# Exposure to Social Engagement Metrics
# Increases Vulnerability to Misinformation

Mihai Avram[1,3], Nicholas Micallef[2], Sameer Patil[3], Filippo Menczer[1,3]
[1]Observatory on Social Media, Indiana University, Bloomington, IN, USA
[2]Center for Cybersecurity, New York University, Abu Dhabi, United Arab Emirates
[3]Luddy School of Informatics, Computing, and Engineering, Indiana University Bloomington, IN, USA

Topics: Misinformation; social media; social engagement; news literacy; game

## Article's lead

News feeds in virtually all social media platforms include engagement metrics, such as the number of times each post is liked and shared. We find that exposure to these social engagement signals increases the vulnerability of users to low-credibility information. This finding has important implications for the design of social media interactions in the misinformation age. To reduce the spread of misinformation, we call for technology platforms to rethink the display of social engagement metrics. Further research is needed to investigate whether and how engagement metrics can be presented without amplifying the spread of low-credibility information.

## Research questions

- What is the effect of exposure to social engagement metrics on people's propensity to share content?

- Does exposure to high engagement metrics increase the chances that people will like and share questionable content and/or make it less likely that people will engage in fact checking of low-credibility sources?

## Essay summary

- We investigated the effect of social engagement metrics on the spread of information from low-credibility sources using Fakey [1], a news literacy game that simulates a social media feed (Figure 1). The game presents users with actual current news articles from mainstream and low-credibility media sources. A randomly generated social engagement metric is displayed with each presented article. Users are instructed to share, like, fact check, or skip articles.

- From a 19-month deployment of the game, we analyzed game sessions by over 8,500 unique users, mostly from the US, involving approximately 120,000 articles, half from low-credibility sources.

- Our findings show that displayed engagement metrics can strongly influence interaction with low-credibility information. The higher the shown engagement, the more prone people were to share questionable content and less to fact check it.

- These findings imply that social media platforms must rethink whether and how engagement metrics should be displayed such that they do not facilitate the spread of misinformation or hinder the spread of legitimate information. Further research is needed to guard against malicious tampering with engagement metrics at an early stage and to design educational interventions that teach users to prioritize trustworthiness of news sources over social engagement metrics.
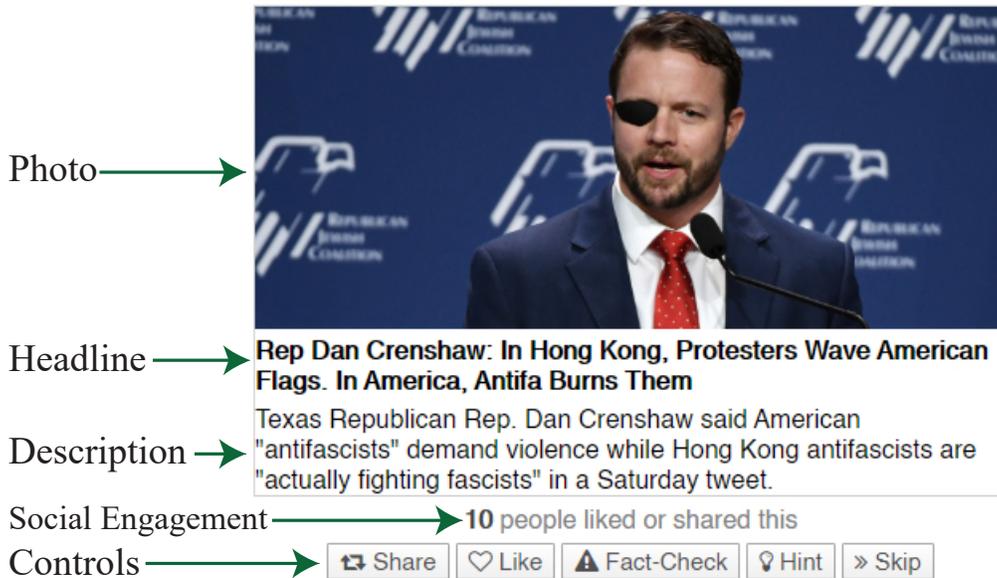
---

[1]https://fakey.iuni.iu.edu/

Figure 1: A news post in the social media feed simulated by the game.

