;11(8):1119-1128. [doi: 10.1177/1948550620940539]
20.  Jamison A, Broniatowski DA, Smith MC, Parikh KS, Malik A, Dredze M, et al. Adapting and Extending a Typology to Identify Vaccine Misinformation on Twitter. Am J Public Health 2020 Oct;110(S3):S331-S339. [doi: 10.2105/AJPH.2020.305940] [Medline: 33001737]

21. Rothgerber H, Wilson T, Whaley D, Rosenfeld DL, Humphrey M, Moore AL, et al. Politicizing the COVID-19 Pandemic: Ideological Differences in Adherence to Social Distancing. PsyArXiv. Preprint posted online April 20, 2020 [FREE Full text] [doi: 10.31234/osf.io/k23cv]

22. Munich security conference. World Health Organization. 2020 Feb 15. URL: https://www.who.int/dg/speeches/detail/munich-security-conference [accessed 2020-12-30]

23. Eysenbach G. How to Fight an Infodemic: The Four Pillars of Infodemic Management. J Med Internet Res 2020 Jun 29;22(6):e21820 [FREE Full text] [doi: 10.2196/21820] [Medline: 32589589]

24. 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]

25. Zarocostas J. How to fight an infodemic. The Lancet 2020 Feb;395(10225):676 [FREE Full text] [doi: 10.1016/s0140-6736(20)30461-x]

26. Abdul-Mageed M, Diab M, Kübler S. SAMAR: Subjectivity and sentiment analysis for Arabic social media. Computer Speech & Language 2014 Jan;28(1):20-37 [FREE Full text] [doi: 10.1016/j.csl.2013.03.001]

27. Chan M, Jamieson K, Albarracin D. Prospective associations of regional social media messages with attitudes and actual vaccination: A big data and survey study of the influenza vaccine in the United States. Vaccine 2020 Sep 11;38(40):6236-6247 [FREE Full text] [doi: 10.1016/j.vaccine.2020.07.054] [Medline: 32792251]

28. Hill HA, Elam-Evans LD, Yankey D, Singleton JA, Kang Y. Vaccination Coverage Among Children Aged 19-35 Months - United States, 2017. MMWR Morb Mortal Wkly Rep 2018 Oct 12;67(40):1123-1128 [FREE Full text] [doi: 10.15585/mmwr.mm6740a4] [Medline: 30307907]

29. More than 140,000 die from measles as cases surge worldwide. World Health Organization. 2019 Dec 05. URL: https://www.who.int/news/item/05-12-2019-more-than-140-000-die-from-measles-as-cases-surge-worldwide [accessed 2020-12-30]

30. Beguerisse-Díaz M, McLennan AK, Garduño-Hernández G, Barahona M, Ulijaszek SJ. The 'who' and 'what' of #diabetes on Twitter. Digit Health 2017;3:2055207616688841 [FREE Full text] [doi: 10.1177/2055207616688841] [Medline: 29942579]

31. Johnson NF, Velásquez N, Restrepo NJ, Leahy R, Gabriel N, El Oud S, et al. The online competition between pro- and anti-vaccination views. Nature 2020 Jun;582(7811):230-233. [doi: 10.1038/s41586-020-2281-1] [Medline: 32499650]

32. Dredze M, Broniatowski D, Hilyard K. Zika vaccine misconceptions: A social media analysis. Vaccine 2016 Jun 24;34(30):3441-3442 [FREE Full text] [doi: 10.1016/j.vaccine.2016.05.008] [Medline: 27216759]

33. Yiannakoulias N, Slavik C, Chase M. Expressions of pro- and anti-vaccine sentiment on YouTube. Vaccine 2019 Apr 03;37(15):2057-2064. [doi: 10.1016/j.vaccine.2019.03.001] [Medline: 30862365]

34. Germani F, Biller-Andorno N. The anti-vaccination infodemic on social media: A behavioral analysis. PLoS One 2021;16(3):e0247642 [FREE Full text] [doi: 10.1371/journal.pone.0247642] [Medline: 33657152]

35. Walter D, Ophir Y, Jamieson KH. Russian Twitter Accounts and the Partisan Polarization of Vaccine Discourse, 2015-2017. Am J Public Health 2020 May;110(5):718-724. [doi: 10.2105/AJPH.2019.305564] [Medline: 32191516]

36. Eysenbach G. Infodemiology and infoveillance: framework for an emerging set of public health informatics methods to analyze search, communication and publication behavior on the Internet. J Med Internet Res 2009 Mar 27;11(1):e11 [FREE Full text] [doi: 10.2196/jmir.1157] [Medline: 19329408]

37. Ceron A, Curini L, Iacus SM, Porro G. Every tweet counts? How sentiment analysis of social media can improve our knowledge of citizens' political preferences with an application to Italy and France. New Media & Society 2013 Apr 04;16(2):340-358. [doi: 10.1177/1461444813480466]

38. Paul M, Dredze M. Discovering health topics in social media using topic models. PLoS One 2014;9(8):e103408 [FREE Full text] [doi: 10.1371/journal.pone.0103408] [Medline: 25084530]

39. Stieglitz S, Dang-Xuan L. Emotions and Information Diffusion in Social Media—Sentiment of Microblogs and Sharing Behavior. Journal of Management Information Systems 2014 Dec 08;29(4):217-248. [doi: 10.2753/MIS0742-1222290408]

40. Eichstaedt JC, Schwartz HA, Kern ML, Park G, Labarthe DR, Merchant RM, et al. Psychological language on Twitter predicts county-level heart disease mortality. Psychol Sci 2015 Feb;26(2):159-169 [FREE Full text] [doi: 10.1177/0956797614557867] [Medline: 25605707]

41. Gibbs M, Meese J, Arnold M, Nansen B, Carter M. Funeral and Instagram: death, social media, and platform vernacular. Information, Communication & Society 2014 Dec 15;18(3):255-268. [doi: 10.1080/1369118X.2014.987152]

42. Calvo Gallardo E, Fernandez de Arroyabe JC, Arranz N. Preventing Internal COVID-19 Outbreaks within Businesses and Institutions: A Methodology Based on Social Networks Analysis for Supporting Occupational Health and Safety Services Decision Making. Sustainability 2020 Jun 06;12(11):4655. [doi: 10.3390/su12114655]

43. Scanfeld D, Scanfeld V, Larson E. Dissemination of health information through social networks: twitter and antibiotics. Am J Infect Control 2010 Apr;38(3):182-188 [FREE Full text] [doi: 10.1016/j.ajic.2009.11.004] [Medline: 20347636]

44. Valente T, Gallaher P, Mouttapa M. Using social networks to understand and prevent substance use: a transdisciplinary perspective. Subst Use Misuse 2004;39(10-12):1685-1712. [doi: 10.1081/ja-200033210] [Medline: 15587948]

XSL•FO
RenderX

45. Wood MJ. Propagating and Debunking Conspiracy Theories on Twitter During the 2015-2016 Zika Virus Outbreak. Cyberpsychol Behav Soc Netw 2018 Aug;21(8):485-490 [FREE Full text] [doi: 10.1089/cyber.2017.0669] [Medline: 30020821]

46. Nawa N, Kogaki S, Takahashi K, Ishida H, Baden H, Katsuragi S, et al. Analysis of public concerns about influenza vaccinations by mining a massive online question dataset in Japan. Vaccine 2016 Jun 08;34(27):3207-3213. [doi: 10.1016/j.vaccine.2016.01.008] [Medline: 26776467]

47. Pruss D, Fujinuma Y, Daughton A, Paul M, Arnot B, Albers Szafir D, et al. Zika discourse in the Americas: A multilingual topic analysis of Twitter. PLoS One 2019;14(5):e0216922 [FREE Full text] [doi: 10.1371/journal.pone.0216922] [Medline: 31120935]

48. Tangherlini TR, Roychowdhury V, Glenn B, Crespi CM, Bandari R, Wadia A, et al. "Mommy Blogs" and the Vaccination Exemption Narrative: Results From A Machine-Learning Approach for Story Aggregation on Parenting Social Media Sites. JMIR Public Health Surveill 2016 Nov 22;2(2):e166 [FREE Full text] [doi: 10.2196/publichealth.6586] [Medline: 27876690]

49. Hu D, Martin C, Dredze M, Broniatowski D. Chinese social media suggest decreased vaccine acceptance in China: An observational study on Weibo following the 2018 Changchun Changsheng vaccine incident. Vaccine 2020 Mar 17;38(13):2764-2770 [FREE Full text] [doi: 10.1016/j.vaccine.2020.02.027] [Medline: 32093982]

50. Stephens M. A geospatial infodemic: Mapping Twitter conspiracy theories of COVID-19. Dialogues in Human Geography 2020 Jun 23;10(2):276-281. [doi: 10.1177/2043820620935683]

51. Ahmed W, Vidal-Alaball J, Downing J, López Seguí F. COVID-19 and the 5G Conspiracy Theory: Social Network Analysis of Twitter Data. J Med Internet Res 2020 May 06;22(5):e19458 [FREE Full text] [doi: 10.2196/19458] [Medline: 32352383]

52. Ahmed W, López Seguí F, Vidal-Alaball J, Katz MS. COVID-19 and the "Film Your Hospital" Conspiracy Theory: Social Network Analysis of Twitter Data. J Med Internet Res 2020 Oct 05;22(10):e22374 [FREE Full text] [doi: 10.2196/22374] [Medline: 32936771]

53. Shahrezaye M, Meckel M, Steinacker L, Suter V. COVID-19's (mis)information ecosystem on Twitter: How partisanship boosts the spread of conspiracy narratives on German speaking Twitter. In: Advances in Information and Communication. Future of Information and Communication Conference 2021. Cham: Springer; Sep 27, 2020:1060-1073.

54. Chan C, Shumaker L, Maler S. Confirmed coronavirus cases in U.S. reach 100,000: Reuters tally. Reuters. 2020 Mar 28. URL: https://www.reuters.com/article/us-health-coronavirus-usa-cases-idUSKBN21E3DA [accessed 2020-12-30]

55. Meltwater. URL: https://www.meltwater.com/en [accessed 2021-09-07]

56. Campbell JC, Hindle A, Stroulia E. Latent dirichlet allocation: extracting topics from software engineering data. In: Bird C, Menzies T, Zimmermann T, editors. The Art and Science of Analyzing Software Data. 225 Wyman Street, Waltham, MA 02451, USA: Elsevier; 2015:139-159.

