

Very Preliminary Draft

July 5, 2021

Using social media data to verify the platforms' regulation policies regarding "misinformation"

1 Introduction

A number of recent studies point towards the idea that "Fake News" or disinformation is a small subset of the total supply of information on online social networking platforms (e.g. Grinberg et al. (2019) [2] and Broniatowski et al. (2020) [1]). Yet, this seemingly small subset is generating great concern in traditional media and in society in a broader sense.¹

Section 230 in the United States Communications Decency Act² provides immunity for website platforms against the content created by users. Nevertheless, there is growing pressure for Mainstream Social Media Platforms (hereafter MSMP), such as Facebook, Twitter or Youtube, to moderate the available content. In particular, platforms seem to take explicit actions when content is in violation of local laws in different jurisdictions, e.g. laws regarding defamation of a racial nature, dissemination of symbols from unconstitutional organizations, privacy protection, digital security, electoral laws. Facebook reports having implemented a total of 64.7 thousand content restrictions based on local law across all countries in 2020.³ Google reports a total of 26 thousand government requests to remove content from July 2020 to December 2020, among which 11.4 thousand concerned Youtube.⁴ Twitter reports having received 42.2 thousand legal demands from third-parties from January to June 2020, and has responded by withholding 82 thousand accounts and 3.1 thousand tweets.⁵

Furthermore, MSMP are increasingly engaging in editorial tasks by implementing targeted policies to insure that each platform's rules are not violated. Community guidelines of Facebook, Twitter and Youtube can be summarized in a handful of categories, regarding safety, privacy and authenticity; which include violence, terrorism, child sexual exploitation, abuse, harassment, hateful conduct, suicide or self-harm, illegal or regulated goods and services, platform manipulation and

¹For example see the February 2020 speech of the Director General of the WHO at the Munich Security Conference, where he says "But we're not just fighting an epidemic; we're fighting an infodemic."

²Similar regulation exists in the European Union, see articles 12 and 15 of the E-commerce Directive (2000).

³See Facebook Transparency Center, Content restrictions based on Local Law: transparency.fb.com/data/contentrestrictions. We summed the count of content restrictions over all countries reported in the table, for *H1* and *H2* of the year 2020.

⁴See Google's Transparency report, government requests to remove content: transparencyreport.google.com/government-removals/overview.

⁵See Twitter Transparency website, Removal requests: transparency.twitter.com/en/reports/removal-requests.html#2020-jan-jun.

spam.⁶ While specific to each platform, the previously cited categories correspond in most cases to well defined concepts that fall into legal frameworks in many countries.

In this article, we focus on MSMP’s policies and actions regarding content with low credibility or false information, commonly referred to as *Fake News*.⁷ The *Fake News* phenomenon is still ill-defined by the academic community as it encompasses several combined features such as spreading inaccurate, false or misleading information, with or without the intention of influencing or manipulating a target pool of audience. The growth of social networking platforms over the last decade in terms of number of users worldwide and volume of content, has modified the information ecosystem in terms of production of information and its mediation. Many users can now produce and share content which includes news related information, without having to abide by strict editorial processes that ensure accuracy of information and reliability of sources. In particular, false or inaccurate content produced and shared on social networking platforms concerning the political life or public health may have a potentially harmful impact on the society, in the rare event that it goes viral. This gave rise to a set of heterogenous fact-checking policies across mainstream platforms.

During the COVID19 global health pandemic platforms have upgraded their guidelines to include a set of rules to tackle the propagation of potentially harmful content.⁸ As each platform is a private company, those *new* policies are not coordinated and are implemented in different ways across platforms. Such targeted policies show the willingness of MSMP to enhance the quality of the online conversation, but also sheds light on the lack of specific policies to tackle misinformation in general. *Say what they say they do + partnerships with fact-checkers + algorithms + community, idea about academic community needs to be able to study this phenomenon and assess the impact of the policies, what effects they have etc.*

In the present article, we will explain how to verify with data mining MSMP’s actions regarding content with low credibility or false information, through a series of examples for different actions and platforms. For the purpose of clarity, we only focus on three platforms: Facebook, Twitter and Youtube. Both Facebook and Youtube are in the top three most popular social media platforms in terms of number of users.⁹ We further choose Twitter because it is a social networking platform with the most news-focused users, according to the Pew Research center (2019) [3]. More specifically, we survey a number of common policies used in order to tackle misinformation, across the three above cited MSMP, Facebook, Twitter and Youtube: temporary or permanent suspension of users, reducing the visibility of some content, introducing flags and notices. We do not provide an exhaustive list of methods on how to investigate the platforms’ policies. We rather provide a methodology to investigate key policies, that can be useful to researchers or journalists interested in implementing external monitoring. We summarize in table 1 the tools used to collect the data from Facebook, Twitter and Youtube, that we use in multiple examples throughout the present article. Finally, we discuss how an increased effort of transparency regarding specific content can help the community of researchers study and assess the impact of platforms’ policies regarding misinformation.

⁶For an exhaustive overview of the community standards of Facebook, see: facebook.com/communitystandards/. For the Twitter Rules see: help.twitter.com/en/rules-and-policies/twitter-rules, and for Youtube community guidelines see: youtube.com/intl/en_us/howyoutubeworks/policies/communityguidelines/.

⁷For an overview on the *concept* of *Fake News*, we refer the reader to the article The science of fake news by Lazer et al. (2018) [4].

⁸For Facebook, Twitter and Youtube see respectively the following updates:

⁹See for example the ranking of the most popular social networks as of April 2021 on Statista: <https://www.statista.com/statistics/272014/global-social-networks-ranked-by-number-of-users/>.