# Implications

Online misinformation is a critical societal threat in the current digital age, and social media platforms are a major vehicle used to spread it (Guess, Nagler, & Tucker, 2019; Lazer et al., 2018; Hameleers, Powell, Meer, & Bos, 2020). As an illustration, the International Fact-Checking Network found more than 3,500 false claims related to the coronavirus in less than 3.5 months.[2] Viral misinformation can cause serious societal harm in multiple ways: affecting public health (Sharma et al., 2020), influencing public policy (Lazer et al., 2018), instigating violence (Arif, Stewart, & Starbird, 2018; Starbird, Maddock, Orand, Achterman, & Mason, 2014), spreading conspiracies (Samory & Mitra, 2018), reducing overall trust in authorities (Gupta, Kumaraguru, Castillo, & Meier, 2014; Shin & Thorson, 2017; Vosoughi, Roy, & Aral, 2018), and increasing polarization and conflict (Stewart, Arif, & Starbird, 2018).

The growing societal impact of misinformation has driven research on technical solutions to detect and in some cases — depending on platform policies — stop actors that generate and spread such content. The techniques have leveraged network analytics (Ratkiewicz et al., 2011; Jin, Dougherty, Saraf, Cao, & Ramakrishnan, 2013), supervised models of automated behavior (Ferrara, Varol, Davis, Menczer, & Flammini, 2016; Varol, Ferrara, Davis, Menczer, & Flammini, 2017; Yang et al., 2019; Yang, Varol, Hui, & Menczer, 2020; Hui et al., 2019), time series analysis to detect promoted campaigns (Varol, Ferrara, Menczer, & Flammini, 2017), and natural language processing for flagging factually incorrect content (Prez-Rosas, Kleinberg, Lefevre, & Mihalcea, 2017; Kumar, West, & Leskovec, 2016). On the user interface side, researchers have explored the use of credibility indicators to flag misinformation and alert users (Clayton et al., 2019). Such credibility indicators can lead to a reduction in sharing the flagged content (Yaqub, Kakhidze, Brockman, Memon, & Patil, 2020; Pennycook, Bear, Collins, & Rand, 2019; Pennycook, McPhetres, Zhang, & Rand, 2020; Nyhan, Porter, Reifler, & Wood, 2019).

Studies have explored the role of environmental, emotional, and individual factors in online contagion (Kramer, Guillory, & Hancock, 2014; Ferrara & Yang, 2015; Coviello et al., 2014; Grinberg, Joseph, Friedland, Swire-Thompson, & Lazer, 2019; Yaqub et al., 2020). However, there has been little empirical research on the effects of current elements of social media feeds on the spread of misinformation (Hameleers et al., 2020; Shen et al., 2019). To address this gap, we empirically investigated how the spread of low-credibility content is affected by exposure to typical social engagement metrics, i.e., the numbers of likes and shares shown for a news article. We found near-perfect correlations between displayed social engagement metrics and user actions related to information from low-credibility sources. We interpret these results as showing that social engagement metrics are doubly problematic for the spread of low-credibility content: high levels of engagement make it less likely that people will scrutinize potential misinformation, at the same time making it more likely that they will like or share it. For example, we recently witnessed a campaign to make the "Plandemic" disinformation video go viral. Our results tell us that people are more likely to endorse

---

[2] https://poynter.org/coronavirusfactsalliance

the video without bothering to verify its content, simply because they see that many other people shared it. In other words, exposure to engagement metrics in social media amplifies our vulnerability to questionable content.

To interpret these findings, consider that the probability of sharing a piece of information grows with the number of times one is exposed to it, a phenomenon called *complex contagion* (Romero, Meeder, & Kleinberg, 2011; Mønsted, Sapieżyński, Ferrara, & Lehmann, 2017). Engagement metrics are proxies for multiple exposures, therefore they are intended to provide signals about the importance, relevance, and reliability of information — all of which contribute to people's decisions to consume and share the information. In other words, being presented with high engagement metrics for an article mimics being exposed to the article multiple times: the brain is likely to assess that the article must be worthy of attention because many independent sources have validated the news article by liking or sharing it.

A key weakness in the cognitive processing of engagement metrics is the assumption of independence; an entity can trick people by maliciously boosting engagement metrics to create the *perception* that many users interacted with an article. In fact, most disinformation campaigns rely on inauthentic social media accounts to tamper with engagement metrics, creating an initial appearance of virality that becomes reality once enough humans are deceived (Shao, Ciampaglia, et al., 2018). To prevent misinformation amplified by fake accounts from going viral, we need sophisticated algorithms capable of early-stage detection of coordinated behaviors that tamper with social engagement metrics (Hui et al., 2019; Yang et al., 2020; Pacheco et al., 2020).

