

Does Debunking Work? Correcting COVID-19 Misinformation on Social Media

Timothy Caulfield^{*}

One of the defining characteristics of this pandemic has been the spread of misinformation.¹ Indeed, the World Health Organization famously called the crisis not just a pandemic, but also an “infodemic.”² It has been suggested, for example, that the coronavirus is both caused by 5G wireless technology and is a bioweapon. Cow urine and bleach have been put forward as cures. And enumerable wellness gurus have pushed immune boosting supplements and diets. All of this is science-free nonsense, of course. It has been suggested that this noise has already, *inter alia*, caused physical harm³ and financial loss,⁴ impacted health and science policy,⁵ added confusion and distraction to an already chaotic information environment,⁶ heightened stigma and prejudice,⁷ and made it more difficult to implement needed health policy initiatives.⁸

^{*} Canada Research Chair in Health Law and Policy, Professor, Faculty of Law and School of Public Health, Research Director, Health Law Institute, University of Alberta. I would like to thank Gordon Pennycook, Areeb Mian & Shujhat Khan, “Coronavirus: The Spread of Misinformation” (2020) 18:89 BMC Med, online: <doi.org/10.1186/s12916-020-01556-3>.

² World Health Organization, “Infodemic Management” (15 April 2020), online: <<https://www.who.int/teams/risk-communication/infodemic-management>>.

³ Alistair Smout & Paul Sandle, “Misinformation Ruins Lives, UK Fact-Checker Says,” *National Post* (30 April 2020), online: <<https://nationalpost.com/pmn/entertainment-pmn/misinformation-ruins-lives-uk-fact-checker-says>>.

⁴ Greg Iacurci, “Americans Have Lost \$13.4 Million to Fraud Linked to Covid-19,” *CNBC* (15 April 2020), online: <<https://www.cnn.com/2020/04/15/americans-have-lost-13point4-million-to-fraud-linked-to-covid-19.html>>.

⁵ Michael Liu et al, “Internet Searches for Unproven COVID-19 Therapies in the United States,” Research Letter (29 April 2020) JAMA Intern Med at E1: “Demand for chloroquine and hydroxychloroquine increased substantially following endorsements by high-profile figures and remained high even after a death attributable to chloroquine-containing products was reported,” DOI: <[10.1001/jamainternmed.2020.1764](https://doi.org/10.1001/jamainternmed.2020.1764)>.

⁶ See generally Amy Mitchell, J. Baxter Oliphant & Elisa Shearer, “About Seven-in-Ten U.S. Adults Say They Need to Take Breaks From COVID-19 News,” Pew Research Center (29 April 2020) at 4, where it was found that 86% believe that misinformation is causing either a great deal (49%) or some (37%) confusion about basic facts, online: <<https://www.journalism.org/2020/04/29/about-seven-in-ten-u-s-adults-say-they-need-to-take-breaks-from-covid-19-news/>>. See also Michael Sean Pepper & Stephanie Burton, “Sheer Volume of Misinformation Risks Diverting Focus from Fighting Coronavirus,” *The Conversation* (29 April 2020), online: <<https://theconversation.com/sheer-volume-of-misinformation-risks-diverting-focus-from-fighting-coronavirus-137408>>.

⁷ Harrison Mantas, “COVID-19 Infodemic Exacerbates Existing Religious and Racial Prejudices,” *Poynter* (1 May 2020) (“COVID-19 has inflamed fears of outsiders across the globe”), online: <<https://www.poynter.org/reporting-editing/2020/covid-19-infodemic-exacerbates-existing-religious-and-racial-prejudices/>>.

⁸ See also Leonardo Bursztyn et al, “Misinformation During a Pandemic” (2020) Becker Friedman Institute [working paper] at abstract: “While our findings cannot yet speak to long-term effects, they indicate that provision of misinformation in the early stages of a pandemic can have important consequences for how a disease ultimately affects the population.” See also Mian & Khan, *supra* note 1 at 2: “Public confusion leaves citizens unprepared for combatting a public health crisis.”

Much of this misinformation is spreading on social media,⁹ which has included the use of bots and strategic disinformation campaigns.¹⁰ It is worth noting that social media has also played a constructive role. It has, for instance, been used as a tool for communicating preventative strategies and mapping the spread of the virus.¹¹ And it has served as a primary source of news for many in the general public.¹² Indeed, more and more people are turning to social media to keep up-to-date on developments surrounding the pandemic.¹³ It has been reported that Twitter had about “12 million more daily users in the first three months of 2020 than in the last three of 2019.”¹⁴

Still, in the context of the “infodemic”, social media platforms have been the focus of much of the concern and policy activity.¹⁵ There is some suggestion that the spread of overt misinformation – that is, misinformation provided by known “fake news” sources – on some platforms, such as Facebook, has decreased since the implementation of platform counter measures, including removing fake accounts and tweaking their algorithm to reduce the reach of debunked articles.¹⁶ But on

⁹ See Soroush Vosoughi, Deb Roy & Sinan Aral, “The Spread of True and False News Online” (2018) 359:6380 *Science* 1141 at 1141, where the authors analyzed millions of social media shares and came to the grim conclusions that, “falsehood diffused significantly farther, faster, deeper, and more broadly than the truth in all categories of information,” DOI: <10.1126/science.aap9559>.

¹⁰ Ryan Ko, “Social Media Is Full of Bots Spreading COVID-19 Anxiety. Don’t Fall For It,” *Science Alert* (2 April 2020): “These fake accounts are common on Twitter, Facebook, and Instagram. They have one goal: to spread fear and fake news,” online: < <https://www.sciencealert.com/bots-are-causing-anxiety-by-spreading-coronavirus-misinformation>>.

¹¹ Katherine Ellison, “Social Media Posts and Online Searches Hold Vital Clues about Pandemic Spread” *Scientific American* (30 March 2020), online: <<https://www.scientificamerican.com/article/social-media-posts-and-online-searches-hold-vital-clues-about-pandemic-spread/>>.