	Application Programming Interface (API)	Web Scraping
Facebook	CrowdTangle: public content insights tool owned and operated by Facebook. Buzzsumo: commercial content database	
Twitter	Twitter API V2	Minet twitter scrape
Youtube	Youtube API V3	

Table 1: Data collection

2 Policies

2.1 Temporary suspension and Permanent suspension

Mainstream social media platforms may suspend the account of a specific user when they deem that the platforms’ rules have been violated. Account suspension can be temporary or permanent.¹⁰ When the suspension is temporary the user is prohibited for a limited period of time from posting content on their account, but created content prior to suspension remains available to the user and their followers. However, when the suspension is permanent, in most cases, followers or subscribers have no longer access to the content prior to the suspension and the user can no longer use the account to create new content. In what follows, we focus on the implementation of this policy by several platforms and provide simple examples to illustrate.

2.1.1 Facebook

When an account is permanently suspended by Facebook, it disappears from the platform. That is, the data can no longer be scrapped and it also disappears from the CrowdTangle API.¹¹ Facebook publishes on monthly basis a *coordinated inauthentic behavior* report, where it informs how many personal accounts, pages or groups were deleted and to which *deceptive network* they may have belonged.¹² But as long as external persons do not have access to deleted accounts data, these reports cannot be verified by independent researchers or journalists.

Facebook can also apply a temporary suspension, and in this case the data can often be collected and analyzed. For example, Donald Trump’s official Facebook page has been suspended following the Capitol attack on January 6, 2021.¹³ Nevertheless the page’s data is still present in the CrowdTangle API. Thus, after manually adding this page to the CrowdTangle dashboard, we collected the 6 083 posts it had published between January 1, 2020 and June 15, 2021 using the *posts* endpoint.¹⁴ We used the minet command line tool [5] to collect data.¹⁵ We can verify on figure 1 that the *Donald J. Trump* page has not published any content since January 6, 2021, and

¹⁰A list of notable Twitter temporary and permanent suspensions can be found on wikipedia: https://en.wikipedia.org/wiki/Twitter_suspensions.

¹¹CrowdTangle is a public insights tool owned and operated by Facebook, that exclusively tracks public content from Facebook public groups and pages.

¹²See the April 2021 report for an example: <https://about.fb.com/news/2021/05/april-2021-coordinated-inauthentic-behavior-report/>

¹³See <https://www.facebook.com/zuck/posts/10112681480907401>

¹⁴(see the endpoint documentation for more details: <https://github.com/CrowdTangle/API/wiki/Posts>).

¹⁵The exact command can be found : [here](#).

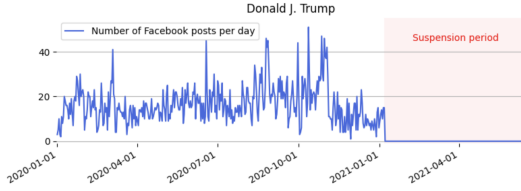


Figure 1: Number of Facebook posts published each day by the Facebook page *Donald J. Trump* between January 1, 2020 and June 15, 2021. The data corresponds to 6 083 posts retrieved from the CrowdTangle API using the *posts* endpoint.

that this behavior is not consistent with the page’s previous activity: an average of 16 posts were published each day on Facebook before the suspension.

2.1.2 Twitter

Twitter has implemented a strike system as part of their Civi Integrity Policy¹⁶ and their COVID19 misleading information policy.¹⁷ Violations of both policies can entail strikes, where two strikes lead to a 12-hour account lock and five or more strikes lead to permanent suspension from the platform. The 12-hour account lock is hard to observe in the data, especially for users who do not have an over the clock tweeting activity. In this section, we provide one example of a temporary suspension¹⁸ of a Twitter account, that seems to be the result of a manual decision concerning a Tweet which violated the rules.

The Twitter account of the website *lifesitenews.com* has been suspended for at least two periods of time: from end of 2019 until fall 2020 for 308 days, then again since January 2021 for having violated Twitter Rules¹⁹. In particular, this website has several failed fact-checks concerning the published articles, according to *Iffy.news*.²⁰ We collected the activity (tweets, replies, quotes, retweets) on their Twitter account *@LifeSite* via the Twitter API, using the historical search endpoint.²¹ We then plotted the number of Tweets, Retweets, Quotes and Replies per day, as shown in panel *a* of figure 2). The two periods of temporary suspension are clearly observed in the data as the user(s) of the account were not allowed to use the functionalities of the Twitter Platform.

To further assess the impact of this double temporary suspension, we collect via Minet Command Line Tool [5], all the tweets that have shared during the same period a url link containing *lifesitenews.com*. Panel (*b*) of figure 2, shows that during both periods of temporary suspension,

¹⁶See help.twitter.com/en/rules-and-policies/election-integrity-policy

¹⁷See help.twitter.com/en/rules-and-policies/medical-misinformation-policy and blog.twitter.com/en_us/topics/company/2021/updates-to-our-work-on-covid-19-vaccine-misinformation

¹⁸See the official documentation on the Twitter’s Help Center regarding account suspension: <https://help.twitter.com/en/managing-your-account/suspended-twitter-accounts>.

¹⁹See *Lifesitenews’s* article discussing the reason for the suspension: <https://www.lifesitenews.com/news/lifesite-is-dumping-twitter-and-so-should-you>. Twitter rules can be found at: <https://help.twitter.com/en/rules-and-policies/twitter-rules>.

²⁰See <https://mediabiasfactcheck.com/life-site-news/>

²¹See the documentation: <https://developer.twitter.com/en/docs/twitter-api/tweets/search/api-reference/get-tweets-search-all>.