Our findings hold important implications for the design of social media platforms. Further research is needed to investigate how alternative designs of social engagement metrics could reduce their effect on misinformation sharing (e.g., by hiding or making engagement less visible for certain posts), without negatively impacting the sharing of legitimate and reliable content. A good trade-off between these two conflicting needs will require a systematic investigation of news properties that can help determine differential display of engagement metrics. Such properties may include the type of sources (e.g., whether claims originate from unknown/distrusted accounts or low-credibility sources) and the type of topics (e.g., highly sensitive or polarizing topics with a significant impact on society).

Further research is also needed to design literacy campaigns (such as Fakey [3], a news literacy game that simulates a social media feed) that teach users to prioritize trustworthiness of sources over engagement signals when consuming content on social media. Studies could investigate the possibility of introducing intermediary pauses when consuming news through a social media feed (Fazio, 2020) and limiting automated or high-speed sharing. A comprehensive literacy approach to reduce the vulnerability of social media users to misinformation may require a combination of these interventions with others, such as inoculation theory (Roozenbeek, van der Linden, & Nygren, 2020; Roozenbeek & van der Linden, 2019b, 2019a; Basol, Roozenbeek, & van der Linden, 2020), civic online reasoning (McGrew, 2020), critical thinking (Lutzke, Drummond, Slovic, & rvai, 2019), and evaluation of news feeds (Nygren, Brounéus, & Svensson, 2019).

# Findings

## Finding 1: High levels of social engagement results in lower fact checking and higher liking/sharing, especially for low-credibility content.

For each article shown in the game, the user is presented a photo, headline, description, and a randomly generated social engagement level. Based on this information, the user can share, like, or fact check the article (Figure 1). To earn points in the game, the user must share or like articles from mainstream sources and/or fact check articles from low-credibility sources. The social engagement level shown with the article provides an opportunity to investigate its effect on behaviors that result in the spread of questionable content.

We measured the correlation between the social engagement metric $\eta$ displayed to users and the rates at which the corresponding articles from low-credibility sources were liked/shared or fact checked by the users. Given the realistically skewed distribution of $\eta$ values, we sorted the data into logarithmic bins based on the shown social engagement levels. For each bin $\lfloor \log_{10}(\eta+1) \rfloor$, we calculated the liking/sharing and fact-checking rates across articles and users. We measured correlation using the non-parametric Spearman test as the data is not normally distributed. We found a significant positive correlation between social engagement level and liking/sharing (Spearman $\rho = 0.97$, $p < 0.001$) and a significant negative correlation between social engagement level and fact checking (Spearman
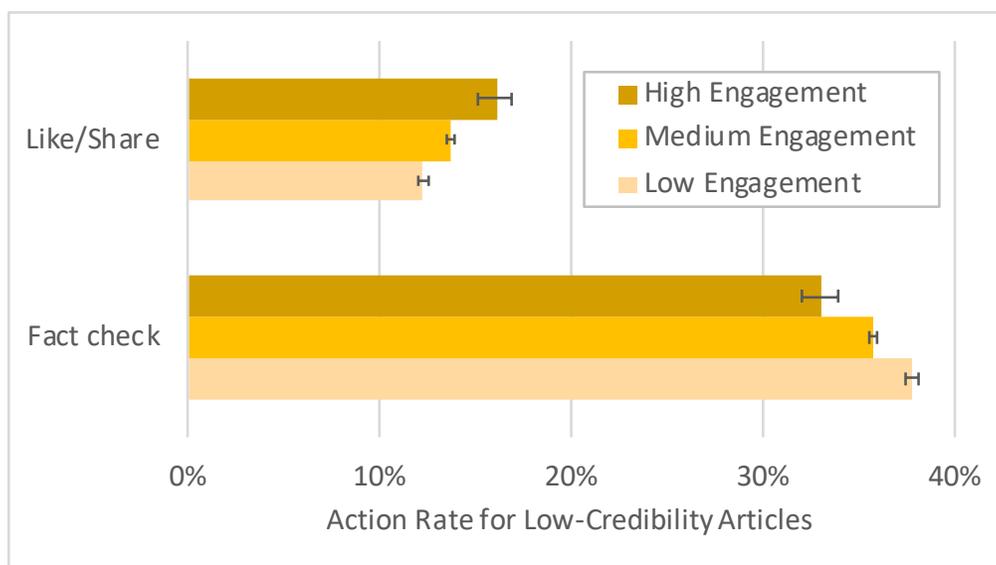
---

[3] https://fakey.iuni.iu.edu/

Figure 2: Mean rates of liking/sharing and fact checking low-credibility articles, categorized by social engagement level (see text). Error bars represent the standard error rate.

