

DECEMBER 2020

Educative Interventions to Combat Misinformation: Evidence From a Field Experiment in India

A CASI Working Paper

SUMITRA BADRINATHAN Doctoral Candidate in Comparative Politics University of Pennsylvania CASI Sobti Family Fellow, 2020-21

CENTER FOR THE ADVANCED STUDY OF INDIA

EDUCATIVE INTERVENTIONS TO COMBAT MISINFORMATION: EVIDENCE FROM A FIELD EXPERIMENT IN INDIA

Sumitra Badrinathan

Doctoral Candidate in Comparative Politics, University of Pennsylvania CASI Sobti Family Fellow, 2020-21

> A CASI Working Paper December 2020

© Copyright 2020 Sumitra Badrinathan & Center for the Advanced Study of India

ABOUT THE AUTHOR

SUMITRA BADRINATHAN is a Doctoral Candidate in Comparative Politics at the University of Pennsylvania. She is a CASI Sobti Family Fellow for the 2020-21 academic year. Sumitra received her MA in Political Science from the University of Chicago and her BA in Psychology from St. Xavier's College, Mumbai. In 2019, Sumitra received the Dean's Award for Distinguished Teaching and the Penn Prize for Excellence in Teaching by Graduate Students. Her research has received prior support from CASI, the Judith Rodin Research Fellowship, and the Russel Ackoff Fellowship.

ABSTRACT

Misinformation makes democratic governance harder, especially in developing countries. Despite its real-world import, little is known about how to combat misinformation outside of the U.S., particularly in places with low education, accelerating Internet access, and encrypted information sharing. This study uses a field experiment in India to test the efficacy of a pedagogical intervention on respondents' ability to identify misinformation during the 2019 elections (N=1224). Treated respondents received in-person media literacy training in which enumerators demonstrated tools and tips to identify misinformation in a coherent learning module. Receiving this hour-long media literacy intervention did not significantly increase respondents' ability to identify misinformation on average. However, treated respondents who support the ruling party became significantly less able to identify pro-attitudinal stories. These findings point to the resilience of misinformation in India and the presence of motivated reasoning in a traditionally non-ideological party system.¹

Keywords: Misinformation, India, Elections, Social Media, Fact-Checking, Literacy Training, WhatsApp

¹ This study was preregistered with Evidence in Governance and Policy (20190916AA) and received IRB approval from the University of Pennsylvania (832833). The author thanks Guy Grossman, Matthew Levendusky, Marc Meredith, Devesh Kapur, Neelanjan Sircar, Brendan Nyhan, Tariq Thachil, Douglas Ahler, Milan Vaishnav, Adam Ziegfeld, Devra Moehler, Jeremy Springman, Emmerich Davies, Simon Chauchard, Simone Dietrich and anonymous referees. Pranav Chaudhary and Sunai Consultancy provided excellent on-ground and implementation assistance. This research was funded by the Center for the Advanced Study of India (CASI) at the University of Pennsylvania and the Judith Rodin Fellowship. For comments and feedback, the author thanks seminar participants at the NYU Experimental Social Sciences Conference, MIT GOV/LAB Political Behavior of Development Conference, Penn Development and Research Initiative, the Harvard Experimental Political Science Conference, APSA 2020, SPSA 2020, and the UCLA Politics of Order and Development lab.

INTRODUCTION

Images of mutilated bodies and lifeless children proliferated across WhatsApp in northern India in 2018, allegedly resulting from an organized kidnapping network. In response to these messages, a young man mistaken for one of the kidnappers was mobbed and brutally beaten by villagers in Meerut, UP. The images, however, were not from a kidnapping network, but rather from a chemical weapons attack in Ghouta, Syria in 2013. Mob lynchings such as this have become a prominent problem in India since 2015, when a Muslim villager in Uttar Pradesh was killed by a mob after rumors spread that he was storing beef in his house. Such misinformation campaigns are often developed and run by political parties with nationwide cyber-armies, targeting political opponents, religious minorities and dissenting individuals (Poonam and Bansal 2019). The consequences of such rumors are as extreme as violence, demonstrating that misinformation is a matter of life and death in India and other developing countries.

What tools, if any, exist to combat the misinformation problem in developing countries? Nearly all of the extant literature on combating misinformation focuses on the U.S. and other developed democracies, where misinformation spreads via public sites such as Facebook and Twitter. Interventions in these contexts are not easily adapted for misinformation distributed on encrypted chat applications such as WhatsApp, where no one, including the app developers themselves, can see, read or analyze messages. Encryption necessitates that the burden of fact-checking fall solely on the user and therefore, the more appropriate solutions in such contexts are bottom-up, user-driven learning and fact-checking to combat misinformation.

This study is one such bottom-up effort to counter misinformation with a broad pedagogical program. I investigate whether improving information processing skills changes actual information processing in a partisan environment. The specific research question asked in this paper is whether in-person, pedagogical training to verify information is effective in combating misinformation in India. To answer this question, I implemented a large-scale field experiment with 1,224 respondents in the state of Bihar in India during the 2019 general elections, when misinformation was arguably at its peak. In an hour-long intervention, treatment group respondents were taught two concrete tools to verify information. They also received a flyer with tips to spot misinformation, along with corrections to four political fake stories. After a two-week period, respondent households were revisited to measure their ability to identify misinformation.

My experiment shows that an hour-long, educative treatment is not sufficient to help respondents combat misinformation. Importantly, the average treatment effect is not significantly distinguishable from zero. Finding that an in-person, hour-long and bottom-up learning intervention does not move people's prior attitudes is testimony to the tenacity and destructive effects of misinformation in low education settings such as India. It challenges conventional findings in American politics that subtle priming treatments, such as disputed tags, can reduce the consumption of misinformation. These findings also confirm qualitative evidence about the distinctive nature of social media consumers in developing states who are new to the Internet, lending them particularly rife and vulnerable to misinformation.

While there is no evidence of a non-zero average treatment effect, there are significant treatment effects among subgroups. Bharatiya Janata Party (BJP) partisans (those self-identifying as supporters of the BJP, the national right-wing party in India) who receive the treatment are less likely to identify pro-attitudinal stories as false. That is, on receiving counter-attitudinal corrections, the treatment backfires for BJP respondents while simultaneously working to improve information processing for non-BJP respondents. This is consistent with findings in American politics on motivated reasoning, demonstrating that respondents seek out information reinforcing prior beliefs, and that partisans cheerlead for their party and are likely to respond expressively to partisan questions (Taber and Lodge 2006; Gerber and Huber 2009; Prior, Sood, and Khanna 2015). These findings also challenge the contention that Indians lack consolidated, strong partisan identities (Chhibber and Verma 2018). I demonstrate that party identity in India is more polarized than previously thought, at least with BJP partisans and during elections.

This study hopes to spark a research agenda on the ways to create an informed citizenry in low-income democracies through testing and implementation of bottom-up measures to fight misinformation. I also seek to contribute to the empirical study of partisan identity in India, revisiting the conventional wisdom of party identities being unconsolidated and fluctuating.

WHAT IS MISINFORMATION AND HOW DO WE FIGHT IT?

I define misinformation as claims that contradict or distort common understanding of verifiable facts (Guess and Lyons 2020) and fabrications that are low in facticity (Tandoc Jr, Lim, and Ling 2018; Li 2020).

The literature on misinformation identifies three key components of false stories:

1) low levels of facticity, 2) journalistic presentation, and 3) intent to deceive (Egelhofer and Lecheler 2019; Farkas and Schou 2018). Given my focus on misinformation in India, my definition does not include the format of the news. In India, misinformation is spread via WhatsApp where much of it is in the form of text messages, with the content copied and pasted into the body of the message, such that it exists as standalone content. Hence this cannot mimic legitimate news websites and is rarely presented in a journalistic format.

Further, while the creation of falsehoods in the Indian context can stem from organized attempts by political parties with the intention to deceive, users in WhatsApp groups who are the victims of such campaigns may further propagate falsehoods inadvertently or unintentionally. Thus my definition also leaves out the intention to deceive, defined in the literature as "disinformation" (Tucker et al. 2018).

The predominant model of misinformation comes from Gentzkow, Shapiro, and Stone (2015). They posit that consumption of misinformation is a result of preferences for confirmatory stories rather than the truth because of the psychological utility from such stories. We tend to seek out information that reinforces our preferences, counterargue information that contradicts preferences, and view pro-attitudinal information as more convincing than counter-attitudinal information (Taber and Lodge 2006). Thus individuals' preexisting beliefs strongly affect their responses to corrections (Flynn, Nyhan, and Reifler 2017). Importantly, a number of contextual and individual moderators of such motivated reasoning predispose subsets of the population to be more vulnerable to misinformation.

The two key political factors that contribute to the vulnerability to misinformation effects are political sophistication and ideology (Wittenberg and Berinsky 2020). Research finds that more politically sophisticated individuals (including political knowledge and education) are more likely to be resistant to corrections (Valenzuela et al. 2019) and are the least amenable to updating beliefs when misinformation supports their existing worldviews (Lodge and Taber 2013). Thus, highly sophisticated partisans have both the motivation and the expertise to counter incongruent corrections (Wittenberg and Berinsky 2020). Further, ideology and partisanship are associated with differences in response to corrections. Although everyone is vulnerable to misinformation to a certain extent, worldview backfire effects are more visible for Republicans but not Democrats, given that the insular nature of the conservative media

system is more conducive to the spread of misinformation (Faris et al. 2017; Ecker and Ang 2019; Nyhan and Reifler 2010).

Apart from political factors, research highlights age as a key demographic variable influencing both exposure to misinformation as well as responses to it. Studies find that older adults are more likely than others to share misinformation (Grinberg et al. 2019) and that the relationship between age and vulnerability to misinformation persists even after controlling for partisanship and ideology.

But despite the growing attention to misinformation in media and scholarship, empirical literature finds that the online audience for misinformation is a small subset of the total online audience. Those consuming false stories are a small, disloyal group of heavy internet users (Nelson and Taneja 2018): Grinberg et al. (2019) find that one percent of twitter users in their sample account for 80 percent of misinformation exposures; Guess, Nyhan, and Reifler (2018) find that almost 6 in 10 visits to fake websites came from the 10 percent of people with the most conservative online information diets. However, though people online are not clamoring for a continuous stream of false stories, misinformation in a multi-faceted and fast paced online environment can command people's limited attention (Guess and Lyons 2020). Such misinformed beliefs are especially troubling when they lead people to action, as these skewed views may well alter political behavior (Hochschild and Einstein 2015).