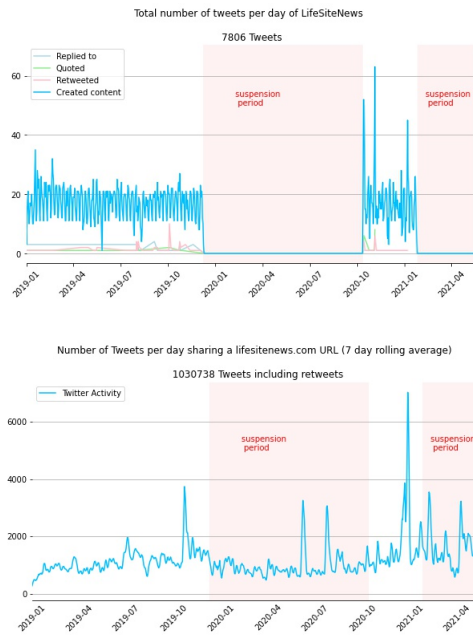


Figure 2: Panel (a): number of Tweets per day of the Twitter account @Lifesite linked to the website lifesitenews.com from January, 2019 until April 2021. Panel (b): number of Tweets per day that have shared a lifesitenews.com URL link from January, 2019 until April 2021.

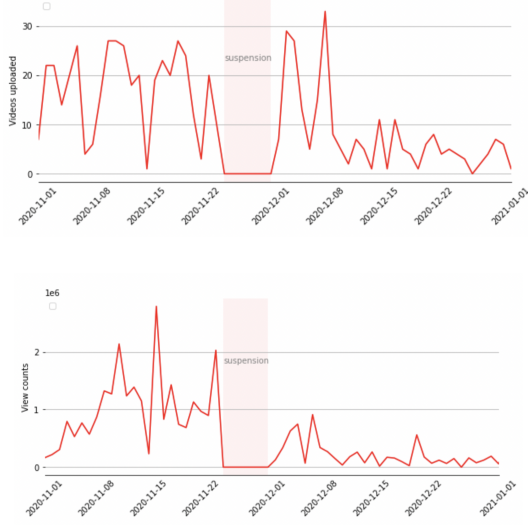


Figure 3: panel (a): Number of Youtube videos uploaded each day by the youtube channel *One America news Network* November 1, 2020 and January 1, 2021. Panel (b): accumulated view counts for videos. The metrics correspond to the videos’ publishing date and the data is retrieved from the youtube API with the *playlists* and *videos* endpoints.

other users still shared *lifesitenews.com* links and that the level was only slightly below the tweeting and retweeting levels prior to the first temporary suspension. More specifically, there was an average of 960 tweets (including retweets) per day over the first temporary suspension period of 308 days from December 9, 2019 until October 12, 2020, against an average of 977 tweets (including retweets) per day during the exact same period one year earlier. Finally, panel (b) points towards the limitations of suspending an account to limit the spread of its content.

2.1.3 Youtube

In this section, we turn to the channel’s temporary or permanent suspension policy of Youtube. Whenever a channel publishes a video that violates the community guidelines for the first time they will usually receive a warning and the content will be removed. For the second time the channel will start receiving strikes. A first strike results in limiting the access of the Youtube channel for one week, like uploading videos, streaming and other activities. Then a second strike is similar but the suspension will be for two weeks. A third strike results in the termination of the channel. The strike count of a channel lasts 90 days. In the special case, where a video is in extreme violation of the guidelines, the publishing channel may get terminated without a warning.²² To illustrate the implementation of this policy we provide two examples for the temporary suspension of the following two Youtube channels: *One America news Network* and *Tony Heller*.

²²See the “Community Guidelines strike basics”. Youtube help, Google Developers, <https://support.google.com/youtube/answer/2802032?hl=en>. Accessed 21 6 2021.



Figure 4: Panel (a): Tweet announcing moving to rumble by OANN (Twitter), Twitter ID 1372238828425998336. Panel (b): Tony Heller’s tweet after getting suspended from Youtube, Twitter ID 1310703852769796097.

First, we investigate the temporary suspension case of the Youtube Channel of *One America News channel*. This channel received a first strike on November 24, 2020 for the promotion of a false cure for COVID19.²³ We collected the activity of the channel OANN (video counts, view counts) using the Youtube API v3, between November 2020 and January 2021. For the video counts, we used the playlist endpoint to retrieve the videos uploaded with their publishing date and for the view count we used the IDs of the videos we had from the playlists and via the videos endpoint we retrieved the view counts on June 2021.²⁴

In addition, as shown in figure 3 when comparing the month before the suspension from 2020/10/24 to 2020/11/24 and one month after from 2020/12/01 to 2021/01/01 it was found that the view count decreased by -73% and the videos uploaded by -55%. Besides that, OANN decided to move officially to Rumble on March 17, 2021 as announced on their Twitter account (see figure 4) and their upload activity on their Youtube channel is close to zero since that announcement.

We now turn to our second example, the temporary suspension of the Youtube channel Tony Heller. This channel got its first strike after posting a video about an anti-covid-lockdown doctor getting arrested (see screenshot in figure 4). The suspension period was for one week from September 29 until October 5. We applied the same methods as in the previous example for the data collection.

²³See nbcnews, “YouTube suspends OANN for violating its Covid-19 policy”. nbcnews, Ahiza García-Hodges, 24 11 2020.

²⁴See the Google documentation <https://developers.google.com/youtube/v3/docs/videos/list> and <https://developers.google.com/youtube/v3/docs/playlists/list>

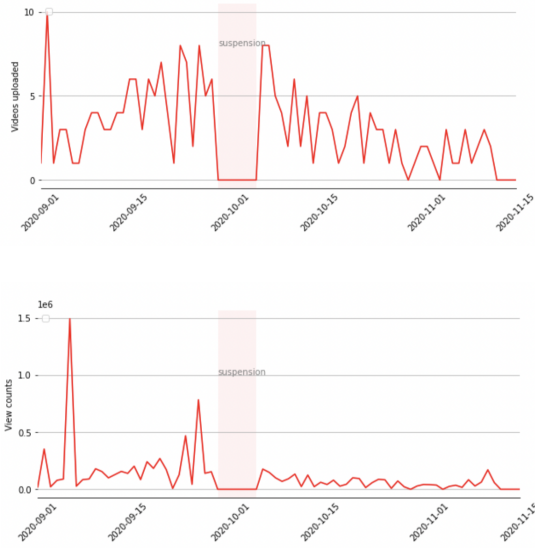


Figure 5: Panel (a): Number of Youtube videos uploaded each day by the Youtube channel *Tony Heller* between September 1, 2020 and November 15, 2020. Panel (b): accumulated view counts for videos uploaded by the same Youtube channel. The date corresponds to the videos’ publishing date.