$\rho = -0.97$, $p < 0.001$) for articles from low-credibility sources.

We found similar relationships between social engagement levels and user behaviors for mainstream news article as well, however the correlations are less strong: $\rho = 0.66$ for liking/sharing and $\rho = -0.62$ for fact checking.

## Finding 2: People are more vulnerable to low-credibility content with high social engagement.

The previous finding is at the population level, aggregating across users. To delve further into the effect of social engagement exposure on individual users, we analyzed whether different social engagement levels influenced each user's liking/sharing and fact-checking rates for articles from low-credibility sources. For this analysis, we treated each user as an independent entity and categorized engagement into three levels: low ($0 \leq \eta < 10^2$), medium ($10^2 \leq \eta < 10^5$), and high ($10^5 \leq \eta \leq 10^6$). For each user, we counted the number of low-credibility articles to which they were exposed within each social engagement bin. We then calculated the corresponding proportions of these articles that each user liked/shared or fact checked. Figure 2 plots the mean liking/sharing and fact-checking rates for low-credibility articles. Although users were more likely to fact check than like or share low-credibility content, Figure 2 shows that the trends observed at the population level held at the individual level as well.

Since the data is not normally distributed ($p < 0.05$ using the Shapiro-Wilk test for normality), we used the Kruskal-Wallis test to compare differences between the three bins of social engagement levels. The test revealed a statistically significant effect of social engagement levels: fact checking ($\chi^2(2) = 214.26$, $p < 0.001$) and liking/sharing ($\chi^2(2) = 417.14$, $p < 0.001$) rates for low-credibility articles differed across the bins. To determine which levels of social engagement impacted the rates at which low-credibility articles were liked/shared or fact checked, we conducted post-hoc Mann-Whitney tests with Bonferroni correction for multiple testing across all pairs of social engagement bins and found that liking/sharing as well as fact-checking rates were statistically significantly different across all pairings ($p < 0.001$).

We employed the same approach to examine liking/sharing and fact-checking rates for mainstream articles across the three bins of social engagement levels. Similar to low-credibility articles, the Kruskal-Wallis test revealed a statistically significant effect of social engagement level on liking/sharing ($\chi^2(2) = 161.80$, $p < 0.001$) and fact checking ($\chi^2(2) = 576.37$, $p < 0.001$) rates for mainstream articles.

# Methods

## Social media simulation

To conduct our experiment investigating the effect of exposure to social engagement metrics on susceptibility to questionable content, we developed and deployed Fakey [4], an online news literacy game that simulates fact checking on a social media feed. The user interface of the game mimics the appearance of Facebook or Twitter feeds for players who log into the game through those platforms. The game provides users with batches of ten news articles in the form of a news feed, as shown in Figure 1. Each article consists of elements that are typically displayed by popular social media platforms: photo, headline, description, and social engagement metrics.

For each article from mainstream as well as low-credibility sources, the game displays a single social engagement metric about the combined number of shares and likes. Not having separate metrics for shares and likes decreases the cognitive workload for game players and simplifies analysis. Engagement values are randomly drawn from an approximately log-normal distribution with a maximum possible value (cutoff) of $\eta = 10^6$. The distribution is such that roughly 69% of the articles would display engagement values $\eta > 10^2$ and roughly 3% would display values $\eta > 10^5$. Although the simulated engagement in the game is not drawn from empirical data, the metric numbers shown have a heavy tail similar to those typically observed in social media (Vosoughi et al., 2018).

Below each article is a set of action buttons to share, like, fact check, or skip the article or use a hint. Before playing the game, users are instructed that clicking *Share* is equivalent to endorsing an article and sharing it with the world, clicking *Like* is equivalent to endorsing the article, and clicking *Fact Check* signals that the article is not trusted. After playing one round of ten articles, users have the option to play another round or check a leader-board to compare their skill with other players.