¹² See, for example, Alaa Abd-Alrazaq et al, “Top Concerns of Tweeters During the COVID-19 Pandemic: Infoveillance Study” (2020) 22:4 *J Med Internet Res* e19016, where the authors analyzed 2.8 million tweets on the pandemic and found tweets on issues such as the source, cause, economic consequences, and treatments and cures, concluding: “Social media provides an opportunity to directly communicate health information to the public,” DOI:10.2196/19016.

¹³ Jeffrey Gottfried & Elisa Shearer, “News Use Across Social Media Platforms 2016” Pew Research Center (16 May 2016), online: < <https://www.journalism.org/2016/05/26/news-use-across-social-media-platforms-2016/>>.

¹⁴ Jon-Patrick Allem, “Social Media Fuels Wave of Coronavirus Misinformation as Users Focus on Popularity, Not Accuracy” *The Conversation* (6 April 2020), online: <<https://theconversation.com/social-media-fuels-wave-of-coronavirus-misinformation-as-users-focus-on-popularity-not-accuracy-135179>>. See also Vengattil Munsif & Dave Paresh, “Twitter Ad Sales Hit by Coronavirus but Active Users Soar” *Reuters* (23 March 2020), online: <<https://www.reuters.com/article/us-health-coronavirus-twitter/twitter-ad-sales-hit-by-coronavirus-but-active-users-soar-idUSKBN21A3HY>>.

¹⁵ Ramez Kouzy et al, “Coronavirus Goes Viral: Quantifying the COVID-19 Misinformation Epidemic on Twitter” (2020) 12:3 *Cureus* e7255, DOI: <10.7759/cureus.7255>.

¹⁶ Hunt Allcott, Matthew Gentzkow & Chuan Yu, “Trends in the Diffusion of Misinformation on Social Media” (2019) 6:2 *Res & Politics* 1 at abstract: “Our results suggest that the relative magnitude of the misinformation problem on Facebook has declined since its peak.” See also Paul Resnick, Aviv Ovadya & Garlin Gilchrist, “Iffy Quotient: A Platform Health Metric for Misinformation” (18 October 2018) School of Information Center for Social Media Responsibility, University of Michigan at 1: “there has been gradual improvement in Facebook’s Iffy Quotient since mid-2017, with a substantial cumulative impact. [...] In 2016 the Iffy sites’ share of attention was about twice as high on Facebook as Twitter; now it is

other platforms, including Twitter, the situation has gotten worse.¹⁷ And much of the misinformation about the coronavirus remains unchecked and continues to circulate, especially on Twitter.¹⁸

Why and how misinformation spreads and has an impact on behaviours and beliefs is a complex and multidimensional phenomenon.¹⁹ And there is an emerging rich academic literature on misinformation, particularly in the context of social media.²⁰ Here, I make no attempt to provide a comprehensive overview of that work. Rather, I focus on two relatively narrow questions: is debunking an effective strategy and, if so, what kind of counter-messaging is most effective? The goal of this article is to bring together relevant empirical research and expert commentary to: 1) serve as a resource and guide in the battle against misinformation (hence the heavy referencing) and 2) to stand as a defence of these efforts.²¹

Is it Worth It?

Let's start with two of most frequently raised arguments *against* vigorously countering the spread of misinformation. One is that correcting misinformation online is simply ineffective. Dumping more science on people has little impact, it is often said, because attempting to correct a misperception can cause individuals to become *more* entrenched in their beliefs. This phenomenon – usually called the “backfire effect” – has received a lot of attention and is often noted whenever there

50% higher on Twitter,” online: <<https://csmr.umich.edu/wp-content/uploads/2018/10/UMSI-CSMR-Iffy-Quotient-Whitepaper-810084.pdf>>.

¹⁷ Allcott, Gentzkow & Yu, *supra* note 16.

¹⁸ J. Scott Brennen et al, “Types, Sources, and Claims of COVID-19 Misinformation” Reuters Institute for the Study of Journalism, University of Oxford (7 April 2020) at 1: “On Twitter, 59% of posts rated as false in our sample by fact-checkers remain up,” online: <<https://reutersinstitute.politics.ox.ac.uk/types-sources-and-claims-covid-19-misinformation>>. See also Craig Timberg, “On Twitter, almost 60 Percent of False Claims about Coronavirus Remain Online — Without a Warning Label” *Washington Post* (7 April 2020), online: <<https://www.washingtonpost.com/technology/2020/04/07/twitter-almost-60-percent-false-claims-about-coronavirus-remain-online-without-warning-label/>>.

¹⁹ Dietram A Scheufele & Nicole M Krause, “Science Audiences, Misinformation, and Fake News” (2019) 116: 16 PNAS 7662 at 7662: “we show how being misinformed is a function of a person’s ability and motivation to spot falsehoods, but also of other group-level and societal factors that increase the chances of citizens to be exposed to correct(ive) information,” DOI: <10.1073/pnas.1805871115>.

²⁰ See generally Yuxi Wang et al, “Systematic Literature Review on the Spread of Health-related Misinformation on Social Media” (2019) 240:112552 Social Science & Medicine 1 at 1: “Overall, we observe an increasing trend in published articles on health-related misinformation and the role of social media in its propagation,” DOI: <10.1016/j.socscimed.2019.112552>. See also Denise-Marie Ordway, “Fake News and the Spread of Misinformation: A Research Roundup” *Journalist’s Resource* (1 September 2017), online: <<https://journalistsresource.org/studies/society/internet/fake-news-conspiracy-theories-journalism-research/>>.