Figure 4 shows the daily number of videos uploaded by the channel. The suspension period can be observed clearly in the historical data of the channel. observing the reach of the audience Figure 3: Tony heller tweet after getting suspended from youtube (Twitter) using view counts one month before the suspension starting from 2020/08/28 to 2020/09/28 and one month after the suspension from 2020/10/05 to 2020/11/05 the channel witnessed a drop of view counts by -69.5% and the videos published in the channel were less by -29% . This drop in views can show that the suspension period may have a good impact on reducing the audience interest or reach to the channel.

2.2 Blocking links

A third measure that mainstream social media platforms can apply is to prevent users from sharing specific types of content, instead of deleting the content after posting, or users accounts. We will show through examples how platforms prevent users from sharing urls coming from specific domain names.

2.2.1 Facebook

In this section, we provide one case study of how Facebook can block a URL as part of its policy enforcement tools. The Beauty of life (thebl.com/) is a US-based media company that shares pro-Trump views and conspiracy theories such as QAnon.²⁵ Facebook has announced on December

²⁵See [https://en.wikipedia.org/wiki/The_Epoch_Times#Removal_of_The.BL_\(The_Beauty_of_Life\)_from_Facebook](https://en.wikipedia.org/wiki/The_Epoch_Times#Removal_of_The.BL_(The_Beauty_of_Life)_from_Facebook).

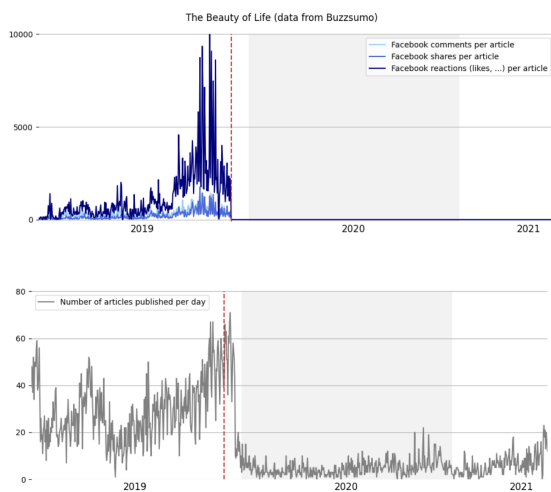


Figure 6: Articles from The Beauty of Life website (thebl.com) published between January 1, 2019 and June 15, 2021 and gathered from the Buzzsumo API. (Top) Facebook engagement metrics (average number of reactions, shares and comments per article). (Bottom) Number of articles published per day. The red line marks the date of December 1, 2019.

20, 2019 that “The BL is now banned from Facebook” for coordinated inauthentic behavior²⁶, which includes using fake accounts that misrepresent one’s identity or using methods to artificially boost the popularity of content. Coordinated inauthentic behavior is a distinct phenomenon from disinformation according to Facebook, as “most of the content shared by coordinated manipulation campaigns isn’t probably false”.²⁷

Nevertheless, for the case of the Beauty of Life, both misinformation and coordinated inauthentic behavior are interlinked, according to the fact-checking organization Snopes which had reported about The BL’s activity to Facebook and other various public articles.²⁸ It took Facebook five months to complete its investigation of The BL before taking action²⁹ and Snopes ended its fact-checking partnership with them in February 2020.³⁰

To verify Facebook’s ban of The BL domain name, we first tested whether we could post a Facebook message containing a url from thebl.com. This turned out to be impossible. But such manual verification cannot inform us whether the ban applies indeed to all Facebook users and accounts (as we used only our own personal accounts), nor when it has started. To further investigate this policy, we collected data from Buzzsumo.³¹ We used the “/search/articles” endpoint

²⁶<https://about.fb.com/news/2019/12/removing-coordinated-inauthentic-behavior-from-georgia-vietnam-and-the-us/>.

²⁷See: <https://about.fb.com/news/2019/10/inauthentic-behavior-policy-update/>.

²⁸<https://www.snopes.com/news/2019/11/12/bl-fake-profiles/>, <https://www.snopes.com/news/2019/11/12/bl-fake-profiles/>, <https://www.snopes.com/news/2019/12/13/facebook-bl-cib/>

²⁹See <https://www.businessinsider.com/facebook-bans-beauty-of-life-for-inauthentic-behavior-2019-12>.

³⁰See <https://www.snopes.com/2019/02/01/snopes-fb-partnership-ends/>.

³¹BuzzSumo is a commercial content database that tracks the volume of user interactions with internet content on Facebook, Twitter, and other social media platforms.

of the Buzzsumo API, to collect the engagement metrics 13 634 articles crawled from the *thebl.com* website between January 1, 2019 and June 15, 2021.³²

The number of Facebook reactions, shares and comments dropped to zero for TheBL’s articles published after December 1, 2019 (see figure 6 top panel), indicating the start of the ban. We can note that although the ban was communicated in an article published³³ on December 20, 2019, it seems to have actually started on December 1, 2019. The number of articles published per day was around 50 until December 20, 2019. The this number decreased drastically to reach a plateau around 5 to 10 articles published per day (see figure 6 bottom panel). The communication around the ban may have discouraged The Beauty of Life to proceed with their activity. Using Buzzsumo data, a global ban over sharing a thebl.com URL link on Facebook was ascertained. It started on December 1, 2019, and was still enforced in June 2021.