A large research agenda has tested interventions to reduce the consumption of misinformation. These interventions can be grouped into reactive or top-down interventions that are implemented after misinformation is seen, and proactive or bottom-up interventions that seek to fight misinformation before it has been encountered.

Examples of top-down interventions include providing corrections, warnings, or factchecking and consequently measuring respondents' perceived accuracy of news stories. For instance, in 2016 Facebook began adding "disputed" tags to stories in its newsfeed that had been previously debunked by fact-checkers (Mosseri 2017); it then switched to providing fact checks underneath suspect stories (Smith, Jackson, and Raj 2017). Chan et al. (2017) find that explicit warnings can reduce the effects of misinformation; Pennycook, Cannon, and Rand (2018) test and find that disputed tags alongside veracity tags can lead to reductions in perceived accuracy; Fridkin, Kenney, and Wintersieck (2015) demonstrate that corrections from professional fact-checkers are more successful at reducing misperceptions.

Bottom-up interventions to combat misinformation rely on inoculation theory, the idea of preparing people for potential misinformation by exposing logical fallacies inherent in misleading communications a priori (Compton 2013). To this end, Tully, Vraga, and Bode (2020) and Vraga, Bode, and Tully (2020) conduct experiments where treatment group respondents were reminded to be critical consumers of the news via tweets

encouraging people to distinguish between high- and low-quality news. Roozenbeek and Van Der Linden (2019) employ games with real-world applications to combat misinformation. Cook, Lewandowsky, and Ecker (2017) inoculated respondents against misinformation by presenting mainstream scientific views alongside contrarian views. Closer in design to the present study, Guess et al. (2020) evaluate a digital literacy intervention in India and the United states utilizing the "tips" provided by WhatsApp to measure whether they are effective at increasing the perceived accuracy of true stories. In Hameleers (2020), similar tips to spot misinformation are paired with fact checks in a bundled treatment.

In sum, research finds that the most effective interventions to correct misinformation come from credible sources and sources that are surprising, such as Republicans correcting Republicans (Porter and Wood 2019). Additionally, strong social connections between fact checkers and rumor spreaders can encourage the latter to be more accepting of corrections (Margolin, Hannak, and Weber 2018), and interventions that come early before a false narrative gains traction can be more effective (Ecker et al. 2015). Finally, corrections that do not directly challenge one's worldview and identity are likely to be more effective (Flynn, Nyhan, and Reifler 2017). Drawing on such findings, metanalyses of misinformation corrections find that fact-checking has an overall positive effect on political beliefs (Walter and Murphy 2018).

Despite the large number of studies in this area, the vast majority of interventions to fight misinformation are conducted in Western contexts. Further, these studies are almost all lab and survey experiments, and hence their success has policy implications limited to populations who are frequently online and use platforms such as Facebook and Mechanical Turk. This does not describe the vast majority of populations in developing countries, who hold varying levels of digital literacy and are less likely to be avid Internet users. The next sections outline the challenge posed by misinformation in developing countries and the need for solutions and interventions specific to those contexts.

DISSEMINATION OF MISINFORMATION IN INDIA: THE SUPPLY

This study was conducted in May 2019 during the general election in India, the largest democratic exercise in the world. The 2019 contest was a reelection bid for Narendra Modi, leader of India's Hindu nationalist Bharatiya Janata Party (BJP). India is a parliamentary system but Narendra Modi's style of politics makes it akin to presidential elections with a high level of polarization, where not unlike Donald Trump, he "inspires either fervent loyalty or deep distrust" (Masih and Slater 2019).

This election was distinctive because it allowed for campaigning to be conducted over the Internet, and chat-based applications such as WhatsApp became a key communication tool for parties. For example, the BJP drew plans to have WhatsApp groups for each of India's 927,533 polling booths. A WhatsApp group can contain a maximum of 256 members, hence this communication strategy potentially reached 700 million voters. This, coupled with WhatsApp being the social media application of choice for over 90 percent of Internet users, led the BJP's social media chief to declare 2019 the year of India's first "WhatsApp elections" (Uttam 2018). Survey data from this period in India finds that one-sixth of respondents said they were members of a WhatsApp group chat started by a political leader or party (Kumar and Kumar 2018).

Unlike the United States where the focus has been on foreign-backed misinformation campaigns, political misinformation circulating in India appears to be largely domestically manufactured. The information spread on such political WhatsApp groups is not only partisan but also hate-filled and often false (Singh 2019). This trend is fueled by party workers themselves: ahead of the 2019 election, national parties hired armies of volunteers "whose job is to sit and forward messages" (Perrigo 2019). Singh (2019) reports that the BJP directed constituency-level volunteers to sort voters into groups created along religious and caste lines, even location, socioeconomic status and age, such that specific messages could be targeted to specific WhatsApp groups. So entrenched is the political misinformation machinery in India that it resembles an industry where spreading false messages is incentivized. Then BJP President Amit Shah underscored these observations during a public address in 2018: "We can keep making messages go viral, whether they are real or fake, sweet or sour" (Wire 2018). Misinformation is inherent political in India, and the creators of viral messages are often parties themselves.

VULNERABILITY TO MISINFORMATION IN INDIA: THE DEMAND

WhatsApp group chats morph into havens for misinformation in India. Four characteristics make their users vulnerable to misinformation.

First, literacy and education rates are low across the developing world. India's literacy rate, along with its rate of formal education, is relatively low compared to other developing countries where misinformation has been shown to affect public opinion (Figure 1). Further, the sample site for this study – the state of Bihar in India – has historically had one of the lowest literacy rates within the country. Research has demonstrated a strong relationship between levels of education and vulnerability to misinformation. While people with higher levels of education have more accurate beliefs (Allcott and Gentzkow 2017), motivated reasoning gives them better tools to argue against counter-attitudinal information (Nyhan et al. 2019). We should thus expect that vulnerability to misinformation is impacted by lower literacy and education.



Figure 1: India Has Low Levels of Literacy and Education

Second, Internet access has exploded in the developing world. India, particularly, is digitizing faster than most mature and emerging economies, driven by the increasing availability and decreasing cost of high-speed connectivity and smartphones, and some of the world's cheapest data plans (Kaka et al. 2019). Internet penetration in India increased exponentially over the past few years and Bihar – the sampling site for this study saw an Internet connectivity growth of over 35 percent in 2018, the highest in the country (Mathur 2019).

81 percent of users in India now own or have access to smartphones and most of these users report obtaining information and news through their phones (Devlin and Johnson 2019). Paradoxically, this leap in development coupled with the novelty and unfamiliarity with the Internet could make new users more vulnerable to information received online. The example of Geeta highlights this aspect. Geeta lives in Arrah, Bihar and recently bought a smartphone with Internet. I asked her if she thought information received over WhatsApp was factually accurate:

"This object [her Redmi phone] is only the size of my palm but is powerful enough to light up my home (...) Previously we would have to walk to the corner shop with a TV for the news. Now when this tiny device shines brightly and tells me what is happening in a city thousands of kilometers away, I feel like God is directly communicating with me" [translated from Hindi].²

² Interview with Geeta, March 27, 2019. Unless noted otherwise, all individual names are changed to protect the confidentiality of focus group participants.

Geeta's example demonstrates that the novelty of digital media could increase vulnerability to all kinds of information. Survey data shows that countries like India have several "unconscious" users who are connected to the Internet without an awareness that they are going online (Silver and Smith 2019). Such users may be unaware of what the Internet is in a variety of ways. The expansion of Internet access and smartphone availability in India thus generate the illusion of a mythic nature of social media, underscoring a belief that if something is on the Internet, it must be true.

Third, online information in developing countries is disproportionately consumed on encrypted chat-based applications such as WhatsApp. India is WhatsApp's biggest market in the world (with about 400 million users in mid-2019), but an important reason contributing to the app's popularity is also at the heart of the misinformation problem: WhatsApp messages are private and protected by encryption. This means that no one, including the app developers and owners themselves, have access to see, read, filter, and analyze text messages. This feature prevents surveillance by design, such that tracing the source or the extent of spread of a message is close to impossible, making WhatsApp akin to a black hole of misinformation. Critically, this means that top-down and platform-driven solutions are impractical in the case of private group chats on WhatsApp, suggesting that bottom-up interventions are more promising.

Finally, the format of misinformation in India is mainly visual: much of what goes viral on WhatsApp constitutes photoshopped images and manufactured videos. Misinformation in graphical and visual form is found to have increased salience, capable of retaining respondent attention to a higher degree (Flynn, Nyhan, and Reifler 2017).



Figure 2: "Cow Urine Cures Cancer" Viral WhatsApp Rumor

My intervention drew from a sampling of false photoshopped images and pseudoscientific narratives that became popular on WhatsApp in India in the months leading up to the election. Among these are false claims relating to the wondrous power of cows, along with rumors targeting minorities for storing beef or illegally slaughtering cows. Killing cows is sacrilege to many Hindus, illegal in some states, and is squarely a political and electoral issue in India (Ali 2020). According to Human Rights Watch, at least 44 people were killed in "cow-related violence" across 12 Indian states between May 2015 and December 2018. Figure 2 is an example of a false story that circulated over WhatsApp prior to the election, claiming that cow urine cures cancer. On WhatsApp, false stories are almost never shared with a link – the image above was forwarded as is to thousands of users, making the original source unknown and difficult to trace.

This image is also partisan in nature, highlighting differences between "Indian liberals," or those who do not support the right-of-center BJP, and others on the political spectrum. Evidence on the power of partisanship and ideology as polarizing social identities in India is mixed. India's party system is not historically viewed as ideologically structured. Research finds that parties as not institutionalized (Chhibber, Jensenius, and Suryanarayan 2014), elections are highly volatile (Heath 2005), and the party system itself is not ideological (Ziegfeld 2016; Kitschelt and Wilkinson 2007; Chandra 2007). More recent literature, however, argues for the idea that Indians are reasonably well sorted ideologically into parties and politics might be becoming more programmatic amongst certain groups (Chhibber and Verma 2018; Thachil 2014). Despite this, we know little about the origins of partisanship in India–whether it stems from transactional relationships with parties, affect for leaders, ties to social groups, ideological leanings–or its stability.