²¹ The word “debunking” is less than ideal as some may feel it fails to capture the need to listen to and engage the public. It can also be associated with a more aggressive, or mocking, approach (a strategy which I criticize below). However, in total, with those critiques noted, I still feel it is a good catchall term that, as defined by Amy Sippitt, can be used to refer to “factual messages which seek to rebut inaccurate factual claims.” See Amy Sippitt, “The Backfire Effect: Does It Exist? And Does It Matter for Factcheckers?” *Full Fact* (March 2019) at 7, online: <<https://fullfact.org/blog/2019/mar/does-backfire-effect-exist/>>.

is a call for more individuals to get actively involved in the countering of misinformation. Debunking doesn't work, it is argued.²²

But how strong is the backfire phenomenon? There are several well-known studies associated with the birth of this concern. Probably the most influential is a study published in 2010 where the researchers explored the impact of corrected news articles that contained a misleading claim by a politician. It was found that "corrections frequently fail to reduce misperceptions among the targeted ideological group" and there were "several instances of a 'backfire effect' in which corrections actually increase misperceptions among the group in question."²³ As a result of this and a few other studies, there now seems to be a widely accepted belief that the backfire effect is a dominant phenomenon that makes debunking a near futile exercise.²⁴

In reality, the backfire effect seems to be a relatively rare occurrence.²⁵ Indeed, the lead author of the 2010 study, Brendan Nyhan, has noted that their results have often "been overstated and oversold,"²⁶ in part because their conclusions may be quite context specific.²⁷ A 2019 comprehensive analysis of the available research concluded that the existing body of evidence, much of it published after the 2010 study, found no backfire effect and that "most recent studies now suggest that generally debunks can make beliefs in specific claims more accurate."²⁸ For example, a study published in 2019 found that "evidence of factual backfire is far more

²² See, for example, Christian Bokhove, "Beware: Debunking Research Myths Can Backfire on You" *Tes* (19 July 2019), online: <<https://www.tes.com/magazine/article/beware-debunking-research-myths-can-backfire-you>>.

²³ Brendan Nyhan & Jason Reifler, "When Corrections Fail: The Persistence of Political Misperceptions," (2010) 32 *Polit Behav* 303, DOI: <10.1007/s11109-010-9112-2>.

²⁴ See, for example, Julie Beck, "This Article Won't Change Your Mind," *The Atlantic* (11 December 2019), online: <<https://www.theatlantic.com/science/archive/2017/03/this-article-wont-change-your-mind/519093/>>; and, "The Backfire Effect: Why Facts Don't Win Arguments," *Big Think* (15 October 2013), online: <https://bigthink.com/think-tank/the-backfire-effect-why-facts-dont-win-arguments>. See also Erin Brodwin, "Facebook's Covid-19 Misinformation Campaign Is Based on Research. The Authors Worry Facebook Missed the Message," *StatNews* (1 May 2020) where it is noted that Facebook's coronavirus misinformation strategy is "designed to avoid what's known as the backfire effect," online: <<https://www.statnews.com/2020/05/01/facebooks-covid-19-misinformation-campaign-is-based-on-research-the-authors-worry-facebook-missed-the-message/>>. Why the "backfire effect" gained so much traction is an interesting question on its own, one which is beyond the scope of this piece. But I think that the fact it feels intuitively correct is a big part of its appeal. It is hard to change opinions.

²⁵ Indeed, some have gone so far as to call its existence a myth. See, for example, Laura Hazard Owen, "The 'Backfire Effect' Is Mostly a Myth" *NiemanLab* (22 March 2019), online: <<https://www.niemanlab.org/2019/03/the-backfire-effect-is-mostly-a-myth-a-broad-look-at-the-research-suggests/>>.

²⁶ See 8 January 2018, tweet by lead author, Brendan Nyhan, where he states: "the research findings, including accounts of my own backfire effect paper with @jasonreifler, have often been overstated and oversold," online: <<https://twitter.com/brendannyhan/status/948544775799607296?lang=en>>.

²⁷ For example, see Sippitt, *supra* note 21 at 10, who notes that the experiment "purposefully covered a highly controversial topic in American politics [WMD in Iraq] where people would have prior beliefs" and as such "it's arguably unsurprising that individuals were unpersuaded by a single news item."

²⁸ See *ibid* at 5.

tenuous than prior research suggests. By and large, citizens heed factual information, even when such information challenges their ideological commitments.”²⁹ Another study from 2019 found that “debunking” works – if done using appropriate strategies (more on that below) – and “no evidence” that “rebutting science denialism in public discussions backfires, not even in vulnerable groups (for example, US conservatives).”³⁰ To be fair, motivated reasoning (constructing rationales to fit a pre-existing position) and other cognitive biases (e.g., confirmation bias) have been shown to influence what information we see online and elsewhere.³¹ Still, for many areas of science, at least some research has found that differences in scientific belief are driven mostly by levels of science knowledge and not motivated reasoning.³² So while a backfire effect may occur in some circumstances – this is an area where more research would be helpful – it certainly isn’t such a robust and measurable phenomenon that it should stop us from mounting efforts to counter misinformation on social media.

The second and perhaps more challenging critique of correcting and debunking is that it may inadvertently help to spread the misinformation.³³ Specifically, there might an “illusory truth” effect.³⁴ Studies have consistently found that merely exposing people to an idea increases the believability of that idea.³⁵ In many ways this is how “fake news” works.³⁶ A study by Gordon Pennycook, et al., for example,

²⁹ Thomas Wood & Ethan Porter, “The Elusive Backfire Effect: Mass Attitudes’ Steadfast Factual Adherence,” (2019) 41 *Polit Behav* 135.

³⁰ Philipp Schmid & Cornelia Betsch, “Effective Strategies for Rebutting Science Denialism in Public Discussions” (2019) 3 *Nat Hum Behav* 931 at abstract.

³¹ For example, see Dan Kahan, “The Politically Motivated Reasoning Paradigm, Part 1: What Politically Motivated Reasoning Is and How to Measure It” in RA Scott and SM Kosslyn eds, *Emerging Trends in the Social & Behavioral Sciences* (Wiley Library Online, 2016), DOI: <10.1002/9781118900772.etrds0417>.