2.3 Reducing the visibility

Mainstream Social Media platforms can reduce the visibility of the content created or shared by specific users, whenever they violate the platforms’ rules. The implementation of this policy varies across platforms and is not easy to verify ex-post. In what follows we provide means to verify this policy on Twitter and Facebook.

2.3.1 Facebook

One of Facebook’s measures to regulate misinformation is to reduce the spread of misleading content through their ranking system. Facebook ranks each post and/or ad by assigning to it a relevancy score, where a high score leads to a high likelihood of the post and/or the ad to appear on a user’s newsfeed. Doing so, Facebook can make a post or a whole account less visible by decreasing the relevancy score of its content; this is precisely the *reduce* measure.³⁴ This measure can be verified by looking at the number of views (reach) of a post, but this metric is not available via the CrowdTangle API or on Buzzsumo. Hence we can indirectly investigate the *reduce* measure by looking at the engagement metrics (likes, comments, shares) related to a given post; which are available via the CrowdTangle API and Buzzsumo. If a post reaches less users because it has a lower ranking, then it is less likely to receive likes, comments and shares, relative to a post with a higher ranking.

To illustrate, we investigate the case of the website *Infowars*. This website appears in the Misinformation Directory of FactCheck.org, among other websites who have posted deceptive content³⁵. Furthermore, the factual reporting of *Infowars* has been rated *very low* by the Media Bias / Fact Check resource of Iffy.news.³⁶

On May 2, 2019, Facebook announced they would prohibit users from sharing Infowars content unless, they are explicitly condemning the material.³⁷ To verify the measure, we used the

³²The command can be found in the following Github repository: https://github.com/medialab/truth-and-trust-online-2021/blob/master/code/collect_facebook_buzzsumo_thebl_data.py.

³³See <https://about.fb.com/news/2019/12/removing-coordinated-inauthentic-behavior-from-georgia-vietnam-and-the-us>.

³⁴Lyons, T. (2018, May 22). The three-part recipe for cleaning up your news feed. Facebook Newsroom. about.fb.com/news/2018/05/inside-feed-reduce-remove-inform/.

³⁵See <https://www.factcheck.org/2017/07/websites-post-fake-satirical-stories/>.

³⁶See the Iffy.news page: <https://mediabiasfactcheck.com/infowars-alex-jones/>.

³⁷See <https://www.wired.com/story/facebook-bans-alex-jones-extremists/> and about.fb.com/news/2018/08/enforcing-our-community-standards/.

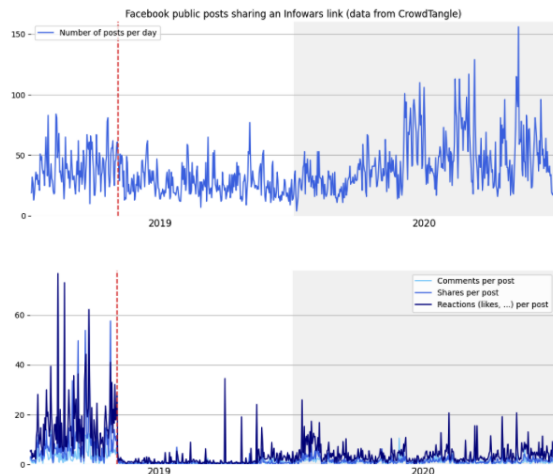


Figure 7: Public Facebook posts sharing an Infowars link in 2019 and 2020 collected from the CrowdTangle API. The red line marks the date of May 2, 2019, when Facebook has announced the ban regarding Infowars. (Top panel) Number of daily posts. (Bottom panel) Engagement metrics: average number of reactions, shares and comments per post.

“/posts/search” endpoint³⁸ of the CrowdTangle API, to collect 37 242 Facebook public posts that had shared a URL link containing “infowars.com”, published between January 1, 2019 and December 31, 2020.³⁹ The command used can be found in the following Github repository: [link](#). The number of public posts sharing an *Infowars* link remained globally stable throughout 2019 (see figure 7 top panel). Thus the measure announced by Facebook doesn’t seem to have prevented users from sharing an *Infowars* link. Nevertheless, a clear drop in engagement was observed on May 2, 2019 (see figure 7 bottom panel). The number of reactions, shares and comments per post have decreased respectively by -94% , -96% and -93% when comparing the two months before and after May 2, 2019. This suggests the measure taken by Facebook in May 2019, is not a ban per se, but rather a *reduce* measure. This is because users could still post *Infowars* links, but these posts generated less engagement. It should be noted that the engagement metrics increased again by the end of 2019 / beginning of 2020, suggesting that the *reduce* measure may have been lifted a few months after its implementation.

As CrowdTangle is tracking posts only from certain public groups and pages, we also used the “/search/articles” endpoint of the Buzzsumo API, to gather a richer Facebook dataset. We collected the engagement data for the 14 232 articles crawled by Buzzsumo from the *Infowars* website between January 1, 2019 and December 31, 2020.⁴⁰ We observe that the articles published after May 2, 2019 received less Facebook engagement than the ones published before (see figure 8),

³⁸see the documentation for more details: github.com/CrowdTangle/API/wiki/Search.

³⁹We found in the collected data some Facebook posts that did not directly share an Infowars link (but rather a YouTube or Facebook video containing an Infowars link in its description), thus we excluded such posts from our data to keep only the 27 721 posts directly sharing an Infowars link.

⁴⁰The command can be found here: https://github.com/medialab/truth-and-trust-online-2021/blob/master/code/collect_facebook_buzzsumo_infowars_data.py.