Despite these findings, I argue that party identities will likely moderate attitudes in India. This is largely because of the nature of the BJP's appeals. The recent BJP administration under the leadership of Prime Minister Narendra Modi represents a departure from traditional models of voting behavior in India, highlighting that Modi's rule is a form of personal politics in which voters prefer to centralize political power in a strong leader, and trust the leader to make good decisions for the polity (Sircar 2020). Some have concluded that under Modi, polarization in India is more toxic than it has been in decades, showing no signs of abating (Sahoo 2020). To add to this, misinformation is India is inherently political in nature, with disinformation campaigns often stemming from party sources themselves (Singh 2019). Finally, partisan identities tend to be more salient during elections, when citizen attachments to parties are heightened (Michelitch and Utych 2018). Taken together, these three factors indicate that BJP partisans are more likely to respond expressively to the partisan treatment and engage in motivated reasoning in the face of counter-attitudinal information.

MEDIA LITERACY INTERVENTION

I designed a pedagogical, in-person media literacy treatment with educational tools to address misinformation in the Indian context. Building on research by Guess et al. (2020), I use concrete tools to spot misinformation along with fact-checking stories and reminding respondents to be critical consumers, all in a coherent learning module.

The concept of media literacy captures the skills and competencies that promote critical engagement with messages produced by the media, needed to successfully navigate a complex information ecosystem (Jones-Jang, Mortensen, and Liu 2019). Research finds that media literacy can bolster skepticism toward false and misleading information, making it particularly suitable to address the spread of misinformation (Kahne and Bowyer 2017). Experimental studies promoting media literacy initiatives against misinformation operationalize media literacy by increasing the salience of critical thinking (Vraga, Bode, and Tully 2020) or by gauging respondent knowledge about media industries and systems (Vraga and Tully 2019) or going a step further by providing tips to spot misinformation (Guess et al. 2020).

But simply nudging respondents to be more critical consumers or providing tips asking them to be more aware may be insufficient to help counter misinformation in contexts where respondents are not armed with the tools to apply such advice to the information they encounter. In contexts such as India where respondents are unconscious Internet users unaware about misinformation, any media literacy initiative must necessarily bridge the gap between critical thinking and desired outcomes by providing concrete tools that can help foster skepticism. That is the premise of this study.

a. Experimental Design

The intervention targeted to treatment group respondents was, by design, a bundled treatment incorporating several elements, drawing on research demonstrating that the most promising tools to fight misinformation are fact-checking combined with media literacy (Hameleers 2020). The intervention consisted of surveying a respondent in their home and undertaking the following activities in a 45-60 minute visit:

- 1) Pre-treatment survey: Field enumerators administered survey modules to measure demographic and pre-treatment covariates including digital literacy, political knowledge, media trust, and prior beliefs about misinformation.³
- 2) Pedagogical intervention: Next, respondents learnt of two concrete tools to identify misinformation. Performing reverse image searches: A large part of misinformation in India comprises of misleading photos and videos, often drawn from one context and used to spread misinformation about another context or

³ Summary statistics for all key variables are included in Table A.1

time. Reverse searching such images is an easy way identify their origins. As one focus group discussion conducted before the experiment revealed: "the time stamp on the photo helped me realize that it is not current news; if this image has existed since 2009, it cannot be about the 2019 election."⁴ Respondents can see the original source and time stamp on an image once it is fed back into Google, making this technique a uniquely useful and compelling tool given the nature of visual misinformation in India. Enumerators demonstrated two examples of this to respondents.

Navigating a fact-checking site: Focus group discussions also revealed that while a minority of those surveyed knew about the existence of fact-checking websites in India, even fewer were able to name one. The second concrete tool involved demonstrating to respondents how to navigate a fact-checking website, www.altnews.in⁵, a non-profit fact-checking service in India. Enumerators explained the layout of the site, showed respondents where to find fact-checked viral false stories, etc.

- 3) Corrections and tips flyer: Enumerators next helped respondents apply these tools to fact-check four false stories. Do to so enumerators displayed a flyer to respondents, the front side of which had descriptions of four recent viral political false stories. For each story, enumerators systematically corrected the false story, explaining in each case why the story was untrue, what the correct version was, and what tools were used to determine veracity. The back side of the flyer contained six tips to reduce the spread of misinformation. The enumerator read and explained each tip to respondents, gave them a copy of the flyer and exhorted them to make use of it. These tools were demonstrated to treatment group respondents only. Control group respondents were shown a placebo demonstration about plastic pollution, and were given a flyer containing tips to reduce plastic usage.
- 4) Comprehension Check: Enumerators lastly administered a comprehension check to measure whether the treatment was effective in the short-term.

For this study, respondents were randomized into one of three groups, two treatment and one placebo control. Table 1 summarizes the three groups.

⁴ Interview with Bharat, March 31, 2019.

⁵ <u>https://www.altnews.in/hindi/</u>

Table 1: Experimental	l Treatments
-----------------------	--------------

Intervention	Goal
T1: Pedagogical Intervention + Pro-BJP flyer	Tools + corrections to 4 pro-BJP false stories
T2: Pedagogical Intervention + Anti-BJP flyer	Tools + corrections to 4 anti-BJP false stories
Control: Plastic Pollution Intervention + flyer	Tools + tips on plastic pollution

Respondents in both treatment groups received the pedagogical intervention. However, one group received corrections to four pro-BJP false stories, the other received corrections to four anti-BJP false stories. Besides differences in the stories that were fact-checked, the tips on the flyer remained the same for both treatment groups. Respondents in the placebo control group received a symmetric treatment where enumerators spoke about plastic pollution and were given a flyer on tips to reduce plastic usage. The false stories included in the treatment group flyers were drawn from a pool of stories fact checked for accuracy by altnews.in and boomlive.in. The partisan slant of each story was determined by a Mechanical Turk pre-test. To ensure balance across both treatment groups, stories with similar salience and subject matter were picked. The back of treatment flyers contained the same tips on how to verify information and spot false stories. The entire intervention was administered in Hindi. Figures C.1, C.2 and C.3 present the English-translated version of flyers distributed to respondents.

To control for potential imbalance in the sample, a randomized block design was used. Those respondents who identified with the BJP were one block, those who identified with any other party were another block. Within each block, respondents were randomly assigned to one of the three experimental groups described in Table 1. This design ensured that each treatment condition had an equal proportion of BJP and non-BJP partisans. Overall, the sample was equally divided between the two treatment and placebo control groups (i.e. one third of the sample in each of the three groups).

b. Sample and Timeline

The sample was drawn from the city of Gaya in the state of Bihar in India. Bihar has both the lowest literacy rate in the country as well as the highest rural penetration of mobile phones, making it a strong test-case for the intervention.

Respondents were selected through a random walk procedure. Within the sampling area, a random sample of polling booths (smallest administrative units) were selected to serve as enumeration areas. Within each enumeration area, enumerators were instructed to survey 10-12 households following a random walk procedure. This method was chosen over traditional survey listing techniques so as to minimize enumerator time

spent in the field during the elections as well as because of a lack of accurate census data for listing (Lupu and Michelitch 2018). Each field enumerator was assigned to only one polling booth, and hence the paths taken during each selection crossed each household only once, increasing the likelihood of a random and unbiased sample.

Once a household was selected, household members could qualify for the study based on three pre-conditions designed to maximize familiarity with the Internet: respondents were required to have their own cellphone (i.e. not a shared household phone), working Internet for 6 months prior to the survey, and WhatsApp was required to be downloaded on the phone. If multiple members of a household qualified based on the pre-conditions, a randomly selected adult member was requested to participate in the study.

Of note, only 20 percent of all households sampled had respondents who met the criteria for recruitment into the study. In Bihar, where only 20-30 percent of citizens have access to the Internet, this is unsurprising. Despite this, the study had a high response rate: of all those who were eligible for the study, 94.5 percent agreed to participate. The final sample comprised of 1,224 respondents.⁶

Trained enumerators administered the intervention in a household visit rolled out in May 2019. Approximately two weeks after the intervention, the same respondents were revisited to conduct an endline survey and measure the outcomes of interest. Critically, respondents voted in the election between the two enumerator visits. Figure 3 summarizes the timeline for this study.



Figure 3: Experimental Timeline (May 2019)

The study took multiple steps in survey design and implementation to minimize exogenous shocks from election results. The timeline ensured that though respondents voted in the general election after the intervention, making voter turnout posttreatment, the endline survey to measure outcomes was conducted before election votes

⁶ Additional details about the sampling process are available in Online Appendix B.

were counted and results were announced.⁷ This timeline had the double advantage of ensuring that outcome measures was not impacted by the exogenous shock of results while also making sure respondents received the intervention before they voted, when political misinformation is arguably at its peak. At the end of the baseline survey, enumerators collected addresses and mobile numbers of respondents for subsequent rounds of the study and then immediately separated this contact information from the main body of the survey to maintain respondent privacy.

c. Dependent Variables

In the endline survey, enumerators revisited respondents after they had voted. The same set of enumerators administered the intervention and the endline survey. However, enumerators were given a random set of household addresses for the endline survey to minimize the possibility of the same enumerator systematically interviewing the same respondent twice. Further, addresses and contact information were separated immediately from baseline survey data to ensure that enumerators only had contact information about respondents. During the baseline survey, 1306 respondents were administered the intervention. The enumerators successfully located 1224 of these respondents, resulting in an attrition rate of 6 percent. Importantly, nobody who was administered the intervention refused to answer the endline survey; the attrited group comprised only of respondents who enumerators were unable to contact at home after three tries.

The key outcome of interest is whether the intervention positively affected respondents' ability to identify misinformation. To this end, respondents were shown a series of fourteen news stories.⁸ These stories varied in content, salience, and critically, partisan slant. Half of the stories were pro-BJP in nature and the other half anti-BJP.⁹ Each respondent saw all the fourteen stories, but the order in which they were shown was randomized.¹⁰

Following each story, two primary dependent variables were measured:

⁷ In India voting is staggered by constituency but ballots are counted after every constituency in the country has voted.

⁸ 12 were false and 2 were true. Given the countless, diverse array of stories that went viral in India during this time with perilous consequences, I chose to maximize on reducing belief in as many false stories as possible. Hence respondents were shown more false stories as part of the outcome measure (rather than a 50-50 split between true and false stories). Two true stories (each of different partisan slant) were included in the measure, and respondents were told that some of the stories were false and some true. More analysis of the true stories is in Online Appendix I.