57. DiMaggio P, Nag M, Blei D. Exploiting affinities between topic modeling and the sociological perspective on culture: Application to newspaper coverage of U.S. government arts funding. Poetics 2013 Dec;41(6):570-606 [FREE Full text] [doi: 10.1016/j.poetic.2013.08.004]

58. Zhao X, Zhan M, Jie C. Examining multiplicity and dynamics of publics' crisis narratives with large-scale Twitter data. Public Relations Review 2018 Nov;44(4):619-632 [FREE Full text] [doi: 10.1016/j.pubrev.2018.07.004]

59. Hung M, Lauren E, Hon ES, Birmingham WC, Xu J, Su S, et al. Social Network Analysis of COVID-19 Sentiments: Application of Artificial Intelligence. J Med Internet Res 2020 Aug 18;22(8):e22590 [FREE Full text] [doi: 10.2196/22590] [Medline: 32750001]

60. Tang J, Meng Z, Nguyen X, Mei Q, Zhang M. Understanding the Limiting Factors of Topic Modeling via Posterior Contraction Analysis. 2014 Presented at: Proceedings of the 31st International Conference on Machine Learning; June 22-24; Beijing, China p. 190-198 URL: http://proceedings.mlr.press/v32/tang14.pdf

61. Grün B, Hornik K. topicmodels: An R Package for Fitting Topic Models. Journal of Statistical Software 2011;40(13):1-30 [FREE Full text] [doi: 10.18637/jss.v040.i13]

62. Natural Language Toolkit. URL: https://www.nltk.org/ [accessed 2021-09-07]

63. Citing spaCy v2 #1555. URL: https://github.com/explosion/spaCy/issues/1555 [accessed 2021-09-07]

64. Lutz C, Carr W, Cohn A, Rodriguez L. Understanding barriers and predictors of maternal immunization: Identifying gaps through an exploratory literature review. Vaccine 2018 Nov 26;36(49):7445-7455. [doi: 10.1016/j.vaccine.2018.10.046] [Medline: 30377064]

65. Mimno D, Wallach H, Talley E, Leenders M, McCallum A. Optimizing semantic coherence in topic models. 2011 Presented at: Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing; July 27–31; Edinburgh, Scotland, UK p. 262-272 URL: https://dl.acm.org/doi/10.5555/2145432.2145462

66. Surian D, Nguyen DQ, Kennedy G, Johnson M, Coiera E, Dunn AG. Characterizing Twitter Discussions About HPV Vaccines Using Topic Modeling and Community Detection. J Med Internet Res 2016 Aug 29;18(8):e232 [FREE Full text] [doi: 10.2196/jmir.6045] [Medline: 27573910]

67. Mehrotra R, Sanner S, Buntine W, Xie L. Improving LDA topic models for microblogs via tweet pooling and automatic labeling. 2013 Presented at: SIGIR'13; July 28-August 1; Dublin, Ireland URL: http://users.cecs.anu.edu.au/~ssanner/Papers/sigir13.pdf

68. Smith N, Graham T. Mapping the anti-vaccination movement on Facebook. Information, Communication & Society 2017 Dec 27;22(9):1310-1327. [doi: 10.1080/1369118x.2017.1418406]

69. Fowler EF, Gollust SE. The Content and Effect of Politicized Health Controversies. The ANNALS of the American Academy of Political and Social Science 2015 Feb 08;658(1):155-171. [doi: 10.1177/0002716214555505]

70. Saulsberry L, Fowler E, Nagler R, Gollust S. Perceptions of politicization and HPV vaccine policy support. Vaccine 2019 Aug 14;37(35):5121-5128. [doi: 10.1016/j.vaccine.2019.05.062] [Medline: 31296376]

71. Iyengar S, Massey DS. Scientific communication in a post-truth society. 2018 Nov 26 Presented at: The Arthur M. Sackler Colloquium of the National Academy of Sciences, "The Science of Science Communication III"; November 16-17; Washington DC, USA p. 7656-7661. [doi: 10.1073/pnas.1805868115]

72. Gollust SE, Dempsey AF, Lantz PM, Ubel PA, Fowler EF. Controversy undermines support for state mandates on the human papillomavirus vaccine. Health Aff (Millwood) 2010 Nov;29(11):2041-2046. [doi: 10.1377/hlthaff.2010.0174] [Medline: 21041746]

73. Dutta S, Kumar A, Dutta M, Walsh C. Tracking COVID-19 vaccine hesitancy and logistical challenges: A machine learning approach. PLoS One 2021;16(6):e0252332 [FREE Full text] [doi: 10.1371/journal.pone.0252332] [Medline: 34077467]

74. Rahul K, Jindal BR, Singh K, Meel P. Analysing Public Sentiments Regarding COVID-19 Vaccine on Twitter. 2021 Presented at: 7th International Conference on Advanced Computing and Communication Systems (ICACCS); 19-20 March; Coimbatore, India p. 488-493 URL: https://ieeexplore.ieee.org/document/9441693 [doi: 10.1109/ICACCS51430.2021.944169]

75. Kwok SWH, Vadde SK, Wang G. Tweet Topics and Sentiments Relating to COVID-19 Vaccination Among Australian Twitter Users: Machine Learning Analysis. J Med Internet Res 2021 May 19;23(5):e26953 [FREE Full text] [doi: 10.2196/26953] [Medline: 33886492]

76. Eysenbach G. SARS and population health technology. J Med Internet Res 2003;5(2):e14 [FREE Full text] [doi: 10.2196/jmir.5.2.e14] [Medline: 12857670]

77. Hameleers M. Separating truth from lies: comparing the effects of news media literacy interventions and fact-checkers in response to political misinformation in the US and Netherlands. Information, Communication & Society 2020 May 18:1-17. [doi: 10.1080/1369118x.2020.1764603]

78. van der Linden S, Dixon G, Clarke C, Cook J. Inoculating against COVID-19 vaccine misinformation. EClinicalMedicine 2021 Mar;33:100772 [FREE Full text] [doi: 10.1016/j.eclinm.2021.100772] [Medline: 33655205]

79. Trethewey SP. Strategies to combat medical misinformation on social media. Postgrad Med J 2020 Jan;96(1131):4-6 [FREE Full text] [doi: 10.1136/postgradmedj-2019-137201] [Medline: 31732511]

80. Jamison AM, Broniatowski DA, Dredze M, Sangraula A, Smith MC, Quinn SC. Not just conspiracy theories: Vaccine opponents and proponents add to the COVID-19 'infodemic' on Twitter. Harvard Kennedy School Misinformation Review 2020 Sep;1(3):1-22 [FREE Full text] [doi: 10.37016/mr-2020-38]

81. Updates to our work on COVID-19 vaccine misinformation. Twitter Safety. 2021 Mar 01. URL: https://blog.twitter.com/en_us/topics/company/2021/updates-to-our-work-on-covid-19-vaccine-misinformation.html [accessed 2021-04-27]

## Abbreviations

**LDA:** Latent Dirichlet Allocation

XSL•FO
RenderX

Original Paper

# Change in Threads on Twitter Regarding Influenza, Vaccines, and Vaccination During the COVID-19 Pandemic: Artificial Intelligence–Based Infodemiology Study

Arriel Benis[1,2], PhD; Anat Chatsubi[1], BSc; Eugene Levner[3], PhD; Shai Ashkenazi[4], MSc, MD

[1]Faculty of Industrial Engineering and Technology Management, Holon Institute of Technology, Holon, Israel

[2]Faculty of Digital Technologies in Medicine, Holon Institute of Technology, Holon, Israel

[3]Faculty of Sciences, Holon Institute of Technology, Holon, Israel

[4]Adelson School of Medicine, Ariel University, Ariel, Israel

**Corresponding Author:**
Arriel Benis, PhD
Faculty of Industrial Engineering and Technology Management
Holon Institute of Technology
Golomb St. 52
Holon, 5810201
Israel
Phone: 972 35026892
Email: arrielb@hit.ac.il

## Abstract

**Background:** Discussions of health issues on social media are a crucial information source reflecting real-world responses regarding events and opinions. They are often important in public health care, since these are influencing pathways that affect vaccination decision-making by hesitant individuals. Artificial intelligence methodologies based on internet search engine queries have been suggested to detect disease outbreaks and population behavior. Among social media, Twitter is a common platform of choice to search and share opinions and (mis)information about health care issues, including vaccination and vaccines.

**Objective:** Our primary objective was to support the design and implementation of future eHealth strategies and interventions on social media to increase the quality of targeted communication campaigns and therefore increase influenza vaccination rates. Our goal was to define an artificial intelligence–based approach to elucidate how threads in Twitter on influenza vaccination changed during the COVID-19 pandemic. Such findings may support adapted vaccination campaigns and could be generalized to other health-related mass communications.

**Methods:** The study comprised the following 5 stages: (1) collecting tweets from Twitter related to influenza, vaccines, and vaccination in the United States; (2) data cleansing and storage using machine learning techniques; (3) identifying terms, hashtags, and topics related to influenza, vaccines, and vaccination; (4) building a dynamic folksonomy of the previously defined vocabulary (terms and topics) to support the understanding of its trends; and (5) labeling and evaluating the folksonomy.

**Results:** We collected and analyzed 2,782,720 tweets of 420,617 unique users between December 30, 2019, and April 30, 2021. These tweets were in English, were from the United States, and included at least one of the following terms: "flu," "influenza," "vaccination," "vaccine," and "vaxx." We noticed that the prevalence of the terms vaccine and vaccination increased over 2020, and that "flu" and "covid" occurrences were inversely correlated as "flu" disappeared over time from the tweets. By combining word embedding and clustering, we then identified a folksonomy built around the following 3 topics dominating the content of the collected tweets: "health and medicine (biological and clinical aspects)," "protection and responsibility," and "politics." By analyzing terms frequently appearing together, we noticed that the tweets were related mainly to COVID-19 pandemic events.

**Conclusions:** This study focused initially on vaccination against influenza and moved to vaccination against COVID-19. Infoveillance supported by machine learning on Twitter and other social media about topics related to vaccines and vaccination against communicable diseases and their trends can lead to the design of personalized messages encouraging targeted subpopulations' engagement in vaccination. A greater likelihood that a targeted population receives a personalized message is associated with higher response, engagement, and proactiveness of the target population for the vaccination process.

## KEYWORDS

## *Introduction*

### Background

As online-mediated communication environments increase, social media platforms enable individuals to discuss diverse issues, express their thoughts, and debate [1-3]. Twitter is a leading social network that provides microblogging services. Users can publish posts, called tweets, with a limited length of 280 characters. Thereby, users can interact with others by responding, sharing, or showing their interest by "liking" a tweet. These interactive abilities are the fundamental building blocks of the connective nature of social networks and serve as an echo of ideas transferred among users on the platform around the globe [4]. Retrieving information in tweets' contents is challenging but is more manageable than in other social media platforms with long messages [5]. Indeed, the amount of structured and unstructured data from social media and Twitter has been increasing exponentially over the years [6,7]. Data mining and text mining enable the discovery of potentially new knowledge and contribute to developing efficient evidence-based decision-making tools [8-10] by extracting meaningful summaries, such as statistical ones, or controlled vocabularies (eg, terminology, folksonomy, taxonomy, and ontology) [11-15].

One of the most critical achievements of modern medicine is the development and widespread use of safe and efficacious vaccines. Nevertheless, their partial acceptance due to vaccine hesitancy and refusal is a significant health threat. Regarding influenza, compliance with the vaccine against it is relatively low compared with other vaccines, mainly because vaccination must be repeated annually [16]. Like other vaccines, influenza generates discussions both in the real world and online [17-20]. The COVID-19 vaccine is no exception.

Moreover, the global spread of the COVID-19 epidemic [21], its significant impact on daily life, and the relatively fast development of a vaccine against it have made the COVID-19 vaccine a critical health topic of discussion on social media. Reducing the incidence of transmissible diseases, such as influenza and COVID-19, requires achieving herd immunity [22,23], preferably by vaccination. This public health objective is achievable only with population engagement [18,19].