³² Jonathon McPhetres & Gordon Pennycook, “Science Beliefs, Political Ideology, And Cognitive Sophistication,” (2020) OSF Preprints at abstract: “We also found very little evidence of motivated reasoning: reasoning ability was instead broadly associated with pro-science beliefs. Finally, one’s level of basic science knowledge was the most consistent predictor of people’s beliefs about science. Results suggest educators and policymakers should focus on increasing basic science literacy and critical thinking rather than the ideologies that purportedly divide people,” online: <<https://osf.io/ad9v7/>>.

³³ This is also often called the backfire effect, though it is different phenomenon than that described in the Nyhan & Reifler, “When Corrections Fail,” which coined the phrase. As such, I usually treat them as distinct and refer to this as the “spreading” concern.

³⁴ Melissa Healy, “Misinformation About the Coronavirus Abounds, but Correcting It Can Backfire” *Los Angeles Times* (8 February 2020), “Sometimes the effort to correct misinformation involves repeating the lie. That repetition seems to establish it in our memories more firmly than the truth,” online: <<https://www.latimes.com/science/story/2020-02-08/coronavirus-outbreak-false-information-psychology>>.

³⁵ See Jonas De keersmaecker, David Dunning & Gordon Pennycook. “Investigating the Robustness of the Illusory Truth Effect Across Individual Differences in Cognitive Ability, Need for Cognitive Closure, and Cognitive Style” (2020) 46:2 *Personality and Social Psychology Bulletin* 204. Indeed, this effect can still have an impact even if the information runs counter to an existing knowledge base. See, for example, Lisa K. Fazio et al, “Knowledge Does Not Protect Against Illusory Truth” (2015) 144 *J Experimental Psychology* 993 at 993: “Contrary to prior suppositions, illusory truth effects occurred even when participants knew better.”

³⁶ See, for example, Danielle C. Polage “Making up History: False Memories of Fake News Stories” (2012) 8:2 *Europe’s J Psychology* 245; and Christopher Paul & Miriam Matthews, “The Russian ‘Firehose of Falsehood’ Propaganda Model: Why It Might Work and Options to Counter It” (2016) RAND, online:

found that even a single exposure to misinformation could increase subsequent perceptions of accuracy.³⁷

So, does this mean that debunking misinformation and conspiracy theories on social media – which often, of necessity, will include a restatement of the problematic belief – has the potential to do more harm than good? While the speculation about the problem of spreading is rooted in evidence about the possible impact of exposure to misinformation, there does not appear to be much direct empirical evidence that debunking actually has this problematic impact. Indeed, a recent study (still in preprint at time of this writing) explored this exact concern by analyzing whether a debunking of a new piece of misinformation – that is, a not widely known and novel myth or conspiracy theory – led to an increase in beliefs about the claim. They found that that corrections that “repeated novel misinformation claims did not lead to stronger misconceptions compared to a control group never exposed to the false claims or corrections.”³⁸ As a result of this finding – which fits with other work on point³⁹ – the authors come to the conclusion that “it is safe to repeat misinformation when correcting it, even when the audience might be unfamiliar with the misinformation.”⁴⁰

The timing of a correction may also be relevant here. Claire Wardle, executive director of an institute dedicated to fighting misinformation, has suggested that if you debunk a bit of misinformation too early you may give it unintended oxygen and allow it to spread further.⁴¹ But once the public awareness of a particular myth,

<<https://www.rand.org/pubs/perspectives/PE198.html>>. I have argued that this is also one reason that celebrities can have such a large impact on the spread of misinformation. See, for example, Timothy Caulfield, “Celebrities like Gwyneth Paltrow Made the 2010s the Decade of Health and Wellness Misinformation” *NBC News* (27 December 2019), online: <<https://www.nbcnews.com/think/opinion/celebrities-gwyneth-paltrow-made-2010s-decade-health-wellness-misinformation-ncna1107501>>. See also Mathew Ingram, “Amplifying the Coronavirus Protests” *Columbia Journalism Review* (22 April 2020) where it is noted that less-than-ideal reporting of lockdown protests may have given them more legitimacy than the objective numbers might have suggested was appropriate, online: <https://www.cjr.org/the_media_today/amplifying-coronavirus-protests.php>.

³⁷ Gordon Pennycook, Tyrone D Cannon & David G Rand, “Prior Exposure Increases Perceived Accuracy of Fake News” (2018) 147:12 *J Exp Psychol Gen* 1865, DOI: <10.1037/xge0000465>.

³⁸ Ullrich KH Ecker, Stephan Lewandowsky & Matthew Chadwick, “Can Corrections Spread Misinformation to New Audiences? Testing for the Elusive Familiarity Backfire Effect” (2020) [working paper], DOI: <10.31219/osf.io/et4p3>.

³⁹ Ullrich KH Ecker et al, “The Effectiveness of Short-Format Refutational Fact-Checks” (2020) 111:1 *British J Psychology* 36 at 36: “we found no evidence for a familiarity-driven backfire effect.”

⁴⁰ *Ibid.*

⁴¹ Claire Wardle, “What Role Should Newsrooms Play in Debunking COVID-19 Misinformation?” *Nieman Reports* (8 April 2020), online: <<https://niemanreports.org/articles/what-role-should-newsrooms-play-in-debunking-covid-19-misinformation/>>. See also Whitney Phillips, *The Oxygen of Amplification: Better Practices for Reporting on Extremists, Antagonists, and Manipulators Online*, Data & Society (2012), online: <<https://datasociety.net/library/oxygen-of-amplification/>>; and Susan Benkelmam, “Getting it Right: Strategies for Truth-Telling in a Time of Misinformation and Polarization” *American Press Institute* (11 December 2019): “Journalists must ask themselves whether a falsehood has become so significant that it needs to be knocked down,” online:

conspiracy theory or item of misinformation hits a tipping point – that is, the item is starting to be shared more widely – it is important to vigorously counter. If we wait too long to attempt a correction, it may become increasingly difficult to stop the momentum of the misinformation.⁴² Once a conspiracy theory gets a strong foothold in the public conscious, it can be difficult to dislodge – as we have seen with issues like the myths surrounding vaccination.