⁹ Partisan slant of the news stories was determined with a Mechanical Turk pre-test.

¹⁰ For field safety reasons, the endline survey was conducted offline and hence the order of appearance of the dependent variable stories was limited to 3 pre-determined random orders. A given enumerator had access to only one of the 3 random orders. As a robustness check, I replicate the main analysis with enumerator fixed effects. Results are presented in Tables E.1 and E.2.

- 1) Perceived accuracy of news stories, with the question "Do you believe this news story is false?" (binary response, 1 if yes, 0 otherwise)
- 2) Confidence in identification of the story as false or real, with the question "How confident are you that the story is real / false?" (4-point scale, 1 = very confident, 4 = not confident at all)

A list of the fourteen stories shown to respondents is presented in Table D.1.11

d. Hypotheses and Estimation

I hypothesize there will be a positive effect of the intervention for respondents assigned to any arm of the treatment group relative to placebo control. I also hypothesize that the individual effect of being assigned to each treatment will be positive relative to placebo control:

Hypothesis 1: Exposure to the media literacy intervention will increase ability to identify misinformation relative to control.

Hypothesis 2a: Exposure to media literacy and pro-BJP corrections will increase ability to identify misinformation.

Hypothesis 2b: Exposure to media literacy and anti-BJP corrections will increase ability to identify misinformation.

I estimate the following equations to test the main effect of the intervention:

$$Misinformation Id_i = \alpha + \beta_1 Intervention_i + \epsilon_i$$
(5.1)

 $MisinformationId_i = \alpha + \beta_1 InterventionPro-BJP_i + \beta_2 InterventionAnti-BJP_i + \epsilon_i$ (5.2)

In the equations, i represents the respondent, the Intervention variable in Equation 5.1 represents pooled assignment to the media literacy intervention (relative to control). In Equation 5.2, the dependent variable is regressed on separate indicators for having received the intervention and pro-BJP corrections, or intervention and anti-BJP corrections, with the control condition as the omitted category. The dependent variable MisinformationId counts the number of stories correctly identified as fake.

MisinformationId has been coded such that a positive estimated b_1 indicates an increase in the ability to identify misinformation.

¹¹ Online Appendix D describes secondary dependent variables measured.

Beyond the average treatment effect, I expect treatment effects to differ conditional on a single factor previously identified in the literature as a significant predictor of information consumption: partisan identity. In line with the literature on partisan motivated reasoning (Nyhan and Reifler 2010), I expect that the treatment effect will be larger for politically incongruent information as compared to politically congruent information, relative to the control condition. A politically congruent condition manifests when corrections are pro-attitudinal, i.e., BJP partisans receiving corrections to pro-BJP false stories.

Hypothesis 3: Effectiveness of the intervention will be higher for politically incongruent information compared to politically congruent information, relative to the control condition.

To determine whether partisan identity moderates treatment effects, I test Hypothesis 3 with the following model:

$$MisinformationId_{i} = \alpha + \beta_{1}Intervention_{i} + \beta_{2}Intervention_{i} * PartyID_{i} + \beta_{3}PartyID_{i} + \epsilon_{i}$$
(5.3)

In Equation 5.3, PartyID is an indicator variable that takes on the value of 1 if the respondent self-identified as a BJP supporter. The choice to code party identity as dichotomous was based on the nature of misinformation in India where false stories are perceived as either favoring or not favoring the BJP. A positive coefficient estimate for b_2 indicates an increase in the ability to identify misinformation among BJP partisans due to the treatment.

However, while partisanship might moderate attitudes, the role of other theoretical moderators such as political sophistication and age is unclear in the Indian context. The context of this experiment is one of low literacy and education, but there is little reason to expect that education or literacy correlate with political knowledge. Indeed, voter turnout rates among low-income groups in India are as high as richer segments of the population, indicating knowledge of and interest in politics despite lower levels of education (Ahuja and Chhibber 2012). Similarly, owing to the lack of priors about effects of age or digital literacy on attitudes, theoretical expectations regarding these variables remain ambiguous.

Thus, while I do not form precise pre-registered hypotheses about these moderators, I examine their relationships with misinformation through the following research questions:

RQ1: What is the relationship between age and vulnerability to misinformation? Does this relationship change as a function of the treatment?

RQ2: Does age correlate negatively with digital literacy, as in the American context? Are more digitally literate respondents likely to learn better from the treatment?

RQ3: What is the relationship between political sophistication (measured both by education and political knowledge) and vulnerability to misinformation?

DATA AND RESULTS

This section begins with descriptive analyses that demonstrate the extent of belief in misinformation as well as partian polarization in this belief.



Figure 4: Percent of Sample Who Believe Rumors

Figure 4 lists the 12 false stories used in the dependent variable measure in this study. This figure plots the share of respondents in the sample who believed each story to be true. Two aspects of the figure are striking. First, general belief in misinformation is low. For half of the 12 false stories, less than 10 percent of the sample thought they were true.

Second, belief in pro-BJP misinformation appears to be stronger, possibly alluding to its increased salience (Jerit and Barabas 2012), frequency of appearance on social media (Sinha, Sheikh, and Sidharth 2019), or to the presence of a higher proportion of BJP supporters in the sample. Overall, across the 12 rumors, respondents correctly classified an average of 9.91 rumors.



Belief in Anti-BJP Rumors by Party

Figure 5: Belief in Rumors by Respondents' Party ID

Soldier

Flag

2000 Note

Pulwama

Attacks

0

Gomutra

19

© Copyright 2020 Sumitra Badrinathan & Center for the Advanced Study of India

Figure 5 plots respondent belief in stories by partisan identity. For 10 out of the 12 partisan stories, we see a correspondence between respondent party identity and pretested political slant of the story. Though there is partisan sorting on belief in political rumors, the gap between BJP and non-BJP partisans in their beliefs is not as large as in the American case: the biggest gap appears in the case of the Unclean Ganga river rumor, where non-BJP partisans showed about 9 percentage points more belief in the rumor relative to BJP supporters. In contrast, Jardina and Traugott (2019) demonstrate that differences between Democrats and Republicans in their belief of the Obama birther rumor can be as large as 80 percentage points.

To identify differences between sub-populations in vulnerability to misinformation, I analyze the correlates of misinformation among control group respondents (N=406). This analysis provides the baseline rate of identification ability in the absence of the intervention. In the regression analysis in Table 2, the dependent variable is the number of stories accurately identified by control group respondents.

First, we observe that BJP supporters were significantly better at identifying false stories. This observational result is striking – on the one hand, pro-BJP rumors are more likely to be believed by respondents, in line with descriptions of a right-wing advantage in producing misinformation (supply side). However, demand side results demonstrate that BJP supporters are better at identifying misinformation.

This finding bodes with observations that incentives to spread partisan misinformation has led parties like the BJP to form "cyber-armies" to disseminate information. Thus, while it is possible that BJP respondents are more aware of party-driven supply of misinformation, thereby being able to identify rumors at greater rates, their partisanship also makes them expressively believe pro-attitudinal rumors. These observational findings suggest the presence of partisan motivated reasoning in the Indian context.

BJP Supporter	0.535***
	(0.196)
Digital Literacy	-1.007**
(Higher = more literate)	(0.423)
Political Knowledge	-0.059
(Higher = more knowledge)	(0.079)
Frequency of WhatsApp Use	0.164*
(Higher = more usage)	(0.085)
Trust in WhatsApp	0.068
(Higher = more trust)	(0.102)
Education	0.079**
	(0.031)
Age	0.033***
0	(0.009)
Male	0.234
	(0.280)
Hindu	-0.141
	(0.253)
Constant	7.724***
	(0.735)
Observations	406
R^2	0.093
Adjusted R ²	0.072
Residual Std. Error	1.545 (df = 396)
F Statistic	4.512^{***} (df = 9; 396)
Note:	*p<0.1; **p<0.05; ***p<0.01

Table 2: Misinformation Identification in Control Group

Dependent variable: Number of Stories Identified as False

Next, we observe interesting findings with respect to age and digital literacy. While findings on misinformation in the United States suggest that older adults are most likely to engage with fake sources (Grinberg et al. 2019), this data demonstrates the opposite:

younger adults are less likely to identify false stories. Similarly, increases in digital literacy are associated with lower identification of misinformation, contrary to findings in the American context that people who are less digitally literate are more likely to fall for misinformation and clickbait (Munger et al. 2018). Finally, I find that while political knowledge is not correlated with misinformation identification, education is associated with significant increases in the identification of false stories.¹²

I now move to discussing experimental results. Enumerators administered a comprehension check at the end of the intervention to measure whether the treatment was effective in the short-term. Respondents were shown two false stories that were debunked by enumerators in the same house visit (as a part of the flyer with corrections). For each story, immediately after the treatment, respondents were asked to identify whether it was false or not. Less than 5 percent of the sample for both stories incorrectly identified them as true, demonstrating that in the short run, respondents were able to successfully identify stories as false after they had been debunked.

I estimate effects of the treatment on outcomes in a between-subjects design. All estimates are ordinary least square (OLS) regressions and empirical models are specified relying on random treatment assignment to control for potential confounders. First, I analyze data for the main effect of the intervention. While research predicts that in-person and field interventions on media effects are likely to have stronger effects (Jerit, Barabas, and Clifford 2013; Flynn, Nyhan, and Reifler 2017), my findings from misinformation-prone India are less encouraging. Even with an in-person intervention, where enumerators spend close to one hour with each respondents to debunk and discuss misinformation and where respondents understood the intervention, I do not see significant increases in the ability to identify misinformation as a function of teaching respondents media literacy tools.

Results are shown in Table 3. The key dependent variable in my analysis counts the number of stories that a respondent correctly identified as false.¹³ Columns 1 and 3 include stories that were classified from the pre-test as having a pro-BJP slant, Columns 2 and 4 include stories that were classified as having an anti-BJP slant. To estimate the pooled effect of the intervention, I construct a variable that takes on the value of 1 if a respondent received any literacy and fact-checking treatment (relative to 0 if the respondent was in the placebo control group). This effect of this pooled treatment is estimated in models (1) and (2). In models (3) and (4), I split the treatment into the pro-BJP corrections and the anti-BJP corrections (note both treatment conditions receive the same literacy intervention).

¹² I explore these results further in Online Appendix H.

¹³ The dependent variable in these models counts the number of stories identified as false out of a total of 12 false stories. I replicate these analyses where the dependent variable is the share of correctly identified stories given all fourteen stories, true and false, and find that the results hold. Analyses are in Tables F.1 and F.2.