Social media platforms, such as Twitter, are a place of choice to share opinions and to search for (mis)information [24,25] about health care issues [26,27], including vaccines [17,18,28]. These open forums can influence opinions and vaccination decisions by hesitant individuals [29]. Discussions between provaccine advocates and "anti-vaxx" militants about vaccines' necessity, effectiveness, and safety are continuous. Moreover, the internet as a whole enables the detection of early warnings of disease outbreaks, their dissemination tracking and resilience [30], and the spread of evidence-based information [31,32]. Artificial intelligence methods and algorithms (ie, data mining, text mining, and natural language processing) have been efficiently used in the last decade to detect outbreaks, such as

influenza, based on emerging trends in internet search engine queries and social media threads [33-36]. There is a need for public health interventions [37] to make drastic stands against the spread of misinformation like that disseminated by vaccine opponents [19,38]. Related tools should be based on artificial intelligence to analyze efficiently and in an automated manner the big data generated over social media [39,40].

Understanding the changes happening during some health-related event discussions is crucial to improving health communication efficiency [41,42]. Disease prevention programs need to incorporate methods to make evidence-based information accessible to widespread populations using online resources and to increase control of biased and misleading announcements. The main focus is on advertising policies and campaigns on social media [30,43].

### Aims, Objectives, and Hypotheses

Our primary objective was to support the design and implementation of future eHealth strategies and interventions on social media to increase the quality of targeted communication campaigns and therefore increase influenza vaccination rates [18,19,44,45].

Our main aim was to define an artificial intelligence–based approach to analyze tweets, including terms related to vaccination against influenza and COVID-19. We focused on detecting co-occurring terms related to influenza vaccination and highlighting the dominant topics related to these terms. Therefore, these results must be used to build a folksonomy [46-49], which may then support the enhancement of vaccination campaigns. The methodology could be generalized to other health-related mass communications. Our research goal was to build a timely and dynamic vocabulary of the various topics related to influenza, vaccines, and vaccination posted in the English language. This vocabulary can be used as a decision support tool for health communication specialists and health policymakers, facilitating the understanding of the variations over time of different topics, such as those suggested in this study (tweets related to "influenza," "vaccines," and "vaccination").

The following 4 hypotheses guided this research:

1. Tweets are a source of understanding the reasoning to take a vaccine.
2. "Influenza," "vaccines," and "vaccination" topics are not linked directly to other topics (such as politics, economics, and fears) but are related to health matters.
3. Actuality and news impact tweet content related to vaccines and vaccination.
4. The terms and hashtags of tweets about influenza, vaccines, and vaccination can be organized in a dynamic vocabulary [50]. It can reflect the main topics and their terms discussed over time on the social media platform.

This research was granted ethical approval by the Ethics Committee of the Faculty of Technology Management of the

Holon Institute of Technology (Israel) (TM/2/2020/AB/004). The information collected on Twitter during this research was stored in a secured encrypted manner, with restricted access provided by the institution to the principal researcher (AB).

## Methods

### Overview

This study included the following 5 stages:

1. data sourcing to collect tweets and related data using the Twitter streaming application programming interface (API) [51];
2. data cleansing and storage;
3. identifying the terms, hashtags, and topics related to "influenza," "vaccines," and "vaccination;"
4. building a dynamic vocabulary, a folksonomy, to support the understanding of the relations between them; and
5. evaluating the vocabulary clusters.

### Data Sourcing

We extracted and collected tweets via the Twitter API for 16 months, between December 30, 2019, and April 30, 2021. These tweets were in English, from North America, and included at least one of the following terms: "flu," "vaccination," "vaccine," and "vaxx" (this last term was used to capture messages related to vaccination opponents as these individuals use it). We selected these terms to maximize the chance to retrieve discussions concerning a vaccine as a product, vaccination as an act or a policy, vaccination hesitancy, and influenza. Moreover, since Twitter participants use informal language, for extracting influenza-related content, we used the popular term "flu." The extraction omitted retweets and likes. The 16-month follow-up period allowed us to capture terms and topics of Twitter threads related to influenza, vaccines, and vaccination. Indeed, in the United States, 2020 involved the COVID-19 pandemic and the presidential elections.

### Data Preprocessing and Cleansing

To ensure efficient use of machine learning methods [52] on the tweet collection [53], we preprocessed it by cleansing and lemmatizing similar words appearing in posts. Data cleansing consisted of removing punctuation marks [54], mentions of users, glyphs, website addresses, and stop words [55]. Moreover, as the tweets were written in a natural language and concise manner (due to the limitation of 280 characters), a word may be written in several ways due to various reasons (eg, typos and short forms), all of which have the same or similar meaning. Lemmatization is one of the methods for overcoming this issue. It consists of replacing words by their root form (eg, "vaccine" for "vaccines"). [56]. For example, due to the COVID-19 pandemic, the tweets retrieved during the collection process contained multiple representations of the term "covid," such as "COVID19," "COVID-19," and "coronavirus." We used the Python Natural Language Toolkit (NLTK) package for lemmatization [55]. Since the nature of tweets is informal, it has been assumed that using a single representation of those words will not significantly change the tweet's context, thus improving the model's accuracy. Therefore, the frequent representations of the term "COVID" were replaced with the single form "covid," and the terms related to "influenza" were lemmatized to "flu" as the popular language used on Twitter. All the lemmas were stored in lowercase.

## Identifying the Terms and Topics Related to Influenza, Vaccines, and Vaccination

We handled the identification of the terms, hashtags, and topics related to influenza, vaccines, and vaccination by a 3-step process as follows: (1) clustering with word embedding and n-grams, (2) building a folksonomy, and (3) evaluating folksonomy clusters.

### Clustering

The objective of clustering is to segregate a set of points into groups, with each one as similar as possible and different from the others [57]. For example, in the context of text mining and specifically mining a tweet corpus, clustering can be used to group terms that are semantically similar or frequently appearing in the same message. Each cluster, according to its content, can then be annotated with a topic.

#### Word Embedding

Handling the high volume of collected tweets over time means dealing with the curse of dimensionality [58]. Therefore, a symbolic-numeric reformulation associated with dimension reduction [59] must be used to handle a large amount of data in a reasonable time and reduce the processing complexity. Word embedding is a relevant approach supporting these 2 goals; it consists of a learned numerical representation of text where words having a similar meaning in a specific context have an equal numerical representation in a vector. Globally, word embedding allows the prediction of words in a specific context. Thus, Word2Vec is a word embedding algorithm based on a neural network model learning from a large corpus of text (ie, a context) the association between words or terms. After the first training step, Word2Vec can detect synonymous words or terms, or suggest complete sentences. This is done by searching for vectors and so words with a close semantic similarity represented by cosine or Euclidean distance (ie, the similarity or the relation) between two vectors (ie, words and terms) in a space (ie, corpus) of $n$ dimensions (ie, number of words or terms in the corpus) [60]. As an example, words related to time, such as "day," "week," "month," "season," and "year," will be used in similar contexts and will be defined as semantically closed. The preprocessed data were used for creating the Gensim Word2Vec model in Python [61]. In order to see each word in the context it has with other words, we produced clusters, with the K-means algorithm [62,63], to assist decision makers in better understanding the public's perceptions of vaccines and vaccination against influenza and COVID-19. As discussions constantly evolve, the word embedding and clustering process was repeated monthly on newly collected tweets.

#### N-grams

As a complementary approach to word embedding, we built an n-gram language model predicting the probability of a sequence of words (after stop-word cleaning) to appear in our corpus of tweets. We extracted the most frequent n-grams comprising between 1 and 4 terms (n) for each week. Moreover, this process

used the Gensim Python library [61]. This approach enables health communication decision makers to learn about new growing or shrinking isolated terms and sets of terms in the discussions related to vaccination and influenza.

### Defining the Numbers of Clusters as Topics

Clustering is an unsupervised learning task and is challenging due to the need to define $k$ and the number of clusters to build. The "silhouette method" allows assessing the quality of clustering, as it determines the similarity of an object (eg, a word also called a unigram) with the content of its cluster and the likeness with the other clusters. A silhouette shows which objects (eg, words, vectors, and values) lie well within a cluster and which are less related. The graphical combination of the silhouettes of an entire clustering (eg, with $k$ clusters) into a single plot allows the appreciation of each cluster's relative quality and the overall clustering itself. The overall average silhouette width (ie, the average silhouette width of each cluster) provides an evaluation of clustering validity. A higher value of the overall average silhouette width (ie, silhouette score) is associated with better clustering with $k$, and therefore, it must be selected as the better partitioning. The silhouette method is independent of the partitioning algorithm used [64]. From our research perspective, each term must have a minimum number of occurrences to be included in the analysis. Moreover, 2 terms must have a maximum distance (number of other terms) between them in a tweet to consider their potential semantic link.

### Cluster Visualization

Cluster visualization is produced by using t-distributed stochastic neighbor embedding (t-SNE), which is a nonlinear dimensionality reduction technique for embedding high-dimensional data and visualizing it in a low-dimensional (ie, 2 or 3) space [65].

### Evaluation of the Terms in the Clusters and as N-grams

To evaluate our approach and the results of identifying the terms, hashtags, and topics related to influenza, vaccines, and vaccination, we implemented a validation process built on complementary approaches. One focused on the word embedding results, the second focused on n-grams, and the third focused on the whole by involving social media users. Thus, the terms were grouped once from a semantic perspective with word embedding on the first hand and once from a high coappearance frequency as n-grams describe the content of the explored Twitter threads in summarized ways.

The second evaluation approach consisted of using Google Trends [66] for getting the relative frequency of search terms during a specific period and in a specific geographic area. In this study, the n-grams (n between 1 to 4) were extracted from the tweets, and their weekly frequency was calculated. Next, the n-grams that appeared in the top 150 list continuously for at least 12 weeks were used as an input for a Google Trends query at the time frame they were published on Twitter. Finally, the n-grams (bi-grams) and the Google Trends query results were normalized. Their Pearson correlation coefficients were calculated by considering the weekly tweet-based n-grams and the weekly relative number of queries (comprising the n-gram terms) on the Google search engine.

The third evaluation consisted of computing Pearson correlations between the weekly frequency (between December 2020 and April 2021) of n-grams specific to vaccines, vaccination, influenza, and COVID-19, and the proportion of the population vaccinated against COVID-19.

### Informed Consent Statement

The social network data were collected in an anonymized way and following Twitter's rules. The participants of the evaluation survey provided anonymous informed consent in an electronic way on the platform before they could proceed to the completion of the questionnaire.