## Content selection

We follow the practice of analyzing content credibility at the domain (website) level rather than the article level (Lazer et al., 2018; Shao, Ciampaglia, et al., 2018; Shao, Hui, et al., 2018; Grinberg et al., 2019; Pennycook & Rand, 2019; Bovet & Makse, 2019). Each article in the game is selected from one of two types of news sources: mainstream and low-credibility.

For mainstream news, we manually selected 32 sources with a balance of moderate liberal, centrist, and moderate conservative views: *ABC News Australia, Al Jazeera English, Ars Technica, Associated Press, BBC News, Bloomberg, Business Insider, Buzzfeed, CNBC, CNN, Engadget, Financial Times, Fortune, Independent, Mashable, National Geographic, New Scientist, Newsweek, New York Magazine, Recode, Reuters, Techcrunch, The Economist, The Guardian, The New York Times, Next Web, Telegraph, Verge, The Wall Street Journal, The Washington Post, Time, USA Today.* Current articles are provided by the News API.[5]

The set of low-credibility sources was selected based on flagging by various reputable news and fact-checking organizations (Shao, Ciampaglia, et al., 2018; Shao, Hui, et al., 2018). These sources tend to publish fake news, conspiracy theories, clickbait, rumors, junk science, and other types of misinformation. The articles are provided by the Hoaxy API.[6]

For each round, the game randomly selects five articles each from mainstream and low-credibility sources. For any given source, any article returned by the News or Hoaxy API is shown to the user regardless of topic, without further selection or filtering except for ensuring that the same story is not shown to the same player multiple times across rounds.

## Data collection

The game is available online through a standard web interface and as a mobile app via the Google Play Store and the Apple App Store. The mobile app is available in English-speaking countries: United States, Canada, United Kingdom, and Australia. People from other countries can still play the game through the web interface.

---

[4]https://fakey.iuni.iu.edu/

[5]https://newsapi.org

[6]http://rapidapi.com/truthy/api/hoaxy

The present analysis is based on data from a 19-month deployment of the game, between May 2018 and November 2019. During this period, we advertised the game through several channels, including social media (Twitter and Facebook), press releases, conferences, keynote presentations, and word of mouth. We recorded game sessions involving approximately 8,606 unique users[7] and 120,000 news articles, approximately half of which from low-credibility sources. We did not collect demographic information, but we collected anonymous data from Google Analytics embedded within the game's hosting service. Participants originated from the United States (78%), Australia (8%), UK (4%), Canada (3%), Germany (3%), and Bulgaria (2%).

## Limitations

Our news literacy app emulates relevant interface elements of popular social media platforms such as Facebook and Twitter, without the ethical concerns of real-world content manipulation (Kramer et al., 2014). Yet, conducting the study in a simulated game environment rather than an actual platform presents clear limitations — the experience and context are not identical. For example, to limit the cognitive burden on players, we capture only Like and Share actions; these were the earliest ones deployed on social media platforms and as such are most common across platform and most familiar to users.

The even mix of articles from mainstream and low-credibility sources is not necessarily representative of the proportion of misinformation to which social media users are exposed in the wild. The fact-checking game also primes users to *expect* misinformation, potentially making it more likely to be spotted. These factors might make users more suspicious within the game compared to the real world, increasing fact-checking rates. However, there is no reason to believe that these effects impact results about engagement signals.

While this study is focused on user interaction elements, other factors related to users and content can affect the spread of questionable content. To respect player privacy, we chose not to collect any user information apart from game analytics. However, knowledge about the background of the users (e.g., education, demographics, political affiliation) might provide further insight into vulnerability to misinformation. Similar refinements in insight would be provided by examining which types of content are more likely to be influenced by engagement metrics. These are important avenues for future research.

## Acknowledgments

## Bibliography

Arif, A., Stewart, L. G., & Starbird, K. (2018, November). Acting the part: Examining information operations within #blacklivesmatter discourse. *Proc. ACM Hum.-Comput. Interact.*, *2*(CSCW). Retrieved from `https://doi.org/10.1145/3274289` doi: 10.1145/3274289

Basol, M., Roozenbeek, J., & van der Linden, S. (2020). Good news about bad news: Gamified inoculation boosts confidence and cognitive immunity against fake news. *Journal of Cognition*, *3*(1).

Bovet, A., & Makse, H. A. (2019). Influence of fake news in twitter during the 2016 us presidential election. *Nature communications*, *10*(1), 1-14.