	Depend	Dependent variable: Number of Stories Identified as False				
	Pro-BJP Stories	Anti-BJP Stories	Pro-BJP Stories	Anti-BJP Stories		
	(1)	(2)	(3)	(4)		
Literacy Intervention	-0.004 (0.067)	0.004 (0.056)				
Literacy + Pro-BJP Fact-Check			0.013 (0.078)	0.013 (0.065)		
Literacy + Anti-BJP Fact-Check			-0.021 (0.078)	-0.006 (0.065)		
Constant	4.569*** (0.055)	5.342*** (0.046)	4.569*** (0.055)	5.342*** (0.046)		
Observations R ² Adjusted R ²	1,224 0.00000 -0.001	1,224 0.00000 -0.001	1,224 0.0002 -0.001	1,224 0.0001 -0.002		
Res. Std. Error	1.110 (df = 1222)	0.925 (df = 1222)	1.111 (df = 1221)	0.926 (df = 1221)		
Note:			*p<0.1; **	p<0.05; ***p<0.01		

Table 3: Effect of Treatment on Ability to Identify Misinformation

Table 3 demonstrates that the intervention did not increase misinformation identification ability on average. Splitting the treatment into its component parts (each compared to placebo control) yields similar results. I find no evidence that an hour-long pedagogical intervention increased ability to identify misinformation among respondents in Bihar, India. The ability to update one's priors in response to factual information is privately and socially valuable, and hence the fact that a strong, in-person treatment does not change opinions demonstrates the resilience of misinformation in India. Priors about misinformation in this context appear resistant to change but, as I demonstrate below, this does not preclude moderating effects of partisan identity.

I now turn to the analysis of heterogeneous effects of partisan identity. Table 4 presents results. In Column 1 I estimate the effect of receiving the treatment for BJP supporters on ability to identify pro-BJP false stories, Column 2 does the same with anti-BJP false stories. The treatment variable for both models pools across receiving any treatment relative to control.

	Dependent variable: Number of Stories Identified as False		
	Pro-BJP Stories	Anti-BJP Stories	
	(1)	(2)	
Literacy Intervention	0.277**	0.091	
	(0.119)	(0.099)	
BJP Supporter	0.226*	0.311***	
	(0.118)	(0.098)	
Literacy Intervention x	-0.412^{***}	-0.130	
BJP Supporter	(0.144)	(0.120)	
Constant	4.415***	5.131***	
	(0.097)	(0.081)	
Observations	1,224	1,224	
R ²	0.007	0.014	
Adjusted R ²	0.005	0.011	
Residual Std. Error ($df = 1220$)	1.107	0.920	
F Statistic (df = 3; 1220)	2.892**	5.651***	

Table 4: Effect of Treatment x Party on Ability to Identify Misinformation

Note:

*p<0.1; **p<0.05; ***p<0.01



Figure 6: Predicted Identification of Pro-BJP Stories

24 © Copyright 2020 Sumitra Badrinathan & Center for the Advanced Study of India Results are striking: while there was no average treatment affect, the interaction effect of the treatment on BJP partisans produces a negative effect on the ability to identify misinformation. For pro-BJP stories, the treatment effect for non-BJP supporters was 0.277, indicating that those who did not support the BJP and received the treatment identified an additional 0.277 stories. However, the treatment effect for BJP supporters was -0.135, indicating that those who supported the BJP and received the treatment identified 0.135 fewer stories.

Visualizing this interaction effect in Figure 6, where I graph the predicted values from the interaction model in Equation 5.3, it appears that the treatment had contradictory effects conditional on party identity (for the set of pro-BJP stories). The intercept for BJP partisans is higher, demonstrating better identification skills ex-ante, in the absence of the treatment. However, treatment group respondents who identify as BJP partisans show a significant decrease in their ability to identify false stories, while treatment group respondents who do not identify as BJP partisans show an increase in their ability to identify false stories. Thus the treatment was successful with non-BJP partisans, and backfired for BJP partisans. Importantly, these effects obtain only for the set of false stories that is pro-BJP in slant (implying that their corrections could be perceived as pro-attitudinal for non-BJP partisans). In Figure 7 I graph the interaction for the set of dependent variable stories that are anti-BJP in slant. While the relationships in this graph are directionally similar, they are smaller in magnitude and not significant. Importantly, fact-checking is much more effective for anti-BJP stories than for pro-BJP stories (note that the effects are much larger). Pro-BJP stories are more likely to be identified as false in the control, but the treatment is weaker for this subset of stories. Taken together, these results imply that non-BJP respondents were able to successfully apply the treatment to identify pro-attitudinal corrections. But for BJP partisans, given that these corrections are not consistent with their partisan identity, the treatment backfires.



Figure 7: Predicted Identification of Anti-BJP Stories

25

© Copyright 2020 Sumitra Badrinathan & Center for the Advanced Study of India

Despite the negative relationship between digital literacy and successful identification in the observational data, heterogeneous effects of the treatment demonstrate that certain sub-populations in the sample could successfully learn from the intervention and improve information processing, suggesting that cognitive detection of real from false news operates orthogonally to digital literacy. However, finding that higher levels of identification (in the control group for BJP respondents) were made worse as a function of the treatment demonstrates the existence of partisan motivated reasoning in the Indian context. I examine this result further in the Discussion.

Moving beyond experimental results, I find that younger adults in the sample are less likely to be able to identify misinformation and that higher levels of digital literacy are associated with greater vulnerability to misinformation, contrary to findings in the United States (Munger et al. 2018; Grinberg et al. 2019). I also find that while political knowledge does not correlate with perceptions of stories, more educated respondents are better at spotting false stories. I explore these associations in Online Appendices G and H.

DISCUSSION

The most striking finding to emerge from this study demonstrates that the intervention improved misinformation identification skills for one set of respondents (non-BJP respondents) but not another (BJP partisans). Paralleling results seen in developed contexts, the perceptual screen (Campbell et al. 1960) of BJP partisanship shaped how respondents interacted with this treatment, with BJP partisans demonstrating a tendency to cheerlead for their party and discredit pro-party stories despite them being false (Gerber and Huber 2009; Prior, Sood, and Khanna 2015). At the same time, non-BJP partisans who learnt from the treatment might also have identified pro-BJP stories as false because this is the response congruent with their identity. These findings of motivated reasoning demonstrate that citizen attachments to political parties are heightened during elections (Michelitch and Utych 2018) and that strong partisans engage in strategic ignorance, pushing away information and facts that get in the way of feelings (McGoey 2012).

This finding is also surprising, given that there is little evidence of such backfire effects in the American context (Wood and Porter 2019). However, several other associations in the American context do not hold in this data: I find a positive correlation between increasing age and vulnerability to misinformation, a negative correlation between increasing digital literacy and vulnerability to misinformation, and no association with political knowledge.¹⁴ The nature of these findings underscores that what we know about misinformation comes largely from Western contexts and may not easily apply to

¹⁴ See Online Appendix G for results.

other settings. It highlights that we need more theorizing and more data from non-Western contexts. Thus while I do find some backfire effects in this data, more needs to be done to establish the robustness of these findings. Future work should examine treatments such as this one in non-electoral contexts where the salience of partisanship may be lower, resulting in smaller differences between parties. Nevertheless, my findings suggest that even in democracies with weaker partisan identification, citizens still engage in motivated reasoning. This has important implications beyond the study of fact-checking and extends more broadly to how Indian citizens make political judgements.

However, it is important to underscore that the intervention worsened misinformation identification only for the pro-BJP set of false stories. This effect does not appear for anti-BJP false stories. This highlights key differences in partisan identities in this data. First, though traditionally India has been described as a non-ideological system, the recent years under the Modi-led BJP governments have led some to conclude that tribalism and psychological attachments to political parties (Westwood et al. 2018) are more salient now than ever before (Sircar 2020). Importantly, such partisan attachments seem to have arisen in response to the personal popularity of Narendra Modi, with no comparable cult of personality on the political left. Thus it stands to reason that partisanship is stronger for BJP supporters. Second, the way the party identity variable is operationalized in my data further emphasizes this point. I cluster BJP supporters into one block and non-BJP supporters into another, but the non-BJP block is a heterogeneous group of respondents from several different parties. Thus we should expect that citizen attachments to political stories, true and false, will be perceived very differently for both political blocks. Third, political disinformation campaigns in India seem to emanate largely from the right-wing. This is underscored in my data by pro-BJP stories being believed to a much greater extent than anti-BJP stories, alluding to the fact that pro-BJP stories are more salient in the minds of respondents (Figure 4).

As a consequence of these factors, there is an inherent lack of symmetry between the two sets of stories that comprise my dependent variable measure. Pro-BJP stories are more salient and believed to a greater extent, hence there is likely more room for the treatment to move attitudes on the stories (as it does, for non-BJP supporters). On the contrary, the majority of anti-BJP stories were believed by less than 10 percent of the sample; this high ceiling might make it difficult for the treatment to work for anti-BJP stories.

CONCLUSION

Misinformation campaigns have the capacity to affect opinions and elections across the world. Purveyors and victims of misinformation and hyper-partisan messaging are no

more just individuals with low digital literacy skills, people who are uninformed, Internet scammers, or Russian trolls. A global rise in polarization has meant that the creators and contributors of misinformation include party workers, stakeholders and politicians themselves. As the world moves to deal with the COVID-19 crisis, we are engulfed in a new deluge of misinformation in hyper-partisan and polarized environments, where traditionally non-political issues are also deeply politicized. The rise of polarization amidst a global pandemic underscores the need to identify robust strategies to counter the pernicious effects of misinformation, especially in societies where it is spread on encrypted platforms and where the stakes are as high as violence.

In this paper, I present new evidence on belief in popular misinformation stories in India in the context of the 2019 general elections. I design a pedagogical intervention to foster bottom-up skills training to identify misinformation. Using tools specifically designed for the Indian context, I administer in-person skills training to 1224 respondents in Bihar, India in a field experiment. I find that this grassroots-level pedagogical intervention in has little effect on respondent ability to identify misinformation on average. But, the partisanship and polarization of BJP supporters appears stickier than that of their out-partisans. Non-BJP supporters in the sample receive the treatment and apply it to identify misinformation at a higher level, demonstrating that cognitive skills can be improved as a function of the treatment. However, for BJP partisans, receiving the treatment leads to a significant decrease in identification ability, but only for pro-attitudinal stories.