The better interpretation of the existing literature, I think, is that while we need to be cognizant of the spreading concern, the evidence is far from definitive and what evidence is available suggests it often doesn't happen. There are, of course, many other challenges associated with efforts to correct misinformation, such as the possibility for a range of additional unintended consequences (e.g., general warning tags skewing how people perceive legitimate news).⁴³ But despite the need for more research, there is nothing in the existing research to suggest debunking is a futile exercise. On the contrary, as we will see, there is a growing body of evidence that tells us correcting misinformation should be viewed as a vitally important science and health policy activity.

What Kind of Counter-Messaging Works?

As with the research on the challenges associated with correcting misinformation, the data surrounding effective debunking strategies is messy and context dependant. More research on how best to deal with misinformation is clearly needed,⁴⁴ but there is little doubt that countering misinformation can have a positive impact.⁴⁵ Indeed, silence in the face of misinformation seems likely to be the

<<https://www.americanpressinstitute.org/publications/reports/strategy-studies/truth-telling-in-a-time-of-misinformation-and-polarization/>>.

⁴² There is some recent evidence to support this view. See, for example, Wasim Ahmed et al, “COVID-19 and the 5G Conspiracy Theory: Social Network Analysis of Twitter Data” (2020) 22:5 J Med Internet Res e19458 at abstract, found that early “there was a lack of an authority figure who was actively combating such [5g] misinformation” on social media. What is needed, they conclude, is the “combination of quick and targeted interventions oriented to delegitimize the sources of fake information is key to reducing their impact” (at abstract).

⁴³ John M Carey et al, “The Effects of Corrective Information about Disease Epidemics and Outbreaks: Evidence from Zika and Yellow Fever in Brazil” (2020) 6:5 Science Advances 1 at 9: “a general warning about the presence of fake news has been found to decrease belief in the accuracy of both false and legitimate news headlines,” DOI: <10.1126/sciadv.aaw7449>. And for a study that found the opposite effect, see Gordon Pennycook et al, “The Implied Truth Effect: Attaching Warnings to a Subset of Fake News Headlines Increases Perceived Accuracy of Headlines Without Warnings” (2020) Management Science [forthcoming], online: < <https://dx.doi.org/10.2139/ssrn.3035384>>. While placing “fake news” warnings on social media content can have a positive impact, this study found that “the presence of warnings caused untagged headlines to be seen as more accurate than in the control” (at abstract).

⁴⁴ See Gordon Pennycook & David Rand, “The Right Way to Fight Fake News” *New York Times* (24 March 2020): “The obvious conclusion to draw from all this evidence is that social media platforms should rigorously test their ideas for combating fake news and not just rely on common sense or intuition about what will work,” online: <<https://www.nytimes.com/2020/03/24/opinion/fake-news-social-media.html>>.

⁴⁵ For the benefits of debunking in the context of a pandemic, see Toni GLA van der Meer & Yan Jin, “Seeking Formula for Misinformation Treatment in Public Health Crises: The Effects of Corrective Information Type and Source” (2020) 35:5 Health Commun 560 at 560: “Results show that, if corrective

worst strategy. A 2019 study, for example, found that not responding to misinformation “has a negative effect on attitudes towards behaviours favoured by science.”⁴⁶ But what kind of social media counter is likely to have the biggest positive result? Below is a list of some of the general themes that have emerged in the research regarding the tone and style of debunking messaging that is relevant to all social media platforms. Here, I am focusing on just the actual content of a social media debunk. Obviously, not every approach will work for every corrective message – a Tweet is, after all, just 280 characters. But these evidence-informed general principles can help to maximize the impact efforts to correct online misinformation.

First, use facts. Despite all the concern regarding the impotence of facts to change minds, most studies have found that providing corrective information can be effective,⁴⁷ especially if the alternative explanation – that is, the science-informed facts – fills in the gap in understanding caused by the debunk and (when appropriate and possible) provides a causal explanation.⁴⁸ This approach can also nudge people to think more critically generally, which may help to shield them against related forms of misinformation.⁴⁹

information is present rather than absent, incorrect beliefs based on misinformation are debunked and the exposure to factual elaboration, compared to simple rebuttal, stimulates intentions to take protective actions.” See generally Nathan Walter & Sheila T Murphy, “How to Unring the Bell: A Meta-Analytic Approach to Correction of Misinformation” (2018) 85:3 *Communications Monographs* 423 a meta-analysis of existing data that concludes: “corrective attempts can reduce misinformation across diverse domains, audiences, and designs” (at 436); Man-pui Sally Chan et al, “Debunking: A Meta-Analysis of the Psychological Efficacy of Messages Countering Misinformation” (2017) 28:11 *Psychological Science* 1531; Brendan Nyhan et al, “Taking Fact-Checks Literally But Not Seriously? The Effects of Journalistic Fact-Checking on Factual Beliefs and Candidate Favorability” (2019) *Polit Behav* [forthcoming], online: <<https://doi.org/10.1007/s11109-019-09528-x>>; and Victoria L Rubin, “Deception Detection and Rumor Debunking for Social Media” in L Sloan & A Quan-Haase eds, *The SAGE Handbook of Social Media Research Methods* (London: SAGE, 2017).

⁴⁶ Schmid and Betsch, *supra* note 30 at abstract.

⁴⁷ Leticia Bode & Emily K Vraga, “In Related News, That Was Wrong: The Correction of Misinformation Through Related Stories Functionality in Social Media” (2015) 65:4 *J Communication* 619 at 630: “Our experimental evidence suggests that attitude change related to GMOs can be achieved with regard to misperceptions by virtue of exposure to corrective information within social media.” See also Emily Falk & Molly Crockett, “You Can Help Slow the Virus if You Talk about it Accurately Online” *Washington Post* (28 April 2020), online: <<https://www.washingtonpost.com/outlook/2020/04/28/you-can-help-slow-virus-if-you-talk-about-it-accurately-online/>>; and *ibid*.