Clayton, K., Blair, S., Busam, J. A., Forstner, S., Glance, J., Green, G., ... Nyhan, B. (2019). Real solutions for fake news? measuring the effectiveness of general warnings and fact-check tags in reducing belief in false stories on social media. *Political Behavior*. Retrieved from `https://doi.org/10.1007/s11109-019-09533-0` doi: 10.1007/s11109-019-09533-0

Coviello, L., Sohn, Y., Kramer, A. D., Marlow, C., Franceschetti, M., Christakis, N. A., & Fowler, J. H. (2014). Detecting emotional contagion in massive social networks. *PloS one*, *9*(3).

Fazio, L. (2020). Pausing to consider why a headline is true or false can help reduce the sharing of false news. *Harvard Kennedy School Misinformation Review*, *1*(2).

---

[7]We used analytics to aggregate anonymous sessions by the same person. However, this approach cannot ascribe anonymous sessions to a single person with complete certainty. Therefore, we cannot provide a precise number of unique users.

Ferrara, E., Varol, O., Davis, C., Menczer, F., & Flammini, A. (2016). The rise of social bots. *Comm. ACM*, *59*(7), 96–104. Retrieved from `http://dl.acm.org/authorize?N16298` doi: 10.1145/2818717

Ferrara, E., & Yang, Z. (2015, 11). Measuring emotional contagion in social media. *PLOS ONE*, *10*(11), 1-14. Retrieved from `https://doi.org/10.1371/journal.pone.0142390` doi: 10.1371/journal.pone.0142390

Grinberg, N., Joseph, K., Friedland, L., Swire-Thompson, B., & Lazer, D. (2019). Fake news on twitter during the 2016 us presidential election. *Science*, *363*(6425), 374-378.

Guess, A., Nagler, J., & Tucker, J. (2019). Less than you think: Prevalence and predictors of fake news dissemination on facebook. *Science Advances*, *5*(1). Retrieved from `https://advances.sciencemag.org/content/5/1/eaau4586` doi: 10.1126/sciadv.aau4586

Gupta, A., Kumaraguru, P., Castillo, C., & Meier, P. (2014). Tweetcred: Real-time credibility assessment of content on Twitter. In L. M. Aiello & D. McFarland (Eds.), *Social informatics: 6th international conference, socinfo 2014, barcelona, spain, november 11-13, 2014. proceedings* (p. 228-243). Cham: Springer International Publishing. Retrieved from `https://doi.org/10.1007/978-3-319-13734-6_16` doi: 10.1007/978-3-319-13734-6_16

Hameleers, M., Powell, T. E., Meer, T. G. V. D., & Bos, L. (2020). A picture paints a thousand lies? the effects and mechanisms of multimodal disinformation and rebuttals disseminated via social media. *Political Communication*, *37*(2), 281-301. Retrieved from `https://doi.org/10.1080/10584609.2019.1674979` doi: 10.1080/10584609.2019.1674979

Hui, P.-M., Yang, K.-C., Torres-Lugo, C., Monroe, Z., McCarty, M., Serrette, B., . . . Menczer, F. (2019). Botslayer: real-time detection of bot amplification on twitter. *Journal of Open Source Software*, *4*(42), 1706.

Jin, F., Dougherty, E., Saraf, P., Cao, Y., & Ramakrishnan, N. (2013). Epidemiological modeling of news and rumors on twitter. In *Proceedings of the 7th workshop on social network mining and analysis.* New York, NY, USA: Association for Computing Machinery. Retrieved from `https://doi.org/10.1145/2501025.2501027` doi: 10.1145/2501025.2501027

Kramer, A. D. I., Guillory, J. E., & Hancock, J. T. (2014). Experimental evidence of massive-scale emotional contagion through social networks. *Proceedings of the National Academy of Sciences*, *111*(24), 8788-8790. Retrieved from `https://www.pnas.org/content/111/24/8788` doi: 10.1073/pnas.1320040111

Kumar, S., West, R., & Leskovec, J. (2016). Disinformation on the web: Impact, characteristics, and detection of wikipedia hoaxes. In *Proceedings of the 25th international conference on world wide web* (p. 591602). Republic and Canton of Geneva, CHE: International World Wide Web Conferences Steering Committee. Retrieved from `https://doi.org/10.1145/2872427.2883085` doi: 10.1145/2872427.2883085