The presence of motivated reasoning is a surprising result in a country with traditionally weak party ties and non-ideological party systems. Democratic citizens have a stake in dispelling rumors and falsehoods, but in societies with polarized social groups, individuals also have a stake in maintaining their personal standing in social groups that matter to them (Kahan et al. 2017). The finding that the intervention worked on a subset of respondents underscores the fact that the training was not strong enough to overcome the effects of group identity for BJP respondents. Theoretically, this result is similar to research that finds that identity protective cognition, a type of motivated reasoning, increases pressure to form group-congruent beliefs and steers individuals away from beliefs that could alienate them from others they are similar to (Sherman and Cohen 2006; Giner-Sorolla and Chaiken 1997). Practically, the result calls for a revision of findings on party identity in India, as it demonstrates the presence of motivated reasoning in electoral settings.

The effects of party identity in this setting are arguably observable because of the rise of hyper-partisan and polarizing parties in India and across the world. It underscores a broader phenomenon of populist parties and narratives, resulting in societies where information is weaponized to divide polarized voters. While elections are times when political discourse is polarized and partisanship salience is heightened, these findings stress the need for more systematic research into motivated reasoning and polarization in societies that have traditionally been non-ideological and where encrypted forms of social media take center stage in the spread of misinformation. Thus, future research should test the effect of long-term learning and skills training to counter misinformation.

Despite these findings on partisanship, I consider some reasons why the average treatment effect was a null, along with some limitations of the study and future avenues for research.

First, it is worth noting that this was an explicitly political intervention. Consistent with recent work (Groenendyk and Krupnikov 2020), the political nature of the treatment itself likely activated motivated reasoning. Next, the timing of the intervention during a contentious election meant that not only were partisan identities more salient (Michelitch and Utych 2018) but also that the presence of several election officials, campaigning party workers, and GOTV efforts meant that respondents in the area had their door knocked on several times a day by different interest groups. Thus it is possible that the marginal effect of an additional house visit by the enumeration team for this study made the in-person intervention less salient. Further, the two-week gap between the intervention and the measurement of outcomes is atypical for studies of this kind, where dependent variables are measured in close proximity to treatments. Thus it's possible that a first-stage effect decayed over time and hence was not captured by the study. Additionally, the design over-sampled false news stories in the outcome measure. While this was done to maximize belief reduction in as many false stories with perilous consequences, future studies can systematically vary the balance of true and false stories to study how this factor shapes the efficacy of these types of campaigns.

In addition, while the study measured the perceived accuracy of news stories, it did not measure whether the participants used fact-checking tools between the intervention and the follow up. An interesting prospect for future work would be to validate the usage and frequency of procedural tools before measuring beliefs. In sum, the results of such a treatment might be different with neutral, apolitical treatments conducted in less partisan times.

The findings from this study are local average treatment effects, dependent heavily on the locality where this experiment was conducted: the low-education, low-internet environment of semi-urban Bihar, at a time where politics was salient and where political misinformation was rife. Whether these findings would hold—or change outside of this locality remains an open empirical question. Consequently, I caution about interpreting these null results to mean that interventions of this kind do not work, as thorough future work must look into replicating such a design in different contexts and times. Thus, while this study was necessarily context-dependent, it is nevertheless an important first step towards tempering the human cost of misinformation in India.

REFERENCES

- Ahuja, Amit, and Pradeep Chhibber. 2012. "Why the poor vote in India: If I don't vote, I am dead to the state." Studies in Comparative International Development 47 (4): 389-410.
- Ali, Mohammed. 2020. "The Rise of a Hindu Vigilante in the Age of WhatsApp and Modi." Wired. April 14, 2020.

https://www.wired.com/story/indias-frightening-descent-social-media-terror/.

- Allcott, Hunt, and Matthew Gentzkow. 2017. "Social Media and Fake News in the 2016 Election." Journal of Economic Perspectives 31 (2): 211–36.
- Campbell, Angus, Philip E Converse, Warren E Miller, and Donald E Stokes. 1960. The American Voter. New York: John Wiley.
- Chan, Man-pui Sally, Christopher R Jones, Kathleen Hall Jamieson, and Dolores Albarrac´ın. 2017. "Debunking: A Meta-Analysis of the Psychological Efficacy of Messages Countering Misinformation." Psychological Science 28 (11): 1531– 1546.
- Chandra, Kanchan. 2007. Why Ethnic Parties Succeed: Patronage and Ethnic Headcounts in India. New York: Cambridge University Press.
- Chhibber, Pradeep, Francesca Refsum Jensenius, and Pavithra Suryanarayan. 2014. "Party organization and party proliferation in India." Party Politics 20 (4): 489–505.
- Chhibber, Pradeep, and Rahul Verma. 2018. Ideology and Identity: The Changing Party Systems of India. New York: Oxford University Press.
- Compton, Josh. 2013. "Inoculation Theory." The Sage handbook of persuasion: Developments in theory and practice 2: 220–237.

- Cook, John, Stephan Lewandowsky, and Ullrich KH Ecker. 2017. "Neutralizing misinformation through inoculation: Exposing misleading argumentation techniques reduces their influence." PloS One 12 (5): e0175799.
- Devlin, Kat, and Courtney Johnson. 2019. "Indian elections nearing amid frustration with politics, concerns about misinformation." Pew Research Center. March 25, 2019. https://www.pewresearch.org/fact-tank/2019/03/25/indian-elections-nearing-amid-frustration-with-politics-concerns-about-misinformation/.
- Ecker, Ullrich KH, and Li Chang Ang. 2019. "Political Attitudes and the Processing of Misinformation Corrections." Political Psychology 40 (2): 241–260.
- Ecker, Ullrich KH, Stephan Lewandowsky, Candy SC Cheung, and Murray T Maybery. 2015. "He did it! She did it! No, she did not! Multiple causal explanations and the continued influence of misinformation." Journal of Memory and Language 85: 101– 115.
- Egelhofer, Jana Laura, and Sophie Lecheler. 2019. "Fake news as a twodimensional phenomenon: a framework and research agenda." Annals of the International Communication Association 43 (2): 97–116.
- Faris, Robert, Hal Roberts, Bruce Etling, Nikki Bourassa, Ethan Zuckerman, and Yochai Benkler. 2017. "Partisanship, Propaganda, and Disinformation: Online media and the 2016 US Presidential Election." Berkman Klein Center Research Publication 6. https://bit.ly/3lmlcOl.
- Farkas, Johan, and Jannick Schou. 2018. "Fake News as a Floating Signifier: Hegemony, Antagonism and the Politics of Falsehood." Javnost-The Public 25 (3): 298–314.

- Flynn, D.J., Brendan Nyhan, and Jason Reifler. 2017. "The Nature and Origins of Misperceptions: Understanding False and Unsupported Beliefs About Politics."
 Political Psychology 38: 127–150.
- Fridkin, Kim, Patrick J Kenney, and Amanda Wintersieck. 2015. "Liar, Liar, Pants on Fire: How Fact-Checking Influences Citizens' Reactions to Negative Advertising." Political Communication 32 (1): 127–151.
- Gentzkow, Matthew, Jesse M Shapiro, and Daniel F Stone. 2015. "Media Bias in the Marketplace: Theory." In Handbook of Media Economics. Vol. 1. Elsevier.
- Gerber, Alan S, and Gregory A Huber. 2009. "Partisanship and Economic Behavior: Do Partisan Differences in Economic Forecasts Predict Real Economic Behavior?" American Political Science Review 103 (3): 407–426.
- Giner-Sorolla, Roger, and Sheily Chaiken. 1997. "Selective Use of Heuristic and Systematic Processing Under Fefense Motivation." Personality and Social Psychology Bulletin 23 (1): 84–97.
- Grinberg, Nir, Kenneth Joseph, Lisa Friedland, Briony Swire-Thompson, and David Lazer. 2019. "Fake news on Twitter during the 2016 U.S. presidential election." Science 363 (6425): 374–378.
- Groenendyk, Eric, and Yanna Krupnikov. 2020. "What Motivates Reasoning? A Theory of Goal-Dependent Political Evaluation." American Journal of Political Science: 1–17.
- Guess, Andrew, Brendan Nyhan, and Jason Reifler. 2018. "Selective exposure to misinformation: Evidence from the consumption of fake news during the 2016 US presidential campaign." Working Paper.

http://www.ask-force.org/web/Fundamentalists/Guess-Selective-Exposure-to-Misinformation-Evidence-Presidential-Campaign-2018.pdf.

- Guess, Andrew M, and Benjamin A Lyons. 2020. "Misinformation, Disinformation, and Online Propaganda." In Social Media and Democracy: The State of the Field, Prospects for Reform, ed. Nathaniel Persily and Joshua A Tucker. Cambridge University Press.
- Guess, Andrew M, Michael Lerner, Benjamin Lyons, Jacob M Montgomery, Brendan Nyhan, Jason Reifler, and Neelanjan Sircar. 2020. "A digital media literacy intervention increases discernment between mainstream and false news in the United States and India." Proceedings of the National Academy of Sciences 117 (27): 15536– 15545.
- Hameleers, Michael. 2020. "Separating truth from lies: comparing the effects of news media literacy interventions and fact-checkers in response to political misinformation in the US and Netherlands." Information, Communication & Society 37 (2): 1–17.
- Heath, Oliver. 2005. "Party systems, political cleavages and electoral volatility in India: A state-wise analysis, 1998–1999." Electoral Studies 24 (2): 177–199.
- Hochschild, Jennifer L, and Katherine Levine Einstein. 2015. Do Facts Matter?: Information and Misinformation in American Politics. Norman, OK: University of Oklahoma Press.
- Jardina, Ashley, and Michael Traugott. 2019. "The Genesis of the Birther Rumor: Partisanship, Racial Attitudes, and Political Knowledge." Journal of Race, Ethnicity and Politics 4 (1): 60–80.
- Jerit, Jennifer, and Jason Barabas. 2012. "Partisan Perceptual Bias and the Information Environment." The Journal of Politics 74 (3): 672–684.