### Data Availability Statement

The Twitter data that support the findings of this study are not available owing to Twitter's rules and regulations. The survey data that support the findings are available from the corresponding author (AB) upon reasonable request, which will need to undergo ethical and legal approvals by the investigators' institutions. The methodology of this research will be reported in the AIMe registry for artificial intelligence in biomedical research [67].
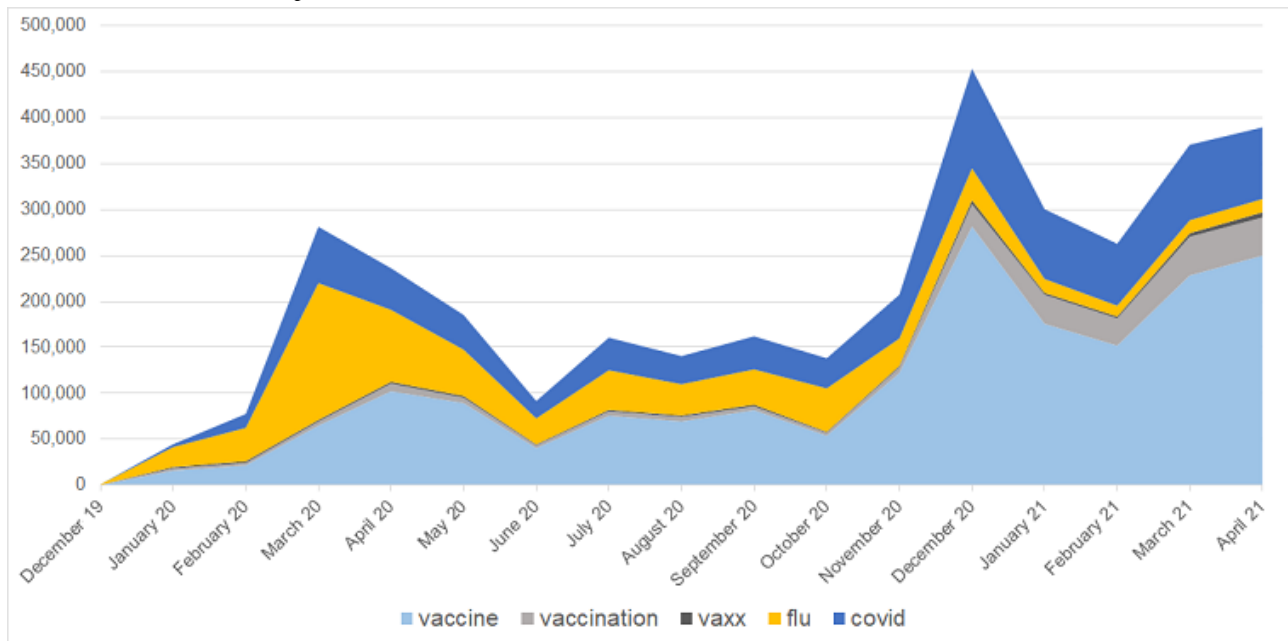
## Results

### Descriptive Statistics

A total of 2,782,720 tweets of 420,617 unique users between December 30, 2019, and April 30, 2021, were collected. The graph in Figure 1 shows the number of tweets per month (bar columns) containing at least one of the following terms (or similar after cleansing and lemmatization): (1) "flu," (2) "vaccination," (3) "vaccine," (4) "vaxx," and (5) "covid." The lines in Figure 1 show the proportion in percentage of each of these terms in the collected tweets. Although the term "covid" and its synonyms were not part of the initial keywords used for querying tweets, its emergence reflects the effect of the COVID-19 pandemic as an important topic in the discussions regarding vaccination and influenza in 2020 and 2021.

Figure 1 also shows that globally the number of tweets comprising at least one of the terms "flu," "vaccination," "vaccine," "vaxx," and "covid" has dramatically increased over the period from December 2020 to April 2021 (see also Multimedia Appendix 1). Two peaks were noticed. The first was in March 2020, with the World Health Organization declaring COVID-19 as a pandemic (March 11, 2020) and President Donald Trump promulgating COVID-19 as a national emergency (March 13, 2020). The second peak in December 2020 was related mainly to "vaccines" in response to the approval of COVID-19 vaccines (Food and Drug Administration [FDA] emergency use authorizations for Pfizer BioNTech vaccine on December 11, 2020, and Moderna vaccine on December 18, 2020). Thus, the term "vaccine" increased from approximately 35% in January 2020 to approximately 80% one year later. In contrast, the term "vaxx" (for the terms "antivaxx," "antivaxxer," "anti-vaxx," and "anti-vaxxer") was stable at 1% to 3% over the whole data collection period. Nevertheless, it is essential to take into account that vaccination opponents used various tools and communication discourse, not evocating the "anti-vaxx" term itself [68-70]. The terms related to influenza

("flu") and COVID-19 ("covid") showed an inverse correlation ($r$=–0.83, $P$<.001) at the monthly level (Multimedia Appendix 1). The use of "covid" increased linearly, starting in January 2020, with the first cases of COVID-19 spreading from China to Europe and the United States [71], until February 2021, when

it was part of approximately 35% of the collected tweets. In parallel, the use of the term "flu" decreased steadily, probably due to the low influenza activity during the 2020-2021 season [72,73].

**Figure 1.** Distribution of the number of tweets by month comprising at least one of the terms "flu," "vaccination," "vaccine," "vaxx," and "covid" between December 30, 2019, and April 30, 2021.



## Identification of the Terms and Topics Related to Influenza, Vaccines, and Vaccination

### Word Embedding

The Word2Vec algorithm was run monthly to find the optimal parameters supporting the finding of the dominant trending topics. Determination of the optimal parameters' values was performed by creating models using a different value for each parameter and calculating the silhouette score for each iteration with the "silhouette_score" function of sklearn.metrics in Python [74]. Multimedia Appendix 2 shows the parameters' values and the silhouette scores of the various models of each month. Moreover, each week, only the terms having the highest occurrence regarding the overall number of terms detected in the tweets collected in the same week were investigated. The values of these attributes were changed over time to consider the dynamic changes in social media users' lexicons impacted by the actuality.

### K-means Clustering

Using the monthly word embedding model as an input, word clusters were generated with the NLTK KMeansClusterer [75]. The clustering method groups together a given data set to a $k$ predetermined number of clusters [66,76]. The partition is performed while aiming to minimize the in-cluster variance and maximize the variance between the elements from different clusters. To determine the optimal number of clusters [77], we computed the silhouette scores of k-means clustering runs with $k \in$ [3;6]. The silhouette scores of the clustering models were generated on the 2,782,720 tweets of 420,617 unique users between December 30, 2019, and April 30, 2021, related to

141,407 n-grams with $n \in$ [2;4]. The highest silhouette score reflects this grouping, wherein the different objects are well affected to their clusters and less linked to neighboring and less relevant clusters. A higher silhouette score ($s$=0.72) was achieved with $k$=3. This score can be considered good as we clustered terms that can relate to different topics and the clusters can overlap partially [78,79]. Furthermore, by computing the Ray-Turi index [80] for $k$ between 2 and 10, and building the curve of the different generated values allowed with the Elbow method, the optimal $k$ was equal to 3 [81].

Indeed, we interpreted the content of the 3 clusters in the tweet collection of the study with consensus of domain experts (public health, infectiology, and informatics). These clusters are the bare bricks of the "vaccination against influenza during the COVID-19 pandemic" folksonomy. We defined the 3 topics dominating the content of the collected tweets as follows:

1. "Health and medicine (biological and clinical aspects)" comprising terms such as "pandemic," "COVID-19," "vaccines," "illness," "die," "variant," "children," "flu," "influenza," and "health;"
2. "Protection and responsibility" with terms such as "protection," "social distancing," "vaccination," "fighting COVID-19," and "responsibility;" and
3. "Politics" supported by terms like "trump," "biden," "lie," "government," "trust," "bill_gate," "free," "money," "president," "politics," "politicians," "elections," "vaccine," and "policy."
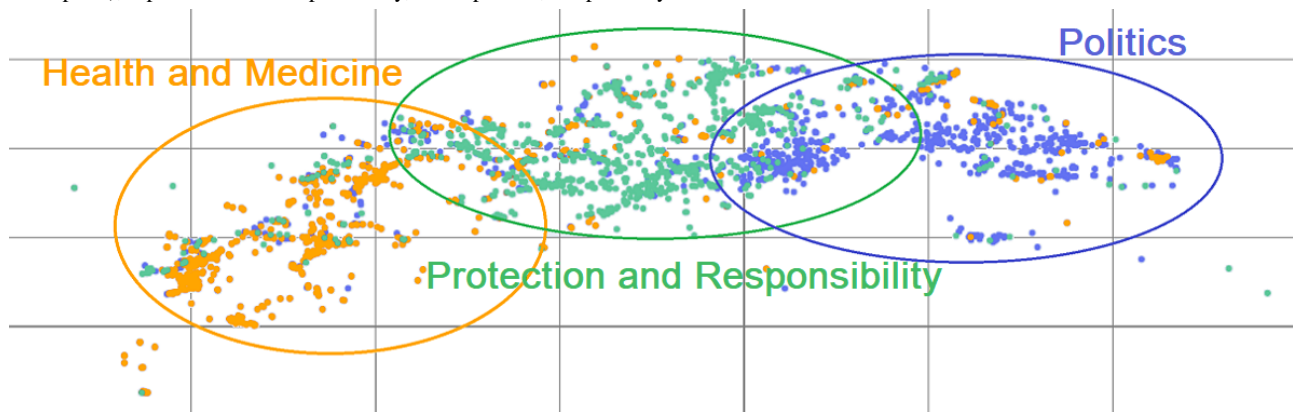
Figure 2 shows a 2-dimensional graphical representation of the 3 clusters with the 1000 most frequent n-grams for each ($n \in$

[1;4]) (Multimedia Appendix 3 and Multimedia Appendix 4), which has been generated by using the t-SNE algorithm [65].

Explicitly, this visualization (Figure 2) allows us to see the first 1000 most used terms in the tweets of each one of the previously computed clusters. It is noticeable that overlaps exist between the clusters, which is quite logical when we realize that the tweets relate in many cases to a few topics at the same time (eg, from an account dealing with *political* issues: "The *vaccines* offer good *protection* with more than 80% effectiveness. Most people will not be *sick* and the ones that will, will not get seriously *ill* or *die*").

**Figure 2.** A t-distributed stochastic neighbor embedding graphical representation of the 3 topic clusters with 1000 most frequent n-grams ($n \in$ [1;4]). Orange, seafoam (green-blue facilitating reading of the figure by color-blind individuals), and violet represent "health and medicine (biological and clinical aspects)," "protection and responsibility," and "politics," respectively.



### N-grams

The preprocessed tweets were used to extract n-grams for each week. Multimedia Appendix 4 shows the 10 most common n-grams for each $n \in$ [1;4]. For example, the words "flu" and "bad" were found close to each other in the word embedding model over the months of this study (Multimedia Appendix 3, list of the 1000 most frequent n-grams for cluster 1). Those 2 words were also a common n-gram, whether a bigram or a part of a higher degree of an n-gram. Although included in the word embedding representation, we see the relations between those 2 words in general, as they get closer to each other and in the same semantic cluster.

Following the extraction, each n-gram received its growth value, indicating an increased or decreased n-gram frequency from the previous week. The growth is used to highlight the significant changes in the n-grams and therefore in general discussions. For example, on November 9, 2020, Pfizer BioNTech published the initial results of the COVID-19 vaccine trial, which showed high efficacy against the disease. The n-grams of the same week showed a significant increase as follows: "*take, vaccine*," 774.6% (1207/51,553 vs 138/51,553) and "*get, vaccine*," 557.9% (1987/149,333 vs 302/149,333) [82].

Moreover, in mid-March 2021, we also noticed a significant increase in n-grams related to vaccination against COVID-19 due to reports on Twitter of individuals being vaccinated or local authorities inviting the population to schedule appointments for taking the vaccine (eg, "vaccine, appointment, available," +264.9% [748/18,678 in the week starting March 15, 2021, vs 205/18,678 in the week starting March 08, 2021] and "code, vaccine, appointment, available," +251.5% [942/6264 in the week starting March 29, 2021, vs 268/6264 in the week starting March 22, 2021]) [83].