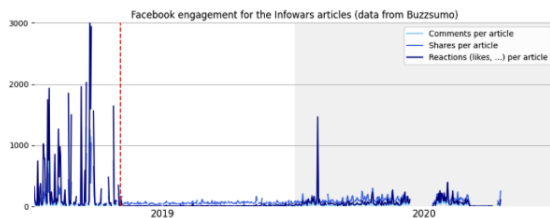


Figure 8: Facebook engagement metrics (average number of reactions, shares and comments per article) for the Infowars articles published in 2019 – 2020 and gathered from the Buzzsumo API. The red line marks the date of May 2, 2019, when Facebook has announced the ban regarding Infowars.

with a percentage change of -97% for the reactions, -59% for the shares and -97% for the comments. An increase in engagement was also observed in 2020. It reinforces the hypothesis that Facebook reduced the reach of posts sharing Infowars links only during a few months in 2019.

Finally, we found irregularities in the number of Infowars articles collected from Buzzsumo. While Infowars usually publishes 20 to 30 articles per day, only 53 articles were collected in the 31-day period between June 11 to July 11, 2020. A temporary crawling problem coming from Buzzsumo may have caused this lack of data. Because no database is perfect, we would like to highlight the importance of cross-checking information between different sources when possible.

2.3.2 Twitter

Twitter can take action against a tweet which violates the Twitter rules⁴¹, by limiting its visibility on users’ timelines and in search results. To illustrate we provide an example for the website *globalresearch.ca*, which has several failed fact-checks according to *iffy.news* - a website which provides a database of websites with low factual reporting levels.⁴²

The website *globalresearch.ca* is linked to the Twitter account @CRG_CRM; which was recently suspended.⁴³ When a user searches via the twitter search-box for any URL link of this website, no results appear as shown in the screenshot in panel (b) of figure 9, taken on June 14, 2021. To further investigate the possible implementation of a reduced visibility measure, we search via the Twitter API for tweets, excluding retweets, containing the query *globalresearch.ca* from January 1, 2021 until June 10, 2021. As shown in panel (a) in figure 10, we find a strictly positive number of tweets containing the URL link *globalresearch.ca* throughout May 2021 and the first week of June 2021. Hence, the visibility of tweets containing this URL link has been reduced because users can no longer access tweets containing the URL link *globalresearch.ca* via the search box. Nevertheless users are not restrained from posting tweets containing this URL, as shown in the

⁴¹See the paragraph *Limiting Tweet visibility*: <https://help.twitter.com/en/rules-and-policies/enforcement-options>.

⁴²For *globalresearch.ca* see <https://mediabiasfactcheck.com/global-research/>.

⁴³We noticed the message about the account suspension on May 25, 2021. But to the best of our knowledge, no official communication by Twitter has announced the suspension nor the exact date at which it was implemented. Hence the account may have gotten suspended anytime between April 15, 2021 and May 25, 2021 (see the suspension screenshot in panel (a) of figure 9). Furthermore, *globalresearch.ca* has multiple failed fact-checks, see <https://mediabiasfactcheck.com/global-research/>

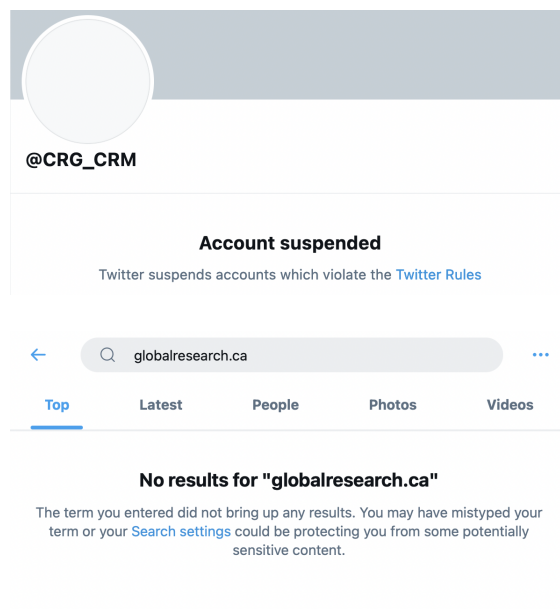


Figure 9: Screenshots taken on June 14, 2021. Top panel: screenshot that shows the account of *globalresearch.ca* suspended on Twitter. Bottom panel: screenshot that shows that no results can be found when searching for *globalresearch.ca*.

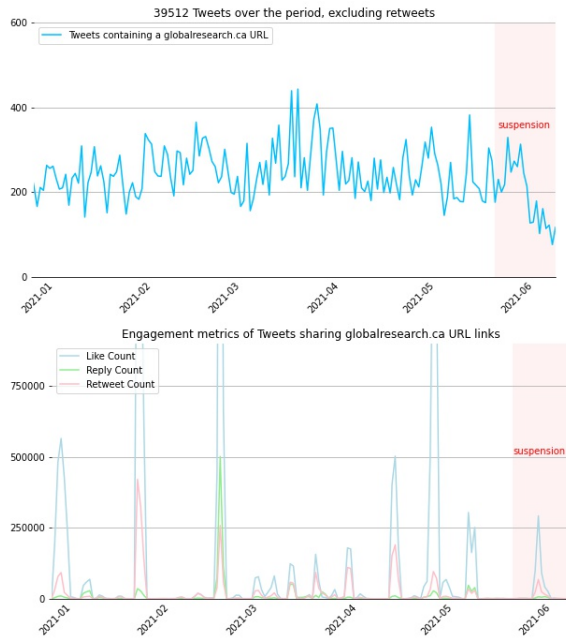


Figure 10: Top panel: daily number of Tweets, excluding retweets, containing the query *globalresearch.ca* from January 1, 2021 until June 10, 2021. Bottom panel: engagement metrics of Tweets containing the query *globalresearch.ca* from January 1, 2021 until June 10, 2021. Data collected via the Twitter API v2 on June 16, 2021.