- Jerit, Jennifer, Jason Barabas, and Scott Clifford. 2013. "Comparing Contemporaneous Laboratory and Field Experiments on Media Effects." Public Opinion Quarterly 77 (1): 256–282.
- Jones-Jang, S Mo, Tara Mortensen, and Jingjing Liu. 2019. "Does media literacy help identification of fake news? Information literacy helps, but other literacies don't." American Behavioral Scientist 00 (0): 1–18.
- Kahan, Dan M, Ellen Peters, Erica Cantrell Dawson, and Paul Slovic. 2017. "Motivated Numeracy and Enlightened Self-Government." Behavioural Public Policy 1 (1): 54– 86.
- Kahne, Joseph, and Benjamin Bowyer. 2017. "Educating for democracy in a partisan age: Confronting the challenges of motivated reasoning and misinformation." American Educational Research Journal 54 (1): 3–34.
- Kaka, Noshir, Anu Madgavkar, Alok Kshirsagar, Rajat Gupta, James Manyika, Kushe Bahl, and Shishir Gupta. 2019. "Digital India: Technology to transform a connected nation." McKinsey Global Institute. March, 2019.

https://www.mckinsey.com/business-functions/mckinsey-digital/ourinsights/digital-india-technology-to-transform-a-connected-nation.

- Kitschelt, Herbert, and Steven I Wilkinson. 2007. Patrons, Clients and Policies: Patterns of Democratic Accountability and Political Competition. Cambridge: Cambridge University Press.
- Kumar, Sanjay, and Pranav Kumar. 2018. "How widespread is WhatsApp's usage in India?" Live Mint. July 18, 2018. https://www.livemint.com/Technology/O6DLmIibCCV5luEG9XuJWL/Howwidespread-is-WhatsApps-usage-in-India.html.

- Li, Jianing. 2020. "Toward a Research Agenda on Political Misinformation and Corrective Information." Political Communication 37 (1): 125–135.
- Lodge, Milton, and Charles S Taber. 2013. The Rationalizing Voter. Cambridge University Press.
- Lupu, Noam, and Kristin Michelitch. 2018. "Advances in Survey Methods for the Developing World." Annual Review of Political Science 21: 195–214.
- Margolin, Drew B, Aniko Hannak, and Ingmar Weber. 2018. "Political Fact-Checking on Twitter: When Do Corrections Have an Effect?" Political Communication 35 (2): 196– 219.
- Masih, Niha, and Joanna Slater. 2019. "U.S.-style polarization has arrived in India. Modi is at the heart of the divide." The Washington Post. May 20, 2019. https://www.washingtonpost.com/world/asiapacific/divided-families-and-tensesilences-us-style-polarization-arrives-in-india/2019/05/18/story.html.
- Mathur, Nandita. 2019. "India's internet base crosses 500 million mark, driven by Rural India." Live Mint. March 11, 2019. https://www.livemint.com/industry/telecom/internet-users-exceed-500-millionrural-india-driving-growth-report-1552300847307.html.
- McGoey, Linsey. 2012. "The logic of strategic ignorance." The British Journal of Sociology 63 (3): 553–576.
- Michelitch, Kristin, and Stephen Utych. 2018. "Electoral Cycle Fluctuations in Partisanship: Global Evidence from Eighty-Six Countries." The Journal of Politics 80 (2): 412–427.

- Mosseri, Adam. 2017. "A New Educational Tool Against Misinformation." Facebook. April 6, 2017. https://about.fb.com/news/2017/04/a-new-educational-tool-againstmisinformation/.
- Munger, Kevin, Mario Luca, Jonathan Nagler, and Joshua Tucker. 2018. "Everyone On Mechanical Turk is Above a Threshold of Digital Literacy: Sampling Strategies for Studying Digital Media Effects." Working Paper. https://csdp.princeton.edu/sites/csdp/files/media/munger-mturk-digital-literacy-

note.pdf.

- Nelson, Jacob L, and Harsh Taneja. 2018. "The small, disloyal fake news audience: The role of audience availability in fake news consumption." New Media & Society 20 (10): 3720–3737.
- Nyhan, Brendan, Ethan Porter, Jason Reifler, and Thomas Wood. 2019. "Taking Fact-Checks Literally But Not Seriously? The Effects of Journalistic Fact-Checking on Factual Beliefs and Candidate Favorability." Political Behavior (42): 939–960.
- Nyhan, Brendan, and Jason Reifler. 2010. "When Corrections Fail: The Persistence of Political Misperceptions." Political Behavior 32 (2): 303–330.
- Pennycook, Gordon, Tyrone D Cannon, and David G Rand. 2018. "Prior Exposure Increases Perceived Accuracy of Fake News." Journal of Experimental Psychology: General 147 (12): 1865–1880.
- Perrigo, Billy. 2019. "How Volunteers for India's Ruling Party Are Using WhatsApp to Fuel Fake News Ahead of Elections." TIME. January 25, 2019. https://time.com/5512032/whatsapp-india-election-2019/.
- Poonam, Snigdha, and Samarth Bansal. 2019. "Misinformation Is Endangering India's Election." The Atlantic. April 1, 2019.

https://www.theatlantic.com/international/archive/2019/04/india-misinformationelection-fake-news/586123/.

- Porter, Ethan, and Thomas J Wood. 2019. False Alarm: The Truth About Political Mistruths in the Trump Era. Cambridge University Press.
- Prior, Markus, Gaurav Sood, and Kabir Khanna. 2015. "You cannot be serious: The impact of accuracy incentives on partisan bias in reports of economic perceptions." Quarterly Journal of Political Science 10 (4): 489–518.
- Roozenbeek, Jon, and Sander Van Der Linden. 2019. "The fake news game: actively inoculating against the risk of misinformation." Journal of Risk Research 22 (5): 570–580.
- Sahoo, Niranjan. 2020. "Mounting Majoritarianism and Political Polarization in India." Carnegie Endowment for International Peace. https://carnegieendowment.org/2020/08/18/mounting-majoritarianism-andpolitical-polarization-in-india-pub-82434.
- Sherman, David K, and Geoffrey L Cohen. 2006. "The Psychology of Self-Defense: Self-Affirmation Theory." Advances in Experimental Social Psychology 38: 183–242.
- Silver, Laura, and Aaron Smith. 2019. "In some countries, many use the internet without realizing it." Pew Research Center. May 02, 2019. https://www.pewresearch.org/fact-tank/2019/05/02/in-some-countries-many-usethe-internet-without-realizing-it/.
- Singh, Shivam Shankar. 2019. *How to win an Indian election: What political parties don't want you to know*. Gurgaon: Penguin Random House.

- Sinha, Pratik, Sumaiya Sheikh, and Arjun Sidharth. 2019. *India Misinformed: The True Story*. Noida: HarperCollins India.
- Sircar, Neelanjan. 2020. "The politics of vishwas: political mobilization in the 2019 national election." Contemporary South Asia 28 (2): 178–194.
- Smith, Jeff, Grace Jackson, and Seetha Raj. 2017. "Designing Against Misinformation." Medium. December 20, 2017. https://medium.com/facebookdesign/designing-against-misinformation-e5846b3aa1e2.
- Taber, Charles S, and Milton Lodge. 2006. "Motivated Skepticism in the Evaluation of Political Beliefs." American Journal of Political Science 50 (3): 755–769.
- Tandoc Jr, Edson C, Zheng Wei Lim, and Richard Ling. 2018. "Defining "Fake News": A typology of scholarly definitions." Digital Journalism 6 (2): 137–153.
- Thachil, Tariq. 2014. "Elite Parties and Poor Voters: Theory and Evidence from India." American Political Science Review 108 (2): 454–477.
- Tucker, Joshua A, Andrew Guess, Pablo Barbera, ´Cristian Vaccari, Alexandra Siegel, Sergey Sanovich, Denis Stukal, and Brendan Nyhan. 2018. "Social Media, Political Polarization, and Political Disinformation: A Review of the Scientific Literature." Hewlett Foundation report. https://eprints.lse.ac.uk/87402/1/Social-Media-Political-Polarization-and-Political-Disinformation-Literature-Review.pdf.
- Tully, Melissa, Emily K Vraga, and Leticia Bode. 2020. "Designing and Testing News Literacy Messages for Social Media." Mass Communication and Society 23 (1): 22–46.
- Uttam, Kumar. 2018. "For PM Modi's 2019 campaign, BJP readies its WhatsApp plan." *Hindustan Times*. September 29, 2018.
 - https://www.hindustantimes.com/india-news/bjp-plans-a-whatsapp-campaign-for-

© Copyright 2020 Sumitra Badrinathan & Center for the Advanced Study of India

2019-lok-sabha-election/story-lHQBYbxwXHaChc7Akk6hcI.html.

- Valenzuela, Sebastian, ´ Daniel Halpern, James E Katz, and Juan Pablo Miranda.
 2019. "The Paradox of Participation Versus Misinformation: Social Media, Political Engagement, and the Spread of Misinformation." Digital Journalism 7 (6): 802–823.
- Vraga, Emily K, Leticia Bode, and Melissa Tully. 2020. "Creating News Literacy Messages to Enhance Expert Corrections of Misinformation on Twitter." Communication Research 00 (0): 1–23.
- Vraga, Emily K, and Melissa Tully. 2019. "News literacy, social media behaviors, and skepticism toward information on social media." Information, Communication & Society: 1–17.
- Walter, Nathan, and Sheila T Murphy. 2018. "How to unring the bell: A meta-analytic approach to correction of misinformation." Communication Monographs 85 (3): 423–441.
- Westwood, Sean J, Shanto Iyengar, Stefaan Walgrave, Rafael Leonisio, Luis Miller, and Oliver Strijbis. 2018. "The tie that divides: Cross-national evidence of the primacy of partyism." European Journal of Political Research 57 (2): 333–354.
- Wire. 2018. "Real or Fake, We Can Make Any Message Go Viral: Amit Shah to BJP Social Media Volunteers." The Wire. September 26, 2018.
 https://thewire.in/politics/amit-shah-bjp-fake-social-media-messages.
- Wittenberg, Chloe, and Adam J Berinsky. 2020. "Misinformation and Its Correction."In Social Media and Democracy: The State of the Field, Prospects for Reform, ed.Nathaniel Persily and Joshua A Tucker. Cambridge University Press.

Wood, Thomas, and Ethan Porter. 2019. "The Elusive Backfire Effect: Mass Attitudes' Steadfast Factual Adherence." Political Behavior 41 (1): 135–163.

Ziegfeld, Adam. 2016. Why Regional Parties? New York: Cambridge University Press.