Another example of Twitter's user response was during the week starting May 11, 2020. The leading n-grams were "*social distancing, flattenthecurve, trump, test*" and "*flattenthecurve, trump, test, vaccine*" (wherein "socialdistancing" and "flattenthecurve" were hashtags). Both demonstrated growth of 693.0% from the previous week (43 occurrences during the week starting May 04, 2020, vs 341 out of 516 occurrences in total). In that week, Forbes magazine published an article reporting that hospitals across the United States are "*not being overwhelmed*," suggesting that the efforts for flattening the curve have succeeded. The overall results show how the tweet threads about influenza, vaccines, vaccination, and COVID-19 dynamically evolved from the end of 2019 to mid-2021.

### Evaluations

#### Google Trends Validation

As a component of the internet, social media like Twitter are a part of how people get and share information and knowledge. Therefore, looking at queries on search engines like Google allows the evaluation of global interests in terms and topics detected on social media. Thus, we computed Pearson correlations between the weekly occurrences of n-grams in tweets and weekly queries in the Google search engine and those reported on Google Trends [84]. As an example of the consistency of the previously disclosed results, the n-gram of "flu, symptom" on Twitter and the number of queries on Google were highly correlated ($r=0.85$, $P<.001$) between January 1, 2020, and March 4, 2021 (Table 1). During these 65 weeks, this n-gram (ie, "flu, symptom") was also used to search for information about "influenza" and "symptoms."

Moreover, as we noticed the decreasing popularity of its use on Twitter, we also noticed similar behavior on Google. Additionally, the n-gram "covid, vaccine" also showed a high correlation between Twitter and Google ($r=0.85$, $P<.001$), and on the 2 platforms, its occurrence increased between January

2020 and January 2021, and then showed a parallel decrease. Globally, the top topics related to vaccines, vaccination, and COVID-19 were similar on social networks and search engines (Table 1). Thus, internet users' queries on search engines relate with the timing of topics defined by analysis of the text of our Twitter message data set.

**Table 1.** Examples of n-grams having high correlations between their trend frequencies in tweets and Google search queries.

| N-gram | Period (start date to end date) | Pearson correlation | P value |
| --- | --- | --- | --- |
| get, second, dose | January 04, 2021, to April 30, 2021 | 0.91 | <.001 |
| get, first, vaccine, shot | January 18, 2021, to April 25, 2021 | 0.89 | <.001 |
| second, vaccine | February 01, 2021, to April 30, 2021 | 0.86 | <.001 |
| flu, symptom | January 01, 2020, to April 04, 2021 | 0.85 | <.001 |
| covid, vaccine | January 20, 2020, to April 30, 2021 | 0.85 | <.001 |
| think, flu | January 01, 2020, to March 30, 2020 | 0.84 | <.001 |
| second, dose, vaccine | January 04, 2021, to April 30, 2021 | 0.84 | <.001 |
| get, second, vaccine | February 01, 2021, to April 30, 2021 | 0.84 | <.001 |
| get, covid, vaccine | March 30, 2020, to April 30, 2021 | 0.84 | <.001 |
| get, vaccine | January 01, 2020, to April 30, 2021 | 0.80 | <.001 |

### Real-World Validation

On December 11, 2020, the FDA issued an emergency use authorization for a COVID-19 vaccine. A few days later, on December 20, 2020, vaccination of the population with the Pfizer BioNTech vaccine was started. We downloaded the daily vaccination rate from Centers for Disease Control and Prevention (CDC) publications and aggregated them at the weekly level [85]. We noticed that starting in December 2020 and ending on April 30, 2021, Pearson correlations between the weekly occurrences of COVID-19 vaccination n-grams and the weekly vaccination rates (Table 2) were high and significant ($r > 0.81$, $P < .001$) [86]. These results demonstrate that the tweets of this study mirror "real-life" significant events during the pandemic.

**Table 2.** Correlations of the 5 highest n-gram trends with the vaccination rate trends reported by the Centers for Disease Control and Prevention between December 20, 2020, and April 30, 2021.

| N-gram | Pearson correlation | Number of occurrences | P value |
| --- | --- | --- | --- |
| get, first | 0.88 | 17,133 | <.001 |
| vaccine, today | 0.87 | 9205 | <.001 |
| first, vaccine | 0.83 | 9260 | <.001 |
| first, dose | 0.82 | 11,357 | <.001 |
| vaccine, shot | 0.81 | 11,113 | <.001 |

## Discussion

### Principal Findings

This research was initiated to elucidate online public perceptions regarding vaccination, mainly against seasonal influenza. However, the COVID-19 pandemic in 2020 was impressively reflected by major changes in the focus of Twitter-based discussions. The most important aspect of this study is the building of a folksonomy based on tweet text analysis, word embedding, and clustering. The 3 topics that were identified in this folksonomy were as follows:

1. General issues from the "health and medicine (biological and clinical aspects)" perspective. The initial terms used for the tweet extraction were "flu," "vaccination," "vaccine," and "vaxx." These terms are de facto strongly related to health and medicine, and generate a large spectrum of threats (ie, from asking/answering questions about symptoms, reporting health conditions, and sharing positions). The presence of terms related to the COVID-19 pandemic is understandable given the period of the data collection.

2. "Protection and responsibility" as a central dimension of the decision to take a vaccine or not. The COVID-19 pandemic showed the need for social distancing and mask wearing to reduce the spread of the virus. For these reasons, tweets related to influenza ("flu") or immunization ("vaccine" and "vaccination") and, by extension, to COVID-19 comprise threads discussing protection measures (like vaccination) and the responsibility to use them (such as taking a vaccine). It is important to highlight, based on prior studies [19,87,88], that the intent to take a vaccine is considered by the younger adult US population as an act of collective responsibility.

3. "Politics" is a cluster showing the divergence of opinions and messages of US political leaders (ie, Republicans and Democrats) about the severity of the crisis and the efforts to reduce disease transmission [89]. Besides this cluster, it

is important to remember that in parallel to the first year and first waves of the COVID-19 pandemic, 2020 was an election year. Thus, the local and national management of this global epidemic was a source of political debates, and support or criticism of governments, administrations, and the health care system.

The mechanisms behind the folksonomy rely on a complex set of factors. First, as pointed out above, the reasons for the emergence of each cluster depend on both culture and real-life events. Second, these mechanisms can be quantified by analyzing terms that frequently appear together (n-grams). Thus, in the context of this research, we observed that the main focus of the tweets related mainly to COVID-19 pandemic events (disease, confinement, politician talks, vaccines approval, and vaccination) and increased over time, like the prevalence of the terms "vaccines" and "vaccination," and this was in contrast with the term "flu," which disappeared over time from the tweets. This reflects that COVID-19 measures, such as social distancing and mask wearing, significantly reduced the seasonal influenza rates in 2020-2021 [73,90,91]. However, a potential major reason and mechanism of these changes in trends and therefore of the folksonomy content may be associated with the diversion of citizens' attention to annual influenza spread, caused by the disruptive and menacing COVID-19 pandemic. These distractions induced different behaviors or feelings, such as devastation, fear, worry, and the need to understand [92,93].

## Strengths and Limitations

Social media and social networks are increasingly being used to disseminate multimodal and multisource-based health-related information in a timely manner. In the context of epidemics and pandemics, such as seasonal influenza and COVID-19, health care organizations and governmental institutions nowadays spread information and run communication campaigns on social media, for example, to increase citizen engagement in vaccination. At the same time, individuals share their positions, even if it is associated with the antivax trend, and sometimes spread misinformation [94]. The strength of our study is its ability to provide health authorities with a weekly, monthly, and long-term folksonomy of the emerging or persisting topics of social media threads related to a health care issue or event, such as vaccination or a virus-related matter. Providing a folksonomy and the co-occurring terms in the same or additional clusters, using these tools, can enhance health-related social media campaigns, focusing on grand public in-time interests and queries, similar to the approaches used in other business fields.

By getting reports in a timely manner, it has been proven possible to point out the various topics, words, and terms frequently used on social media, thereby enabling health communication specialists, and more specifically those dealing with social media, to focus on up-to-date campaigns to increase population engagement, such as that done in other business fields [95], and actions related to health promotion, especially during epidemics and crises [96] (eg, H1N1 [97] and Ebola [98]), as has been suggested in prior research not dealing with terms, topics, and target population discovery or designation [99].

Exploring social media, and more particularly social networks, is limited by the passive exclusion of nonusers of these communication channels or inactive users who only read posts but do not post by themselves or respond to the messages of other users.

Another limitation of this study is that it was based only on tweets in English and posted from North America. This filtering limits the generalization of the results. The diversity of the US population suggests that running this kind of study in the United States in other languages will enable fine-tuning of health communication and increase vaccination compliance in non-English speaking communities (ie, around 22.0% of the US population) [19,100].

In parallel with our study, another study dealing specifically and strictly with vaccination and COVID-19 was performed among Australian Twitter users (versus US Twitter users in our study) between January and October 2020 (versus between December 2019 and April 2021 in our study) and collected 31,100 tweets (versus 2,782,720 tweets collected by us). The analysis was based on latent Dirichlet allocation, which is an unsupervised learning approach that can be large-scale intensively system resource consuming [101]. The Australian tweet analysis revealed the following 3 dominant topics: (1) "COVID-19 and its vaccination," (2) "advocacy for infection control measures and vaccine trials," and (3) "conspiracy theories, complaints, and misinformation" [102]. Even though some convergence exists, these results are distinct from ours by focusing more specifically on COVID-19–related issues.

Moreover, the set of words initially used for extracting the tweets ("influenza" OR "vaccine" OR "vaccination" OR "vaxx") allowed us to capture a larger spectrum of threads related to each one of the terms that we were interested in focusing on and not in a strict filtering approach, as in other prior research [101]. Nevertheless, without extending the extraction word set, with terms of the COVID-19 pandemic, tweets potentially interesting but not comprising one of these terms would not have been extracted. For example, the following tweet published in mid-April 2021 that included words detected in the n-gram analysis but not explicitly the words used for the tweet's extraction failed to be retrieved: "I am excited, I am in my county seat to get my first injection of the Pfizer." A future perspective for enhancing the dynamic of trend tracking can be considered to update the terms of the tweet extraction query with other disrupting terms due to actuality (eg, "covid," "dose," "injection," and trade names of vaccines). This enhancement can be achieved by a domain expert (ie, human action) or by automatically selecting words emerging as trending in a cluster of the folksonomy and co-occurrence frequency analysis (ie, n-grams) [95].