⁴⁸ See Walter & Murphy, *supra* note 45 at 436: “corrective messages that integrate retractions with alternative explanations (i.e., coherence) emerge as an effective strategy to debunk falsehoods.” See also Briony Swire & Ullrich Ecker, “Misinformation and its Correction: Cognitive Mechanisms and Recommendations for Mass Communication” in Brian G. Southwell, Emily A. Thorson & Laura Sheble eds, *Misinformation and Mass Audiences* (Austin: University of Texas Press, 2018): The alternative explanation effectively plugs the model gap left by the retraction. See also Brendan Nyhan & Jason Reifler, “Displacing Misinformation about Events: An Experimental Test of Causal Corrections” (2015) 2:1 *J Experimental Political Science* 81.

⁴⁹ See Ecker et al, *supra* note 39 at 49: “We can thus conclude that embedding a rebuttal in a fact-oriented context has beneficial implications beyond specific belief reduction, fostering a more sceptical and evidence-based approach to the issue at hand.”

Second, provide clear, straightforward and shareable content.⁵⁰ Studies have shown that the use of scientific jargon will cause people to disengage, even if explanatory language is also provided in the text.⁵¹

Third, use trustworthy and independent sources. Evidence perceived to be removed from an agenda (and the profit motive) is more likely to be trusted and persuasive.⁵² While it can be a challenge to find sources that are trusted by all – there has been a significant erosion in trust in many public institutions⁵³ – public health authorities and independent scientists still retain a relatively high level of trustworthiness, particularly during times of crisis.⁵⁴

Fourth, if applicable and available, emphasize the scientific consensus.⁵⁵ Ideally, this tactic should be accompanied by a recognition that science evolves and, as such, the consensus can change.

⁵⁰ Samantha Yammine, “Going Viral: How to Boost the Spread of Coronavirus Science on Social Media,” *Nature Careers Community* (5 May 2020), online: <<https://www.nature.com/articles/d41586-020-01356-y>>.

⁵¹ See, for example, Hillary C Shulman et al, “The Effects of Jargon on Processing Fluency, Self-Perceptions, and Scientific Engagement” (2020) *J Language and Social Psychology* 1 at 13: “Jargon can then serve as exclusionary language that disengages meaningful relationships between public and expert communities from forming,” online: <<https://doi.org/10.1177%2F0261927X20902177>>.

⁵² Susan T Fiske & Cydney Dupree, “Gaining Trust as Well as Respect in Communicating to Motivated Audiences about Science Topics” (2014) 111:4 *PNAS* 13593.

⁵³ Timothy Caulfield, “Now More Than Ever, We Must Fight Misinformation. Trust in Science Is Essential” *Globe and Mail* (20 March 2020). Not surprisingly, studies have found that debunking has a more modest effect if people view the original source of misinformation favourably. But even in this situation, debunking efforts can help. See Jeong-woo Jang, Eun-Ju Lee & Soo Yun Shin, “What Debunking of Misinformation Does and Doesn’t” (2019) 22:6 *Cyberpsychology, Behavior, & Social Networking* 423 at 426: “Overall, the results showed that when the falsehood of information was exposed, participants became less favorable toward the immediate source who shared the misinformation, but their initial source attitude also moderated their reactions by inducing different attribution processes,” online: <<https://doi.org/10.1089/cyber.2018.0608>>. For another commentary on the impact of low trust see Mike Caulfield, “Cynicism, Not Gullibility, Will Kill Our Humanity,” *Hapgood* (27 November 2018) Digital Polarization Initiative, online: <<https://hapgood.us/2018/11/27/cynicism-not-gullibility-will-kill-our-humanity/>>.

⁵⁴ See Pew Research Centre, “Public Holds Broadly Favorable Views of Many Federal Agencies, Including CDC and HHS” (9 April 2020) “Currently, 79% of U.S. adults express a favorable opinion of the CDC...”; and Hannah Fingerhut, “AP-NORC poll: High use, mild trust of news media on COVID-19” (30 April 2020) Associated Press: “Americans are especially likely to trust information about the coronavirus that comes from the CDC or from personal health care providers,” online: <<https://www.people-press.org/2020/04/09/public-holds-broadly-favorable-views-of-many-federal-agencies-including-cdc-and-hhs/>>. See van der Meer & Jin, *supra* note 45 at 560 where it is summarized that during times of crisis “government agency and news media sources are found to be more successful in improving belief accuracy compared to social peers.”

⁵⁵ See Sander L van der Linden, Chris E Clarke & Edward W Maibach, “Highlighting Consensus among Medical Scientists Increases Public Support for Vaccines: Evidence from a Randomized Experiment” (2015) 15:1207 *BMC Public Health*; Jeremy D Sloane & Jason R Wiles, “Communicating the Consensus on Climate Change to College Biology Majors: The Importance of Preaching to the Choir” (2020) 10:2 *Ecology and Evolution* 594; Sander L van der Linden et al, “The Scientific Consensus on Climate Change as a Gateway Belief: Experimental Evidence” 10:2 *PLoS ONE* e0118489, DOI: <10.1371/journal.pone.0118489>; and Sander L van der Linden, “Why Doctors Should Convey the Medical Consensus on Vaccine Safety” (2016) 21:3 *Evidence Based Medicine* 119,

Fifth, be nice and be authentic. Research has found that an aggressive language style is perceived to be both less credible and less trustworthy.⁵⁶ Don't shame, ridicule or marginalize members of the public who are looking for answers (though I have less patience for those pushing bunk for profit, brand enhancement and ideological spin).⁵⁷ In addition, messaging that comes from someone that is seen to be a unique and authentic individual – that is, not just a talking head associated with an institution – can also enhance trust, credibility, and the persuasiveness of the message.⁵⁸