Lazer, D. M. J., Baum, M. A., Benkler, Y., Berinsky, A. J., Greenhill, K. M., Menczer, F., . . . Zittrain, J. L. (2018). The science of fake news. *Science*, *359*(6380), 1094–1096. Retrieved from `https://science.sciencemag.org/content/359/6380/1094` doi: 10.1126/science.aao2998

Lutzke, L., Drummond, C., Slovic, P., & rvai, J. (2019). Priming critical thinking: Simple interventions limit the influence of fake news about climate change on facebook. *Global Environmental Change*, *58*, 101964. Retrieved from `http://www.sciencedirect.com/science/article/pii/S0959378019307009` doi: doi.org/10.1016/j.gloenvcha.2019.101964

McGrew, S. (2020). Learning to evaluate: An intervention in civic online reasoning. *Computers & Education*, *145*, 103711. Retrieved from `http://www.sciencedirect.com/science/article/pii/S0360131519302647` doi: doi.org/10.1016/j.compedu.2019.103711

Mønsted, B., Sapieżyński, P., Ferrara, E., & Lehmann, S. (2017). Evidence of complex contagion of information in social media: An experiment using twitter bots. *PloS ONE*, *12*(9). doi: 10.1371/journal.pone.0184148

Nygren, T., Brounéus, F., & Svensson, G. (2019). Diversity and credibility in young people's news feeds: A foundation for teaching and learning citizenship ina digital era. *Journal of Social Science Education*, *18*(2), 87-109.

Nyhan, B., Porter, E., Reifler, J., & Wood, T. J. (2019). Taking fact-checks literally but not seriously? the effects of journalistic fact-checking on factual beliefs and candidate favorability. *Political Behavior*, 1-22.

Pacheco, D., Hui, P.-M., Torres-Lugo, C., Truong, B. T., Flammini, A., & Menczer, F. (2020). *Uncovering coordinated networks on social media* (preprint No. arXiv:2001.05658). Retrieved from `https://arxiv.org/abs/2001.05658`

Pennycook, G., Bear, A., Collins, E. T., & Rand, D. G. (2019). The implied truth effect: Attaching warnings to a subset of fake news headlines increases perceived accuracy of headlines without warnings. *Management Science*, *0*(0), null. Retrieved from `https://doi.org/10.1287/mnsc.2019.3478` doi: 10.1287/mnsc.2019.3478

Pennycook, G., McPhetres, J., Zhang, Y., & Rand, D. (2020). *Fighting covid-19 misinformation on social media: Experimental evidence for a scalable accuracy nudge intervention* (Preprint). PsyArXiv. Retrieved from `https://doi.org/10.31234/osf.io/uhbk9` doi: 10.31234/osf.io/uhbk9

Pennycook, G., & Rand, D. G. (2019). Fighting misinformation on social media using crowdsourced judgments of news source quality. *Proceedings of the National Academy of Sciences*, *116*(7), 2521-2526.

Prez-Rosas, V., Kleinberg, B., Lefevre, A., & Mihalcea, R. (2017). *Automatic detection of fake news.*

Ratkiewicz, J., Conover, M., Meiss, M., Gonçalves, B., Flammini, A., & Menczer, F. (2011). Detecting and tracking political abuse in social media. In *Proc. 5th international aaai conference on weblogs and social media (icwsm).* Retrieved from `http://www.aaai.org/ocs/index.php/ICWSM/ICWSM11/paper/view/2850`

Romero, D. M., Meeder, B., & Kleinberg, J. (2011). Differences in the mechanics of information diffusion across topics: Idioms, political hashtags, and complex contagion on twitter. In *Proceedings of the 20th international conference on world wide web* (p. 695704). New York, NY, USA: Association for Computing Machinery. Retrieved from `https://doi.org/10.1145/1963405.1963503` doi: 10.1145/1963405.1963503

Roozenbeek, J., & van der Linden, S. (2019a). The fake news game: actively inoculating against the risk of misinformation. *Journal of Risk Research*, *22*(5), 570-580. Retrieved from `https://doi.org/10.1080/13669877.2018.1443491` doi: 10.1080/13669877.2018.1443491

Roozenbeek, J., & van der Linden, S. (2019b). Fake news game confers psychological resistance against online misinformation. *Palgrave Communications*, *5*(1), 65. Retrieved from `https://doi.org/10.1057/s41599-019-0279-9` doi: 10.1057/s41599-019-0279-9

Roozenbeek, J., van der Linden, S., & Nygren, T. (2020). Prebunking interventions based on "inoculation" theory can reduce susceptibility to misinformation across cultures. *Harvard Kennedy School Misinformation Review*, *1*(2).