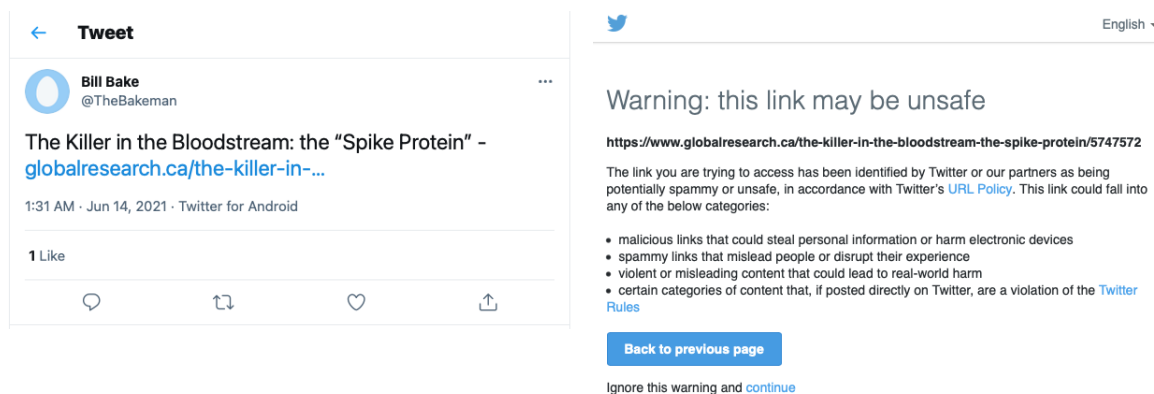


Figure 11

screenshot in panel (a) of figure 11, found by taking the tweet ID of one of the collected tweets via the Twitter API. Furthermore, those collected Tweets have strictly positive engagement metrics as shown in panel (b) of figure 10. Hence, the users who tweet articles from the *globalresearch.ca* website receive tweet level engagement from their own followers. Finally, when a user attempts to click on the URL link *globalresearch.ca* contained in the Tweet, a warning message appears and indicates that the link may be unsafe (see screenshot in panel (b) in figure 11).

2.3.3 Youtube

What about something about recommendations ? authoritative content

2.4 Flags, Notices and labels

2.4.1 Facebook

2.4.2 Twitter

Alongside other social networking platforms, when the content of a tweet violates the Twitter rules, a notice can be added to provide more context according to Twitter’s Help Center.⁴⁴ At the tweet level, notices take the form of a label or an interstitial. Labels are context specific (e.g. COVID19 or presidential elections) and redirect users to a URL link to get more context. Interstitials are presented as a greyed box on top of a tweet, which indicate sensitive content, violations of Twitter rules, withheld tweets for violation of local laws or even tweets from suspended accounts. At the account level, notices can indicate whether an account has been temporarily or permanently suspended.

In this section, we take a deeper look at how labels and notices are introduced by Twitter, to indicate content which is inaccurate or false. To that end, we gathered a set of 3094 URL links of articles which were marked as *False* by Science Feedback, a fact-checking organization verifying the credibility of science-related viral information. As a second step, we collected (on June 30,

⁴⁴See Notices on Twitter and what they mean: <https://help.twitter.com/en/rules-and-policies/notices-on-twitter>.



Figure 12: Two tweets sharing the same URL link marked as False by a Fact-Checker, screenshots taken on June 20, 2021. Panel (a): Tweet ID 1241088065462026242 without a label. Panel (b): Tweet ID 1261876171584745472 containing a label.

2021) via Minet Command line tool [5] all the tweets that have shared a URL link which belongs to the set of 3094 links marked as *False*. This data collection resulted in 323 938 tweets, excluding retweets. Only 28 tweets contained the label “Get the facts about COVID-19”, 5 tweets contained the label “Learn about US 2020 election security efforts” and only 1 had the following label “This claim about election fraud is disputed”. Furthermore, we noticed that the labeling rule might not be applied uniformly on a given set of tweets sharing the exact same URL link, among the set of collected tweets. More specifically, exactly 657 tweets had shared a URL link redirecting to a video on Bitchute, entitled “Important information on coronavirus 5G Kung Flu”. Among those 657 tweets only 3 contained the label “Get the facts about COVID-19” (see figure 12). This points towards the non-automation of the tweet labelling process and that it might be that only 3 tweets got labelled after being reported by a user.

We further examine the placement of interstitials that indicate a possibly sensitive content. We find that only 2.97% out of 323 938 tweets containing a URL marked as False, have an interstitial “potentially sensitive content”. In particular, many speak about COVID19 and do not contain a label to provide users with more context from authoritative sources. Figure 13 provides an example of a tweet sharing a URL marked as False by a Fact-checker and which contains an interstitial “potentially sensitive content”. We find 32 other tweets, among our set of collected tweets, who share the exact same URL link in our the previous example and among those 32 only 5 tweets had the interstitial “potentially sensitive content”. Again this points towards the non-automation of the interstitials placement.

Finally, to the best of our knowledge, when using the Twitter API v2, there is no field which indicates whether a tweet is labeled or not; while the interstitial “possibly sensitive” and “withheld” are both Tweet fields that can be recovered⁴⁵ from the API.

2.4.3 Youtube

Youtube may provide an information panel for videos with topics that are prone to misinformation like COVID19, moon landing, and climate change.⁴⁶ The information panel contains main information about the topic from independent third party partners. Youtube states that these panels exist regardless of the video point of view. In addition, these panels are not available in all languages and countries yet. (“Information panel giving topical context”)

Based on a video list (171 videos) checked by fact checkers and marked to have misinformation. We collected the information panels in each video by scrapping the content of the web and connecting to a US server. Four types of information panels were found in the list of visited videos. The first information panel gives info about COVID19 as in Fig 14, the second is about Vaccine as in Fig 14, the third is about Climate change as in Fig 14 and last information panel showed on a COVID19 video about the US voting as in Fig 14.