Sixth, consider using a narrative. Humans are wired to respond to stories.⁵⁹ Indeed, there is some evidence that an engaging anecdote can overwhelm our ability to think scientifically.⁶⁰ This is one reason that testimonials are such an effective strategy for the marketing of unproven therapies.⁶¹ But a narrative can also be used to convey science – and information about critical thinking and the scientific process⁶² – in a way that is compelling and memorable.⁶³

Seventh, emphasize the gaps in logic and the flawed strategies used by those pushing misinformation. Several studies have found that using rational arguments, such as highlighting the rhetorical tools used to spread misinformation (reliance on conspiracy theories, misrepresentation of risks, use of false “experts”, etc.), can be an effective debunking strategy.⁶⁴

DOI: <10.1136/ebmed-2016-110435>.

⁵⁶ See Lars König & Regina Jucks, “Hot Topics in Science Communication: Aggressive Language Decreases Trustworthiness and Credibility in Scientific Debates” (2019) 28:4 Public Understanding of Science 401; see also Fisk & Dupree, *supra* note 52.

⁵⁷ Anand Ram, “How to (Tactfully) Discourage Spread of False Pandemic Information,” *CBCNews* (19 April 2020) where misinformation expert, Claire Wardle, notes the value of being empathetic and using words that “put yourself in the same perspective,” online: <<https://www.cbc.ca/news/canada/covid-19-misinformation-rumour-1.5532302>>.

⁵⁸ See Lise Saffran et al, “Constructing and Influencing Perceived Authenticity in Science Communication” (2020) 15:1 PLoS ONE e0226711; and Sara Reardon, “Adding a Personal Backstory Could Boost Your Scientific Credibility with the Public,” *Nature Career News* (2020), DOI: <10.1038/d41586-020-00857-0>.

⁵⁹ Michael F Dahlstrom, “Using Narratives and Storytelling to Communicate Science with Nonexpert Audiences” (2014) 111:4 PNAS 13614.

⁶⁰ Fernando Rodriguez et al, “Examining The Influence Of Anecdotal Stories And The Interplay Of Individual Differences On Reasoning,” (2016) 22:3 Thinking & Reasoning 274 at 274: “anecdotal stories decreased the ability to reason scientifically even when controlling for education level and thinking dispositions.”

⁶¹ Bethany Hawke et al, “How to Peddle Hope: An Analysis of YouTube Patient Testimonials of Unproven Stem Cell Treatments” (2019) 12:6 Stem Cell Reports 1186.

⁶² See Michael F Dahlstrom & Dietram A Scheufele, “(Escaping) the Paradox of Scientific Storytelling” (2018) 16:10 PLoS Biology e2006720: “narratives might have most of their power not in conveying facts or building excitement but in rebuilding the foundation of understanding scientific reasoning,” online: <<https://doi.org/10.1371/journal.pbio.2006720>>.

⁶³ For an overview of the evidence on point, see Timothy Caulfield et al, “Health Misinformation and the Power of Narrative Messaging in the Public Sphere” (2019) 2:2 Can J Bioethics 52.

⁶⁴ See Schmid & Betsch, *supra* note 30; Stephan Lewandowsky and John Cook, *The Conspiracy Theory Handbook* (Fairfax: George Mason University, 2020); and Gábor Orosz et al. “Changing Conspiracy

Eighth, lead with the correct information, not the misinformation. While the evidence on the spreading concern is mixed, it makes sense to frame debunking in a manner that makes the correct information – not the misinformation, myth or conspiracy theory – the memorable part of the messaging.⁶⁵ Make sure the misinformation is clearly flagged as wrong so the debunk is the key takeaway.

Finally, the audience should be the general public, not the hard-core believer. And this should be the case even if the debunk is triggered by information circulated by hard-core believers or those pushing misinformation for personal gain.⁶⁶ It is very difficult to change the mind of someone who is heavily invested in a particular myth or conspiracy theory. As noted by the World Health Organization, the probability of changing a vocal science denier is very low.⁶⁷ As such, the corrective information should be framed as if the general public is listening.

Empowering Users

Fighting the spread of misinformation will, of course, require more than just carefully crafted debunks on social media. We need to come at this issue from every angle.⁶⁸ We need, for instance, social media platforms to adopt evidence-informed strategies that will both remove the most harmful content and heighten user vigilance. Studies have found, for example, that the use of warning tags – like “rated false” – on social media posts can be an effective strategy to inform the public about potential problems with accuracy with specific content.⁶⁹ And we need a more

Beliefs through Rationality and Ridiculing” (2016) 7:1525 *Frontiers in Psychology* at 8: “uncovering arguments regarding the logical inconsistencies of CT beliefs can be an effective way to discredit them.”

⁶⁵ Some have called this the “truth sandwich” strategy. See Benkelmam, *supra* note 41 at sum: “There are a number of strategies for reporting on falsehoods without amplifying them. One is the ‘truth sandwich,’ which involves stating a true fact, then the falsehood, then the true fact again.”

⁶⁶ I will often use a pop culture moment – the spread of misinformation by a celebrity, for example – as an opportunity to create sharable content about science and the problems associated with the spread of health misinformation.

⁶⁷ World Health Organization, “Best Practices Guidance: How to Respond to Vocal Vaccine Deniers in Public” (Copenhagen: Regional Office for Europe of the World Health Organization, 2016), “Rule 1: The general public is your target audience, not the vocal vaccine denier.”

⁶⁸ See, for example, Kate Starbird, “Disinformation’s Spread: Bots, Trolls and All of Us” (2019) 571 *Nature World View* 449: “But effective disinformation campaigns involve diverse participants; they might even include a majority of ‘unwitting agents’ who are unaware of their role,” DOI: <10.1038/d41586-019-02235-x>.