Samory, M., & Mitra, T. (2018). Conspiracies online: User discussions in a conspiracy community following dramatic events. In *Twelfth international aaai conference on web and social media.*

Shao, C., Ciampaglia, G. L., Varol, O., Yang, K.-C., Flammini, A., & Menczer, F. (2018). The spread of low-credibility content by social bots. *Nature communications*, *9*(1), 1-9. doi: doi.org/10.1038/s41467-018-06930-7

Shao, C., Hui, P.-M., Wang, L., Jiang, X., Flammini, A., Menczer, F., & Ciampaglia, G. L. (2018). Anatomy of an online misinformation network. *PLoS ONE*, *13*(4), e0196087. Retrieved from `https://doi.org/10.1371/journal.pone.0196087` doi: 10.1371/journal.pone.0196087

Sharma, K., Seo, S., Meng, C., Rambhatla, S., Dua, A., & Liu, Y. (2020). Coronavirus on social media: Analyzing misinformation in twitter conversations. *arXiv preprint arXiv:2003.12309*.

Shen, C., Kasra, M., Pan, W., Bassett, G. A., Malloch, Y., & OBrien, J. F. (2019). Fake images: The effects of source, intermediary, and digital media literacy on contextual assessment of image credibility online. *New Media & Society*, *21*(2), 438-463. Retrieved from `https://doi.org/10.1177/1461444818799526` doi: 10.1177/1461444818799526

Shin, J., & Thorson, K. (2017, 02). Partisan Selective Sharing: The Biased Diffusion of Fact-Checking Messages on Social Media. *Journal of Communication*, *67*(2), 233-255. Retrieved from `https://doi.org/10.1111/jcom.12284` doi: 10.1111/jcom.12284

Starbird, K., Maddock, J., Orand, M., Achterman, P., & Mason, R. M. (2014). Rumors, false flags, and digital vigilantes: Misinformation on twitter after the 2013 boston marathon bombing. *IConference 2014 Proceedings*.

Stewart, L. G., Arif, A., & Starbird, K. (2018). Examining trolls and polarization with a retweet network. In *Proc. acm wsdm, workshop on misinformation and misbehavior mining on the web.*

Varol, O., Ferrara, E., Davis, C. A., Menczer, F., & Flammini, A. (2017). Online human-bot interactions: Detection, estimation, and characterization. In *Proc. intl. aaai conf. on web and social media (icwsm)* (pp. 280–289). Retrieved from `https://aaai.org/ocs/index.php/ICWSM/ICWSM17/paper/view/15587`

Varol, O., Ferrara, E., Menczer, F., & Flammini, A. (2017). Early detection of promoted campaigns on social media. *EPJ Data Science*, *6*(13). Retrieved from `http://rdcu.be/tWZN` doi: 10.1140/epjds/s13688-017-0111-y

Vosoughi, S., Roy, D., & Aral, S. (2018). The spread of true and false news online. *Science*, *359*(6380), 1146-1151. Retrieved from `https://science.sciencemag.org/content/359/6380/1146` doi: 10.1126/science.aap9559

Yang, K.-C., Varol, O., Davis, C. A., Ferrara, E., Flammini, A., & Menczer, F. (2019). Arming the public with artificial intelligence to counter social bots. *Human Behavior and Emerging Technologies*, *1*(1), 48-61. Retrieved from `https://onlinelibrary.wiley.com/doi/abs/10.1002/hbe2.115` doi: 10.1002/hbe2.115

Yang, K.-C., Varol, O., Hui, P.-M., & Menczer, F. (2020). Scalable and generalizable social bot detection through data selection. In *Proc. 34th aaai conf. on artificial intelligence (aaai).*

Yaqub, W., Kakhidze, O., Brockman, M. L., Memon, N., & Patil, S. (2020). Effects of credibility indicators on social media news sharing intent. In *Proceedings of the 2020 chi conference on human factors in computing systems* (p. 114). New York, NY, USA: Association for Computing Machinery. Retrieved from `https://doi.org/10.1145/3313831.3376213` doi: 10.1145/3313831.3376213

## Funding

## Competing interests

The authors have no competing interests to declare.

## Ethics

The mechanisms and procedures reported in this article were reviewed and approved by the Institutional Review Board (IRB) of the authors' institution.