Additionally, when dealing with the large volume of tweets generated each minute, looking at all tweets in real time is impossible without deploying a high computational infrastructure, which is available in dedicated centers. Accordingly, the objective of this research was to define a framework enabling health system decision makers to focus on specific issues in order to enhance their social media campaigns by understanding the topics discussed in a particular context

(ie, vaccination and influenza). Furthermore, the tweets are collected daily (due to Twitter constraints, without using a paying platform) and analyzed, with the machine learning flow described in the methodology, at the weekly, monthly, and all-time levels. To deal with others' terms of interest, changing the terms of the tweet extraction query will allow the expansion of the current data set or the start of new research with the same methodology. This study shows that combining social media data, such as tweets, and artificial intelligence approaches, such as machine learning algorithms for text and data mining, enables an infodemiology and infoveillance study as a whole. More specifically, in this study, we noticed the strength of this combined approach by following the changes in the contents and topics of the tweets over time and the influence of the actual events. Like other Twitter-based public health research, the approach of collecting, analyzing, and assessing in near real time the content of messages provides powerful indications to health decision makers for adapting and enhancing communication as an emergency response and in planning [103]. In other words, these forewarnings must support social media–based health information in targeting advertisements of recommendations, instructions, and directives, according to social media user' interests and focuses (ie, terms appearing in the clusters of the folksonomy) disclosed passively in previous posts, shares, or likes. Moreover, social media platforms allow accurate targeting by stratifying advertising campaigns on sociodemographic attributes, such as age, gender, marital status, location, spoken language, and educational and professional background [104]. Thus, social media–based health information is intended to increase population adherence to health policies, such as vaccination against epidemic or pandemic diseases (eg, influenza and COVID-19), by delivering personalized messages taking into account both sociodemographics and domains of interest. For example, a young person playing basketball, living in an area with recurrent high acute influenza incidence in a young population, following social media groups dealing with basketball, and sharing posts related to vaccination hesitancy will get advertisements with personalized content targeting young vaccination-hesitant individuals playing collective sports

and emphasizing that vaccination is the best solution to continue this activity during an epidemic [105].

## Conclusions

Twitter is one of the leading social network platforms allowing anyone to share positions and information in any domain. Therefore, any kind of information published and spread about influenza and COVID-19, and the vaccines against each, can be perceived as reliable and can influence social media users. Specifically, during the COVID-19 pandemic, world leaders have widely used Twitter to communicate public health information with citizens. These messages had a strong effect on vaccination compliance [106], with the ability to dynamically improve the content and target health communication campaigns on social media.

This study allowed us to validate our initial hypothesis. Tweets are a source of information for understanding why it is recommended to take a vaccine and the public perception about it [107-109]. Indeed, we defined a folksonomy of the 3 main topics coexisting in the collected messages over 16 months. Accordingly, the terms and hashtags of tweets concerning "influenza," "vaccines," and "vaccination" can be organized in a dynamic vocabulary, such as a folksonomy, reflecting the main topics and their terms discussed over time on the social media platform. Additionally, the emergence and dominance of terms related to COVID-19 over time, reported in the folksonomy with frequently co-occurring words, shows that although the study did not initially focus on this thematic, the health changes are reflected in the Twitter threads related to vaccines and vaccination.

This study focused initially on vaccination against influenza and moved to vaccination against COVID-19. Infoveillance on Twitter (and other social media) about the topics related to vaccines and vaccination against communicable diseases can create opportunities to design and convey personalized messages encouraging specific targeted subpopulations' engagement in vaccination. A greater likelihood that a targeted population receives a personalized message is associated with a higher response, engagement, and proactiveness of the target population for vaccination or other public health measures [110].

## Authors' Contributions

All authors attest that they meet the International Committee of Medical Journal Editors criteria for authorship, have reviewed the manuscript version to be submitted, and agreed with its content and submission. AB was responsible for project supervision; the conception, design, and conduct of the study; the preparation and submission of the relevant documents to the ethics committee; data analysis; data interpretation; writing of the first draft of the manuscript; and critical review and revision of the manuscript for important intellectual content. AC is a Master of Science in Technology Management student, Holon Institute of Technology, Israel (under the supervision of AB), and was responsible for the conception, design, and conduct of the study; data collection; data curation; data analysis; data interpretation; writing of the first draft of the manuscript; and critical review and revision of the manuscript for important intellectual content. EL was responsible for data interpretation, and critical review and revision of the manuscript for important intellectual content. SA was responsible for data analysis, data interpretation, and critical review and revision of the manuscript for important intellectual content.

[XSL•FO]

**RenderX**

## Conflicts of Interest

None declared.

---

Multimedia Appendix 1

Number of tweets by month comprising at least one of the terms "flu," "vaccination," "vaccine," "vaxx," and "covid.".

[PDF File (Adobe PDF File), 44 KB - infodemiology_v1i1e31983_app1.pdf ]

---

Multimedia Appendix 2

Optimal number of clusters, the related maximal silhouette score for each month, and the parameters used for creating the topic clustering.

[PDF File (Adobe PDF File), 53 KB - infodemiology_v1i1e31983_app2.pdf ]

---

Multimedia Appendix 3

List of the 1000 most frequent n-grams in the 3 clusters.

[PDF File (Adobe PDF File), 99 KB - infodemiology_v1i1e31983_app3.pdf ]

---

Multimedia Appendix 4

N-grams having the highest increase in occurrence week over week.

[PDF File (Adobe PDF File), 53 KB - infodemiology_v1i1e31983_app4.pdf ]

---

## References

1. Bello-Orgaz G, Hernandez-Castro J, Camacho D. Detecting discussion communities on vaccination in twitter. Future Generation Computer Systems 2017 Jan;66:125-136. [doi: 10.1016/j.future.2016.06.032]
2. Grajales FJ, Sheps S, Ho K, Novak-Lauscher H, Eysenbach G. Social media: a review and tutorial of applications in medicine and health care. J Med Internet Res 2014 Feb 11;16(2):e13 [FREE Full text] [doi: 10.2196/jmir.2912] [Medline: 24518354]
3. Choi S. The Two-Step Flow of Communication in Twitter-Based Public Forums. Social Science Computer Review 2014 Nov 07;33(6):696-711. [doi: 10.1177/0894439314556599]
4. Mosleh M, Pennycook G, Arechar AA, Rand DG. Cognitive reflection correlates with behavior on Twitter. Nat Commun 2021 Feb 10;12(1):921 [FREE Full text] [doi: 10.1038/s41467-020-20043-0] [Medline: 33568667]
5. Trye D, Calude AS, Bravo-Marquez F, Keegan TT. Hybrid Hashtags: #YouKnowYoureAKiwiWhen Your Tweet Contains Māori and English. Front Artif Intell 2020;3:15 [FREE Full text] [doi: 10.3389/frai.2020.00015] [Medline: 33733134]
6. Inmon WH. Data Architecture: a Primer for the Data Scientist || The Data Infrastructure. Cambridge, MA: Academic Press; 2015.
7. Big Data. Gartner. URL: https://www.gartner.com/en/information-technology/glossary/big-data [accessed 2021-05-02]
8. Zhou S, Qiao Z, Du Q, Wang GA, Fan W, Yan X. Measuring Customer Agility from Online Reviews Using Big Data Text Analytics. Journal of Management Information Systems 2018 May 15;35(2):510-539. [doi: 10.1080/07421222.2018.1451956]
9. Abbas A, Zhou Y, Deng S, Zhang P. Text Analytics to Support Sense-Making in Social Media: A Language-Action Perspective. MISQ 2018 Feb 2;42(2):427-464. [doi: 10.25300/MISQ/2018/13239]
10. Hoogeveen D, Wang L, Baldwin T, Verspoor K. Web Forum Retrieval and Text Analytics: A Survey. FNT in Information Retrieval 2018;12(1):1-163. [doi: 10.1561/1500000062]
11. Mokkink LB, Terwee CB, Patrick DL, Alonso J, Stratford PW, Knol DL, et al. The COSMIN study reached international consensus on taxonomy, terminology, and definitions of measurement properties for health-related patient-reported outcomes. J Clin Epidemiol 2010 Jul;63(7):737-745. [doi: 10.1016/j.jclinepi.2010.02.006] [Medline: 20494804]
12. Vrijens B, De Geest S, Hughes DA, Przemyslaw K, Demonceau J, Ruppar T, et al. A new taxonomy for describing and defining adherence to medications. Br J Clin Pharmacol 2012 May;73(5):691-705 [FREE Full text] [doi: 10.1111/j.1365-2125.2012.04167.x] [Medline: 22486599]
13. Šmite D, Wohlin C, Galviņa Z, Prikladnicki R. An empirically based terminology and taxonomy for global software engineering. Empir Software Eng 2012 Jul 18;19(1):105-153. [doi: 10.1007/s10664-012-9217-9]
14. Zimbra D, Abbasi A, Zeng D, Chen H. The State-of-the-Art in Twitter Sentiment Analysis. ACM Trans. Manage. Inf. Syst 2018 Sep 05;9(2):1-29. [doi: 10.1145/3185045]
15. de Lusignan S, Liyanage H, McGagh D, Jani BD, Bauwens J, Byford R, et al. COVID-19 Surveillance in a Primary Care Sentinel Network: In-Pandemic Development of an Application Ontology. JMIR Public Health Surveill 2020 Nov 17;6(4):e21434 [FREE Full text] [doi: 10.2196/21434] [Medline: 33112762]
16. Immunization coverage. World Health Organization. URL: https://www.who.int/en/news-room/fact-sheets/detail/immunization-coverage [accessed 2021-05-30]

XSL·FO
RenderX

17. Ashkenazi S, Livni G, Klein A, Kremer N, Havlin A, Berkowitz O. The relationship between parental source of information and knowledge about measles / measles vaccine and vaccine hesitancy. Vaccine 2020 Oct 27;38(46):7292-7298. [doi: 10.1016/j.vaccine.2020.09.044] [Medline: 32981777]

18. Benis A, Khodos A, Ran S, Levner E, Ashkenazi S. Social Media Engagement and Influenza Vaccination During the COVID-19 Pandemic: Cross-sectional Survey Study. J Med Internet Res 2021 Mar 16;23(3):e25977 [FREE Full text] [doi: 10.2196/25977] [Medline: 33651709]

19. Benis A, Seidmann A, Ashkenazi S. Reasons for Taking the COVID-19 Vaccine by US Social Media Users. Vaccines (Basel) 2021 Mar 29;9(4) [FREE Full text] [doi: 10.3390/vaccines9040315] [Medline: 33805283]

20. Deiner MS, Fathy C, Kim J, Niemeyer K, Ramirez D, Ackley SF, et al. Facebook and Twitter vaccine sentiment in response to measles outbreaks. Health Informatics J 2019 Sep 01;25(3):1116-1132 [FREE Full text] [doi: 10.1177/1460458217740723] [Medline: 29148313]

21. WHO Coronavirus (COVID-19) Dashboard. World Health Organization. URL: https://covid19.who.int/ [accessed 2021-05-30]

22. Randolph HE, Barreiro LB. Herd Immunity: Understanding COVID-19. Immunity 2020 May 19;52(5):737-741 [FREE Full text] [doi: 10.1016/j.immuni.2020.04.012] [Medline: 32433946]

23. McDermott A. Core Concept: Herd immunity is an important-and often misunderstood-public health phenomenon. Proc Natl Acad Sci U S A 2021 May 25;118(21) [FREE Full text] [doi: 10.1073/pnas.2107692118] [Medline: 34011611]

24. Rainie L, Wellman B. Networked: The New Social Operating System. Cambridge, MA: The MIT Press; 2014.

25. Wu L, Morstatter F, Carley KM, Liu H. Misinformation in Social Media. SIGKDD Explor. Newsl 2019 Nov 26;21(2):80-90. [doi: 10.1145/3373464.3373475]

26. Eysenbach G. Infodemiology: The epidemiology of (mis)information. Am J Med 2002 Dec 15;113(9):763-765. [doi: 10.1016/s0002-9343(02)01473-0] [Medline: 12517369]