⁶⁹ Katherine Clayton et al, “Real Solutions for Fake News? Measuring the Effectiveness of General Warnings and Fact-Check Tags in Reducing Belief in False Stories on Social Media” (2019) *Polit Behav* at abstract: “indicate that false headlines are perceived as less accurate when people receive a general warning.” While warning tags seem to have a role to play, they need to be deployed sensibly. Research has found, for example, that general warnings telling readers to beware of misinformation can have an unintended spillover of effect of decreasing “belief in the accuracy of true headlines,” online: <<https://doi.org/10.1007/s11109-019-09533-0>>. See Pennycook et al, *supra* note 43 highlights that using warning tags can lead to an inappropriate implication that posts without warnings are *more* accurate. See also Melanie Freeze et al, “Fake Claims of Fake News: Political Misinformation, Warnings, and the Tainted Truth Effect” (2020) *Polit Behav*, online: <<https://doi.org/10.1007/s11109-020-09597-3>>.

robust policy response against those pushing unproven products and ideas on social media in a manner that infringes existing laws and regulations.⁷⁰

Perhaps the most important strategy will be to empower people with the tools necessary to be more critical information consumers. This should incorporate teaching both critical thinking skills and media literacy,⁷¹ including inoculating (or “pre-bunking”) people against misinformation⁷² and simply reminding them to think about accuracy before sharing.⁷³ A growing body of literature has found that, in general, people want to be accurate and want to share only factual material.⁷⁴ Most users do not fall for or share misinformation due to a monovalent agenda or, even, a partisan bias.⁷⁵ As such, if we can nudge people to think about accuracy prior to sharing social media content we may be able to have a significant impact on the spread of misinformation.⁷⁶ A 2020 study that specifically looked at misinformation in the context of the coronavirus found exactly this effect, concluding that “nudging

⁷⁰For example regulatory action see Health Canada, Advisory, RA-72659, “Health Products that Make False or Misleading Claims to Prevent, Treat or Cure COVID-19 May Put Your Health at Risk” (27 March 2020); Federal Trade Commission, Press Release, “FTC Sends 45 More Letters Warning Marketers to Stop Making Unsupported Claims That Their Products and Therapies Can Effectively Prevent or Treat COVID-19” (7 May 2020).

⁷¹ See, for example, Michelle A Amazeen & Erik P Bucy, “Conferring Resistance to Digital Disinformation: The Inoculating Influence of Procedural News Knowledge” (2019) 63:3 J Broadcasting & Electronic Media 415 at 429: “additional educational campaigns to inform citizens about mainstream news media operations could yield significant benefits.” And see Viren Swami et al, “Analytic Thinking Reduces Belief in Conspiracy Theories” (2014) 133:3 Cognition 572.

⁷² See, for example, Jon Roozenbeek & Sander van der Linden, “The New Science of Prebunking: How to Inoculate against the Spread of Misinformation” (7 October 2019) BMC On Society, online: <<http://blogs.biomedcentral.com/on-society/2019/10/07/the-new-science-of-prebunking-how-to-inoculate-against-the-spread-of-misinformation/>>; and Jon Roozenbeek & Sander van der Linden, “Fake News Game Confers Psychological Resistance against Online Misinformation” (2019) 5:65 Palgrave Commun at abstract: “We provide initial evidence that people’s ability to spot and resist misinformation improves after gameplay [which teaching about misinformation], irrespective of education, age, political ideology, and cognitive style,” online: <<https://doi.org/10.1057/s41599-019-0279-9>>.

⁷³ Bence Bago, David G Rand & Gordon Pennycook, “Fake News, Fast and Slow: Deliberation Reduces Belief in False (But Not True) News Headlines” J Experimental Psychology: General. Advance online publication, at abstract: “Our data suggest that, in the context of fake news, deliberation facilitates accurate belief formation and not partisan bias,” online: <<https://www.ncbi.nlm.nih.gov/pubmed/31916834>>.

⁷⁴ Emma Young, “Most People Who Share ‘Fake News’ Do Care About the Accuracy of News Items — They’re Just Distracted” (16 January 2020) Research Digest (The British Psychological Society), online: <<https://digest.bps.org.uk/2020/01/16/most-people-who-share-fake-news-do-care-about-the-accuracy-of-news-items-theyre-just-distracted/>>.

⁷⁵ Gordon Pennycook and David G Rand, “Lazy, Not Biased: Susceptibility to Partisan Fake News is Better Explained by Lack of Reasoning Than By Motivated Reasoning” (2019) 188 Cognition 39 at abstract: “Our findings therefore suggest that susceptibility to fake news is driven more by lazy thinking than it is by partisan bias per se – a finding that opens potential avenues for fighting fake news.”

⁷⁶ See, for example, Lisa Fazio, “Pausing to Consider Why a Headline is True or False Can Help Reduce the Sharing of False News” (10 February 2020) Misinformation Review: “This research suggests that forcing people to pause and think can reduce shares of false information,” online: <<https://misinfoeview.hks.harvard.edu/article/pausing-reduce-false-news/>>; and Gordon Pennycook et al, “Understanding and Reducing the Spread of Misinformation Online” (25 November 2019) [working paper] at abstract: “we find that subtly inducing people to think about the concept of accuracy increases the quality of the news they share,” online: <<https://psyarxiv.com/3n9u8/>>.

people to think about accuracy is a simple way to improve choices about what to share on social media.”⁷⁷

Conclusion

There is a growing body of research on both the phenomenon of online misinformation and the best way counter it. While the data remains complex and, at times, contradictory, there is little doubt that efforts to correct misinformation are worthwhile. In fact, fighting the spread of misinformation should be viewed as vitally important health and science policy priority.

⁷⁷ Gordon Pennycook, “Fighting COVID-19 Misinformation on Social Media: Experimental Evidence for a Scalable Accuracy Nudge Intervention” (2020) [working paper], online: <<https://psyarxiv.com/uwbk9/>>.