We classified the list of videos examined for the information panel based on its content as in Table 2. The Table shows the number of videos with a panel in each category. For COVID and vaccine related videos more than half of the videos had a panel, nevertheless, we noticed that some videos get uploaded multiple times by different channels under different video titles but others didn’t even though the video was the same. In addition, we noticed for COVID when the video title doesn’t include keywords like (Testing, Pandemic, COVID, coronavirus), the video wouldn’t have

⁴⁵See <https://developer.twitter.com/en/docs/twitter-api/data-dictionary/object-model/tweet>.

⁴⁶See the section “Information panel giving topical context” on Youtube Help, Google accessed on June 28, 2021: support.google.com/youtube/answer/9004474?hl=en.



Figure 13: Tweet ID 1285866521533861888. Screenshots taken on July 1, 2021.

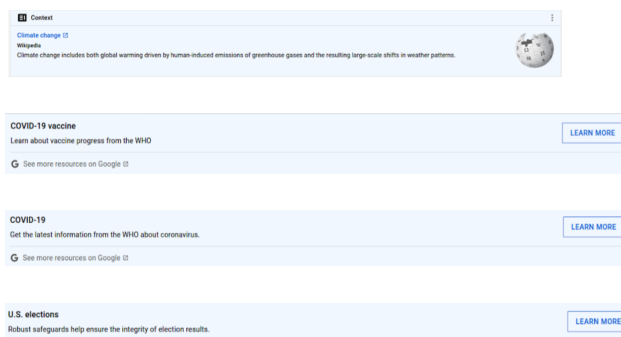


Figure 14: Top to bottom respectively: information panel displayed under some videos related to climate change, COVID19 vaccine, COVID19 and elections.

Category	With information panel	Without Information panel	Total	IP percentage
General Health	0	11	11	0
Vaccine	21	7	28	75
COVID	73	43	116	63
Climate change	5	13	18	28













Table 2

a panel associated with it, and in some cases when the video title have variations of word COVID like (C.O.V.I.D or Cv19) it wouldn't include a panel either. Therefore, we suspect that youtube is automatically adding panels to the videos based on the video title not the content of the video since many of the videos that didn't have a panel had the word COVID or coronavirus or vaccine mentioned in the video itself. For climate change we only examined 18 misleading videos and we found 72% of them didn't have a panel. Lastly , for the general health category, we included the videos that weren't related to COVID but had misinformation in topics like fake cures for cancer, abortion, and viruses; there was no information panel in all these videos. Therefore, we suspect that youtube just add information panels to the most controversial topics that can have misleading information like the climate change, COVID19, flat earth and vaccine.

3 Discussion

- For a previous research project⁴⁷, we searched on CrowdTangle for public accounts sharing specific content associated with misinformation in November 2020, and selected 94 Facebook pages corresponding to our criteria. We then tried to collect these pages' posts in January 2021, and discovered that 11 pages could not be found anymore. This highlights an important issue when studying misinformation trends on Facebook: some data disappears from the CrowdTangle API as accounts are deleted or changed to *private*.
- To facilitate the verification in the policy applications, we would generally recommend for the platforms to be more transparent. But too much transparency on how the regulation policies are exactly implemented can actually backfire. For example YouTube is certainly applying an 'information' banner on all videos mentioning Covid and related terms in their title. Misinformation accounts are trying to avoid the official banners by using terms as 'C.O.V.I.D' or 'C O V I D'. If YouTube was totally transparent on that matter and published the list of 'dangerous' words that leads to an information banner, this list would of course help us to understand YouTube's policies but it would also help the misinformation actors to escape the regulation. There is thus a balance between communicating enough so the public can know precisely how the platforms are regulating their content, but without giving too much information that would allow the policies to be bypassed.
- There are other ways to collect data from platforms, and besides Buzzsumo, other API are also aggregating data from multiple social platforms. For example Newsguard, blablabla... In

⁴⁷reference?

Rules		https://www.facebook.com/communitystandards/recentupdates/
		https://help.twitter.com/en/rules-and-policies/twitter-rules
		youtube.com/intl/en_us/howyoutubeworks/policies/community-guidelines/
Rules enforcement		https://transparency.fb.com/data/community-standards-enforcement/
		https://transparency.twitter.com/en/reports/rules-enforcement.html
		https://transparencyreport.google.com/youtube-policy/
Transparency center		https://transparency.fb.com/data/
		https://transparency.twitter.com/en/reports.html
		https://transparencyreport.google.com/?hl=en
Policy regarding Covid-19		https://www.facebook.com/help/230764881494641/
		https://help.twitter.com/en/rules-and-policies/medical-misinformation-policy
		https://support.google.com/youtube/answer/9891785

References

- [1] David A. Broniatowski, Daniel Kerchner, Fouzia Farooq, Xiaolei Huang, Amelia M. Jamison, Mark Dredze, and Sandra Crouse Quinn. Debunking the misinfodemic: Coronavirus social media contains more, not less, credible content. *mimeo*, 2020.
- [2] Nir Grinberg, Kenneth Joseph, Lisa Friedland, Briony Swire-Thompson, and David Lazer. Fake news on twitter during the 2016 u.s. presidential election. *Science*, 263(274), 2019.
- [3] Adam Hughes and Stefan Wojcik. 10 facts about americans and twitter. *Pew Research Center*, 2019.
- [4] David Lazer, Matthew Baum, Yochai Benkler, Adam Berinsky, Kelly Greenhill, Filippo Menczer, Miriam Metzger, Brendan Nyhan, Gordon Pennycook, David Rothschild, Michael Schudson, Steven Sloman, Cass Sunstein, Emily Thorson, Duncan Watts, and Jonathan Zittrain. The science of fake news. *Science*, 359(6380), 2018.
- [5] Guillaume Plique, Pauline Breteau, Jules Farjas, Héloïse Théro, and Jean Descamps. Minet, a webmining cli tool and library for python zenodo. <http://doi.org/10.5281/zenodo.4564399>. 2019.

4 Appendix