27. Morley J, Cowls J, Taddeo M, Floridi L. Public Health in the Information Age: Recognizing the Infosphere as a Social Determinant of Health. J Med Internet Res 2020 Aug 03;22(8):e19311 [FREE Full text] [doi: 10.2196/19311] [Medline: 32648850]

28. Cordina M, Lauri MA, Lauri J. Attitudes towards COVID-19 vaccination, vaccine hesitancy and intention to take the vaccine. Pharm Pract (Granada) 2021;19(1):2317 [FREE Full text] [doi: 10.18549/PharmPract.2021.1.2317] [Medline: 33828623]

29. MacDonald NE, SAGE Working Group on Vaccine Hesitancy. Vaccine hesitancy: Definition, scope and determinants. Vaccine 2015 Aug 14;33(34):4161-4164 [FREE Full text] [doi: 10.1016/j.vaccine.2015.04.036] [Medline: 25896383]

30. Gupta A, Katarya R. Social media based surveillance systems for healthcare using machine learning: A systematic review. J Biomed Inform 2020 Aug;108:103500 [FREE Full text] [doi: 10.1016/j.jbi.2020.103500] [Medline: 32622833]

31. Thomson A, Vallée-Tourangeau G, Suggs LS. Strategies to increase vaccine acceptance and uptake: From behavioral insights to context-specific, culturally-appropriate, evidence-based communications and interventions. Vaccine 2018 Oct 22;36(44):6457-6458 [FREE Full text] [doi: 10.1016/j.vaccine.2018.08.031] [Medline: 30201305]

32. Alessa A, Faezipour M. A review of influenza detection and prediction through social networking sites. Theor Biol Med Model 2018 Feb 01;15(1):2 [FREE Full text] [doi: 10.1186/s12976-017-0074-5] [Medline: 29386017]

33. Aramaki E, Maskawa S, Morita M. Twitter catches the flu: detecting influenza epidemics using Twitter. In: EMNLP '11: Proceedings of the Conference on Empirical Methods in Natural Language Processing. 2011 Presented at: Conference on Empirical Methods in Natural Language Processing; July 27-31, 2011; Edinburgh, United Kingdom p. 1568-1576 URL: https://dl.acm.org/doi/10.5555/2145432.2145600

34. Talvis K, Chorianopoulos K, Kermanidis K. Real-Time Monitoring of Flu Epidemics through Linguistic and Statistical Analysis of Twitter Messages. 2014 Presented at: 9th International Workshop on Semantic and Social Media Adaptation and Personalization; November 6-7, 2014; Corfu, Greece URL: https://ieeexplore.ieee.org/document/6978958

35. Wakamiya S, Kawai Y, Aramaki E. Twitter-Based Influenza Detection After Flu Peak via Tweets With Indirect Information: Text Mining Study. JMIR Public Health Surveill 2018 Sep 25;4(3):e65 [FREE Full text] [doi: 10.2196/publichealth.8627] [Medline: 30274968]

36. Hassan Zadeh A, Zolbanin HM, Sharda R, Delen D. Social Media for Nowcasting Flu Activity: Spatio-Temporal Big Data Analysis. Inf Syst Front 2019 Jan 5;21(4):743-760. [doi: 10.1007/s10796-018-9893-0]

37. Faasse K, Chatman CJ, Martin LR. A comparison of language use in pro- and anti-vaccination comments in response to a high profile Facebook post. Vaccine 2016 Nov 11;34(47):5808-5814. [doi: 10.1016/j.vaccine.2016.09.029] [Medline: 27707558]

38. Sturm L, Kasting ML, Head KJ, Hartsock JA, Zimet GD. Influenza vaccination in the time of COVID-19: A national U.S. survey of adults. Vaccine 2021 Apr 01;39(14):1921-1928 [FREE Full text] [doi: 10.1016/j.vaccine.2021.03.003] [Medline: 33715898]

39. Gandomi A, Haider M. Beyond the hype: Big data concepts, methods, and analytics. International Journal of Information Management 2015 Apr;35(2):137-144. [doi: 10.1016/j.ijinfomgt.2014.10.007]

40. Secinaro S, Calandra D, Secinaro A, Muthurangu V, Biancone P. The role of artificial intelligence in healthcare: a structured literature review. BMC Med Inform Decis Mak 2021 Apr 10;21(1):125 [FREE Full text] [doi: 10.1186/s12911-021-01488-9] [Medline: 33836752]

41. Schraeder TL. Physician Communication Connecting with Patients, Peers, and the Public. Oxford, United Kingdom: Oxford University Press; 2019.

42. Benis A, Barak Barkan R, Sela T, Harel N. Communication Behavior Changes Between Patients With Diabetes and Healthcare Providers Over 9 Years: Retrospective Cohort Study. J Med Internet Res 2020 Aug 11;22(8):e17186 [FREE Full text] [doi: 10.2196/17186] [Medline: 32648555]

43. Dai X, Bikdash M, Meyer B. From social media to public health surveillance: Word embedding based clustering method for twitter classification. 2017 Presented at: SoutheastCon 2017; March 30-April 2, 2017; Concord, NC p. 1-7. [doi: 10.1109/secon.2017.7925400]

44. Henrich NJ. Increasing pandemic vaccination rates with effective communication. Hum Vaccin 2011 Jun;7(6):663-666 [FREE Full text] [doi: 10.4161/hv.7.6.15007] [Medline: 21445004]

45. Feemster KA. Building vaccine acceptance through communication and advocacy. Hum Vaccin Immunother 2020 May 03;16(5):1004-1006 [FREE Full text] [doi: 10.1080/21645515.2020.1746603] [Medline: 32401681]

46. Sinclair J, Cardew-Hall M. The folksonomy tag cloud: when is it useful? Journal of Information Science 2007 May 31;34(1):15-29. [doi: 10.1177/0165551506078083]

47. Robu V, Halpin H, Shepherd H. Emergence of consensus and shared vocabularies in collaborative tagging systems. ACM Trans. Web 2009 Sep;3(4):1-34. [doi: 10.1145/1594173.1594176]

48. Wetzker R, Zimmermann C, Bauckhage C, Albayrak S. I tag, you tag: translating tags for advanced user models. In: WSDM '10: Proceedings of the Third ACM International Conference on Web Search and Data Mining. 2010 Presented at: Third ACM International Conference on Web Search and Data Mining; February 4-6, 2010; New York, NY p. 71-80. [doi: 10.1145/1718487.1718497]

49. Hönings H, Knapp D, Nguy n BC, Richter D, Williams K, Dorsch I, et al. Health information diffusion on Twitter: The content and design of WHO tweets matter. Health Info Libr J 2021 Mar 08. [doi: 10.1111/hir.12361] [Medline: 33682996]

50. Sacco G. Dynamic taxonomies: a model for large information bases. IEEE Trans. Knowl. Data Eng 2000;12(3):468-479. [doi: 10.1109/69.846296]

51. Twitter API. Twitter Developer Platform. URL: https://developer.twitter.com/en/docs/twitter-api [accessed 2021-05-04]

52. Zanin M, Aitya NAA, Basilio J, Baumbach J, Benis A, Behera CK, et al. An Early Stage Researcher's Primer on Systems Medicine Terminology. Netw Syst Med 2021 Feb;4(1):2-50 [FREE Full text] [doi: 10.1089/nsm.2020.0003] [Medline: 33659919]

53. Russell S, Norvig P. Artificial Intelligence: A Modern Approach. New York, NY: Pearson; 2020.

54. Straus J, Kaufman L, Stern T. The Blue Book of Grammar and Punctuation: An Easy-to-Use Guide with Clear Rules, Real-World Examples, and Reproducible Quizzes. Hoboken, NJ: Wiley; 2014.

55. Bird S, Klein E, Loper E. Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit. Sebastopol, CA: O'Reilly Media; 2009.

56. Bergmanis T, Goldwater S. Context Sensitive Neural Lemmatization with Lematus. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers). 2018 Presented at: Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies; 2018; New Orleans, LO p. 1391-1400. [doi: 10.18653/v1/N18-1126]

57. Aggarwal C, Reddy C. Data Clustering: Algorithms and Applications. Boca Raton, FL: Chapman and Hall/CRC; 2014.

58. Keogh E, Mueen A. Curse of Dimensionality. In: Sammut C, Webb GI, editors. Encyclopedia of Machine Learning and Data Mining. Boston, MA: Springer; 2017.

59. Vlachos M. Dimensionality Reduction. In: Sammut C, Webb GI, editors. Encyclopedia of Machine Learning and Data Mining. Boston, MA: Springer; 2017.

60. Mikolov T, Chen K, Corrado G, Dean J. Efficient Estimation of Word Representations in Vector Space. arXiv. 2013. URL: https://arxiv.org/abs/1301.3781 [accessed 2021-10-02]

61. gensim 4.1.2. Python Package Index (PyPI). URL: https://pypi.org/project/gensim/ [accessed 2021-06-04]

62. Butnaru A, Ionescu R. From Image to Text Classification: A Novel Approach based on Clustering Word Embeddings. Procedia Computer Science 2017;112:1783-1792. [doi: 10.1016/j.procs.2017.08.211]

63. Li B, Drozd A, Guo Y, Liu T, Matsuoka S, Du X. Scaling Word2Vec on Big Corpus. Data Sci. Eng 2019 Jun 25;4(2):157-175. [doi: 10.1007/s41019-019-0096-6]

64. Rousseeuw PJ. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. Journal of Computational and Applied Mathematics 1987 Nov;20:53-65. [doi: 10.1016/0377-0427(87)90125-7]

65. van der Maaten L, Hinton G. Visualizing Data using t-SNE. Journal of Machine Learning Research 2008;9(86):2579-2605 [FREE Full text]

66. Google Trends. URL: https://trends.google.com/trends/ [accessed 2021-06-13]

67. Matschinske J, Alcaraz N, Benis A, Golebiewski M, Grimm DG, Heumos L, et al. The AIMe registry for artificial intelligence in biomedical research. Nat Methods 2021 Aug 25. [doi: 10.1038/s41592-021-01241-0] [Medline: 34433960]

68. Hoffman BL, Felter EM, Chu K, Shensa A, Hermann C, Wolynn T, et al. It's not all about autism: The emerging landscape of anti-vaccination sentiment on Facebook. Vaccine 2019 Apr 10;37(16):2216-2223. [doi: 10.1016/j.vaccine.2019.03.003] [Medline: 30905530]

XSL·FO
RenderX

69. Burki T. Vaccine misinformation and social media. The Lancet Digital Health 2019 Oct;1(6):e258-e259. [doi: 10.1016/S2589-7500(19)30136-0]

70. Ahmed I. Dismantling the anti-vaxx industry. Nat Med 2021 Mar;27(3):366. [doi: 10.1038/s41591-021-01260-6] [Medline: 33723446]

71. Lopreite M, Panzarasa P, Puliga M, Riccaboni M. Early warnings of COVID-19 outbreaks across Europe from social media. Sci Rep 2021 Jan 25;11(1):2147 [FREE Full text] [doi: 10.1038/s41598-021-81333-1] [Medline: 33495534]

72. Weekly U.S. Influenza Surveillance Report. Centers for Disease Control and Prevention. URL: https://www.cdc.gov/flu/weekly/index.htm [accessed 2021-06-21]

73. Uyeki TM, Wentworth DE, Jernigan DB. Influenza Activity in the US During the 2020-2021 Season. JAMA 2021 Jun 08;325(22):2247-2248. [doi: 10.1001/jama.2021.6125] [Medline: 34028492]

74. Selecting the number of clusters with silhouette analysis on KMeans clustering. Scikit-learn developers. URL: https://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html [accessed 2021-07-27]

75. NLTK 3.6 documentation. NLTK. URL: https://www.nltk.org/_modules/nltk/cluster/kmeans.html [accessed 2021-06-20]

76. Capó M, Pérez A, Lozano JA. An efficient approximation to the K-means clustering for massive data. Knowledge-Based Systems 2017 Feb;117:56-69. [doi: 10.1016/j.knosys.2016.06.031]

77. Škrlj B, Kralj J, Lavrač N. Embedding-based Silhouette community detection. Mach Learn 2020;109(11):2161-2193 [FREE Full text] [doi: 10.1007/s10994-020-05882-8] [Medline: 33191975]

78. Lovmar L, Ahlford A, Jonsson M, Syvänen AC. Silhouette scores for assessment of SNP genotype clusters. BMC Genomics 2005 Mar 10;6:35 [FREE Full text] [doi: 10.1186/1471-2164-6-35] [Medline: 15760469]

79. Maugeri A, Barchitta M, Agodi A. A Clustering Approach to Classify Italian Regions and Provinces Based on Prevalence and Trend of SARS-CoV-2 Cases. Int J Environ Res Public Health 2020 Jul 22;17(15) [FREE Full text] [doi: 10.3390/ijerph17155286] [Medline: 32707989]

80. Ray S, Turi RH. Determination of number of clusters in K-means clustering and application in colour segmentation. 1999 Presented at: 4th International Conference on Advances in Pattern Recognition and Digital Techniques (ICAPRDT'99); December 28-31, 1999; Calcutta, India p. 137-143.

81. Thorndike RL. Who belongs in the family? Psychometrika 1953 Dec;18(4):267-276. [doi: 10.1007/BF02289263]

82. Pfizer and BioNTech announce vaccine candidate against COVID-19 achieved success in first interim analysis from phase 3 study. Pfizer. 2020 Nov. URL: https://www.pfizer.com/news/press-release/press-release-detail/pfizer-and-biontech-announce-vaccine-candidate-against [accessed 2021-06-27]

83. U.S. administers 293.7 mln doses of COVID-19 vaccines -CDC. Reuters. 2021 May. URL: https://www.reuters.com/world/us/us-administers-2937-mln-doses-covid-19-vaccines-cdc-2021-05-29/ [accessed 2021-06-27]

84. pytrends 4.7.3. Python Package Index (PyPI). URL: https://pypi.org/project/pytrends/ [accessed 2021-06-14]

85. Reporting COVID-19 Vaccinations in the United States. Centers for Disease Control and Prevention. URL: https://www.cdc.gov/coronavirus/2019-ncov/vaccines/reporting-vaccinations.html [accessed 2021-06-14]

86. Dong E, Du H, Gardner L. An interactive web-based dashboard to track COVID-19 in real time. The Lancet Infectious Diseases 2020 May;20(5):533-534 [FREE Full text] [doi: 10.1016/S1473-3099(20)30120-1]

87. Baumgaertner B, Carlisle JE, Justwan F. The influence of political ideology and trust on willingness to vaccinate. PLoS One 2018 Jan 25;13(1):e0191728 [FREE Full text] [doi: 10.1371/journal.pone.0191728] [Medline: 29370265]

88. Debus M, Tosun J. Political ideology and vaccination willingness: implications for policy design. Policy Sci 2021 Jun 16:1-15 [FREE Full text] [doi: 10.1007/s11077-021-09428-0] [Medline: 34149102]

89. Allcott H, Boxell L, Conway J, Gentzkow M, Thaler M, Yang D. Polarization and public health: Partisan differences in social distancing during the coronavirus pandemic. J Public Econ 2020 Nov;191:104254 [FREE Full text] [doi: 10.1016/j.jpubeco.2020.104254] [Medline: 32836504]

90. Servick K. COVID-19 measures also suppress flu-for now. Science 2021 Jan 15;371(6526):224. [doi: 10.1126/science.371.6526.224] [Medline: 33446538]

91. Feng L, Zhang T, Wang Q, Xie Y, Peng Z, Zheng J, et al. Impact of COVID-19 outbreaks and interventions on influenza in China and the United States. Nat Commun 2021 May 31;12(1):3249 [FREE Full text] [doi: 10.1038/s41467-021-23440-1] [Medline: 34059675]

92. Razai M, Doerholt K, Ladhani S, Oakeshott P. Coronavirus disease 2019 (covid-19): a guide for UK GPs. BMJ 2020 Mar 05;368:m800 [FREE Full text] [doi: 10.1136/bmj.m800] [Medline: 32144127]

93. Bagus P, Peña-Ramos JA, Sánchez-Bayón A. COVID-19 and the Political Economy of Mass Hysteria. Int J Environ Res Public Health 2021 Feb 03;18(4) [FREE Full text] [doi: 10.3390/ijerph18041376] [Medline: 33546144]

94. Ortiz-Sánchez E, Velando-Soriano A, Pradas-Hernández L, Vargas-Román K, Gómez-Urquiza JL, Cañadas-De la Fuente GA, et al. Analysis of the Anti-Vaccine Movement in Social Networks: A Systematic Review. Int J Environ Res Public Health 2020 Jul 27;17(15) [FREE Full text] [doi: 10.3390/ijerph17155394] [Medline: 32727024]

95. US 2015-0235246 A1 - Cross-channel audience segmentation. Patent Center. 2015. URL: https://patentcenter.uspto.gov/#!/applications/14623738 [accessed 2021-05-02]

96. Wendling C, Radisch J, Jacobzone S. The Use of Social Media in Risk and Crisis Communication. OECD Working Papers on Public Governance 2013;24. [doi: 10.1787/5k3v01fskp9s-en]

97. Freberg K, Palenchar MJ, Veil SR. Managing and sharing H1N1 crisis information using social media bookmarking services. Public Relations Review 2013 Sep;39(3):178-184. [doi: 10.1016/j.pubrev.2013.02.007]

98. Guidry JP, Jin Y, Orr CA, Messner M, Meganck S. Ebola on Instagram and Twitter: How health organizations address the health crisis in their social media engagement. Public Relations Review 2017 Sep;43(3):477-486. [doi: 10.1016/j.pubrev.2017.04.009]

99. Ozawa S, Clark S, Portnoy A, Grewal S, Stack ML, Sinha A, et al. Estimated economic impact of vaccinations in 73 low- and middle-income countries, 2001-2020. Bull World Health Organ 2017 Sep 01;95(9):629-638 [FREE Full text] [doi: 10.2471/BLT.16.178475] [Medline: 28867843]

100. Selected social characteristics in the United States. United States Census Bureau. URL: https://data.census.gov/cedsci/table?tid=ACSDP5Y2019.DP02 [accessed 2021-06-24]

101. Xie X, Liang Y, Li X, Tan W. CuLDA: Solving Large-scale LDA Problems on GPUs. In: HPDC '19: Proceedings of the 28th International Symposium on High-Performance Parallel and Distributed Computing. 2019 Presented at: 28th International Symposium on High-Performance Parallel and Distributed Computing; June 22-29, 2019; New York, NY p. 195-205. [doi: 10.1145/3307681.3325407]

102. Kwok SWH, Vadde SK, Wang G. Tweet Topics and Sentiments Relating to COVID-19 Vaccination Among Australian Twitter Users: Machine Learning Analysis. J Med Internet Res 2021 May 19;23(5):e26953 [FREE Full text] [doi: 10.2196/26953] [Medline: 33886492]

103. Xue J, Chen J, Hu R, Chen C, Zheng C, Su Y, et al. Twitter Discussions and Emotions About the COVID-19 Pandemic: Machine Learning Approach. J Med Internet Res 2020 Nov 25;22(11):e20550 [FREE Full text] [doi: 10.2196/20550] [Medline: 33119535]

104. US-20150088636-A1 - Classification of Geographic Performance Data. Patent Center. 2015. URL: https://patentcenter.uspto.gov/#!/applications/14555758 [accessed 2021-05-02]

105. US 2014-0236715 A1 - Targeted Advertising in Social Media Networks. Patent Center. 2014. URL: https://patentcenter.uspto.gov/#!/applications/14036494 [accessed 2021-05-02]

106. Rufai S, Bunce C. World leaders' usage of Twitter in response to the COVID-19 pandemic: a content analysis. J Public Health (Oxf) 2020 Aug 18;42(3):510-516 [FREE Full text] [doi: 10.1093/pubmed/fdaa049] [Medline: 32309854]

107. Read W, Robertson N, McQuilken L, Ferdous A. Consumer engagement on Twitter: perceptions of the brand matter. EJM 2019 Sep 09;53(9):1905-1933. [doi: 10.1108/ejm-10-2017-0772]

108. Dyer J, Kolic B. Public risk perception and emotion on Twitter during the Covid-19 pandemic. Appl Netw Sci 2020;5(1):99 [FREE Full text] [doi: 10.1007/s41109-020-00334-7] [Medline: 33344760]

109. Saleh SN, Lehmann CU, McDonald SA, Basit MA, Medford RJ. Understanding public perception of coronavirus disease 2019 (COVID-19) social distancing on Twitter. Infect Control Hosp Epidemiol 2021 Feb;42(2):131-138 [FREE Full text] [doi: 10.1017/ice.2020.406] [Medline: 32758315]

110. Benis A, Tamburis O, Chronaki C, Moen A. One Digital Health: A Unified Framework for Future Health Ecosystems. J Med Internet Res 2021 Feb 05;23(2):e22189 [FREE Full text] [doi: 10.2196/22189] [Medline: 33492240]

## Abbreviations

**API:** application programming interface
**FDA:** Food and Drug Administration
**NLTK:** Natural Language Toolkit
**t-SNE:** t-distributed stochastic neighbor embedding

<u>Corrigenda and Addenda</u>

# Metadata Correction: "Desensitization to Fear-Inducing COVID-19 Health News on Twitter: Observational Study"

Hannah R Stevens[1], BA; Yoo Jung Oh[1], MA; Laramie D Taylor[1], PhD

Department of Communication, University of California, Davis, Davis, CA, United States

**Corresponding Author:**
Hannah R Stevens, BA
Department of Communication
University of California, Davis
1 Shields Ave
Davis, CA, 95616
United States
Phone: 1 530 752 1011
Email: hrstevens@ucdavis.edu

**Related Article:**

Correction of: https://infodemiology.jmir.org/2021/1/e26876/

In "Desensitization to Fear-Inducing COVID-19 Health News on Twitter: Observational Study" (JMIR Infodemiology 2021;1(1):e26876) the authors noted one error.

One author's name was displayed as follows:

*Laramie R Taylor*

The middle initial in the name has now been corrected as follows:

*Laramie D Taylor*

The correction will appear in the online version of the paper on the JMIR Publications website on July 28, 2021, together with the publication of this correction notice. Because this was made after submission to PubMed, PubMed Central, and other full-text repositories, the corrected article has also been resubmitted to those repositories.

XSL•FO
**RenderX**