Online Appendices for Educative Interventions to Combat Misinformation: Evidence From a Field Experiment in India

Contents

A	Summary Statistics	2
B	Survey and Sampling Design	3
С	Flyers	6
D	Dependent Variables	10
E	Enumerator Fixed Effects	15
F	All Stories as DV	17
G	Correlates of Misinformation	19
н	Age and Digital Literacy	21
Ι	True Stories	23
J	Predicted Story Identification	25

A Summary Statistics

Table A.1 provides summary statistics for key variables in this study. Literacy Intervention is a dummy variable indicating random assignment to both treatment groups relative to control. BJP Supporter is a dummy variable indicating respondents' self-reported support for the BJP relative to all other parties. Accurate Priors measures prior beliefs in veracity of news with a battery of four stories (two true and two false); for each story respondents are asked to discern the veracity on a 3-point scale. The variable Accurate Priors calculates the mean accuracy rating across all four stories. Digital Literacy is measured through eight five-point (self-reported) ratings of degree of understanding of WhatsApp-related items. The variable Digital Literacy calculates the mean level of literacy across the eight items. Political Knowledge is measured by a battery of 6 questions of varying difficulty on local and national politics in India; the variable Political Knowledge counts the number of correct answers. WhatsApp Use Frequency measures how frequently respondents use WhatsApp on a 7-point scale ranging from a few times a month to a few times a day. Trust in WhatsApp measures respondents' level of trust in WhatsApp as an accurate medium of receiving news about politics, on a four-point scale.

Statistic	Ν	Mean	St. Dev.	Min	Median	Max
Literacy Intervention	1,224	0.668	0.471	0	1	1
BJP Supporter	1,224	0.684	0.465	0	1	1
Accurate Priors	1,158	0.695	0.196	0	0.750	1
Digital Literacy	1,224	0.758	0.194	0.083	0.833	1
Political Knowledge	1,224	5.000	1.135	0	5	6
WhatsApp Use Frequency	1,224	6.068	0.952	1	6	7
Trust in WhatsApp	1,224	2.729	0.821	1	3	4
Education	1,224	9.388	2.652	1	9	13
Age	1,224	26.646	9.182	18	24	68
Male	1,224	0.911	0.285	0	1	1
Hindu	1,224	0.837	0.369	0	1	1

Table A.1: Summary Statistics

B Survey and Sampling Design

The primary sampling unit, the city of Gaya in Bihar, consists of several electoral polling booths (smallest administrative units). Out of the total number of polling booths, a random sample of 85 polling booths were selected (through a random number generator in the statistical framework R) to serve as enumeration areas.

Within each enumeration area, enumerators were instructed to survey 10-12 households following a random walk procedure. This methodology has the benefits of fast implementation and unpredictability of movement and was chosen over traditional listing methods so that enumerators could spend as little time in the field as possible given the potential for electoral violence. It was also chosen over traditional listing methods due to lack of accurate census data and reliable addresses in the area.

Surveying households within each chosen polling booth area involved choosing a starting point and then proceeding along a path, selecting every kth household. I followed the method similar to that used by the Afrobarometer surveys of picking a sample starting point and then choosing a landmark as near as possible to the sample starting point. Landmarks could be street corners, schools, or water sources, and field enumerators were instructed to randomly rotate the choice of such landmarks. From the landmark starting point, the field enumerator walked in a designated direction away from the landmark and selected the tenth household for the survey, counting houses on both the left and the right. Once they left their first interview they continued in the same direction, selecting the next household after another interval of 10. If the settlement came to an end and there were no more houses, the field enumerator turned at right angles to the right and kept walking, continuing to count until finding the tenth dwelling. Each field enumerator was assigned to only one polling booth, and hence the paths taken during each selection crossed each household only once, increasing the likelihood of a random and unbiased sample.

Once a household is selected, a randomly chosen adult member of the household

was chosen to answer our survey questions after they qualified based on pre-conditions. The three pre-conditions of the survey were (1) access to a personal smartphone (i.e. not a shared household cellphone), (2) connectivity of the phone to working Internet for the past 6 months, (3) usage of WhatsApp on the phone.

Importantly, these qualification conditions resulted in only 20 percent of all houses knocked on having a respondent who was eligible for the study. This is not atypical for Bihar, where only 20-30 percent of citizens have access to the internet. Despite this, the study also had a high response rate. Of all those who were eligible for the study, 94.5 percent agreed to participate. The high participation response rate corresponds to research in face-to-face surveys and in developing countries where response rates tend to be typically higher than in developed countries.

Of the 5.5 percent who refused, enumerator notes suggest that these respondents tended to be older women who (despite having a phone and internet) indicated they would be comfortable if the survey was conducted with a younger member of the household; in some cases they suggested enumerators wait inside the house until a younger member came back home. Once respondents consented to the survey and invited enumerators in their house, no respondent terminated the intervention early or asked that enumerators leave and come back at a different time. Thus, all respondents in the first wave who met the criteria and agreed to the survey completed the intervention in one setting.

The survey pre-conditions ensured that access to WhatsApp and other social media accounts was by the respondent alone, and these restrictions were put into place to ensure that respondents in the study were likely to be exposed to political misinformation over WhatsApp in the months leading up to the election. Sharing mobile phones is especially common among adults in semi-urban and rural India. Further, it is also more common for women than it is for men. Pew survey data from 2019 finds that women are less likely than men to own their own mobile phones, and consequently, significantly more women (20 percent) than men (5 percent) report sharing a device with someone else.

4

These sampling conditions resulted in an uneven age distribution for the study, with about 35 percent of respondents below age 22 and only about 6 percent of the sample above age 45. It also resulted in an uneven gender distribution. Focus group discussions with men and women above the age of 45 showed that people in this age group largely did not own their own cellphones; they reported having shared cellphones used by the entire house or not having access to a phone with working Internet at all. Women, particularly, reported using their husbands' cellphones to communicate and did not report owning their own social media accounts. As a result, only 6 of the women in this sample were above the age of 40.

C Flyers

Respondents were given flyers as part of the intervention. For treatment group respondents, the front side of the flyer included four false political stories that went viral on social media in the months before the 2019 election. The flyer included the photos / screen grabs associated with these false stories along with an explanation for what the correct version of the story is. The back of the flyer contained 6 general tips to spot misinformation. Enumerators explained each bit of information in the flyer and then finally handed the flyers over to respondents. Treatment 1 flyer has pro-BJP false stories, Treatment 2 flyer has anti-BJP false stories, the control flyer is a placebo and has information on plastic pollution. All materials were in Hindi and the survey and intervention were also administered in Hindi. Below I include English translations of the survey materials.

Figure C.1: Treatment 1 – Pro-BJP Flyer (front and back)



Figure C.2: Treatment 2 – Anti-BJP Flyer (front and back)



Figure C.3: Placebo Control Flyer (front and back)



D Dependent Variables

To measure key outcomes of interest, respondents were shown a series of fourteen news stories. These stories varied in content, salience, and critically, partisan slant. Half of the stories were pro-BJP in nature and the other half anti-BJP. Each respondent saw all the fourteen stories, but the order in which they were shown was randomized. Table D.1 lists the fourteen stories shown to respondents. Following each story, two primary dependent variables were measured:

- 1. Perceived accuracy of news stories, with the question "Do you believe this news story is false?" (binary response, 1 if yes, 0 otherwise)
- 2. Confidence in identification of the story as false or real, with the question "How confident are you that the story is real / false?" (4-point scale, 1 = very confident, 4 = not confident at all)

	Story	Party Slant	Veracity
1	Cow urine cures cancer	Pro-BJP	False
2	Photos of militant bloodshed in Kashmir w/ pro-army message	Pro-BJP	False
3	India has not experienced a single foreign terror attack since 2014	Pro-BJP	False
4	Photoshopped image of war hero in BJP attire	Pro-BJP	False
5	Images of the Indian flag projected onto the Statue of Liberty	Pro-BJP	False
6	Rumor that new Indian notes have tracking chips embedded	Pro-BJP	False
7	Rumor that the govt. has installed CCTV cameras in voting booths	Anti-BJP	False
8	Photoshopped images of BJP workers littering the Ganga river	Anti-BJP	False
9	Rumor that BJP workers use duplicate votes to rig elections	Anti-BJP	False
10	Rumors on lack of policing by govt. leading to child kidnapping	Anti-BJP	False
11	Photoshopped image of govt. built Patel statue developing cracks	Anti-BJP	False
12	Rumors of BJP voters hacking voting machines to rig elections	Anti-BJP	False
13	PM Modi has a new radio show on air called Mann Ki Baat	Pro-BJP	True
14	A recent attack killed 40 Indian CRPF soldiers in Kashmir's Pulwama	Anti-BJP	True

Table D.1: Dependent Variable Stories

After the fourteen political stories, two additional dependent variables were measured: self-perceived efficacy of the treatment, and self-reported media literacy. Self-perceived efficacy was measured by asking respondents "How confident are you that you can spot false news from real news?" (4-point scale, 1 = very confident, 4 = notconfident at all). Media literacy was measured in two ways: trust in news received over WhatsApp (4-point scale); and how frequently they forwarded political messages over WhatsApp (6-point scale). Self-reported literacy and efficacy were measured to determine whether the intervention was successful at generating awareness of the problem of misinformation, arguably demonstrated by decreased trust in WhatsApp and forwarding of political stories. Finally, voter turnout was measured. This was done by asking respondents to show enumerators the index finger of their left hand, which, if they voted, would be marked with purple indelible ink. Because respondents were surveyed within a few days of having voted, the presence of an inked finger is a clean and near-perfect measure of voter turnout. Though this may not be true for instances where respondents refuse to show their ink, in this study every respondent willingly showed enumerators their index finger and no one refused.

The analysis in Table D.2 measures the effect of the treatment on self-reported confidence that respondents had in each story being true or false. Confidence was measured on a four-point scale between 0 and 1 for each story with higher numbers indicating more expressed confidence. The dependent variable was calculated as the average confidence level across all stories. While there is no main effect of the treatment on confidence, there is an effect with certain subgroups. Respondents who were more educated and received the intervention were significantly less likely to be confident in their responses. By contrast, men who received the intervention were more likely to be confident in their responses relative to women.

Tables below identify the effect of the intervention on secondary dependent variables measured for this study. The first column estimates the effect of the intervention on

11

	Dependent variable: Confidence in Story Veracity				
	Ave	erage Confidence Le	evel		
	(1)	(2)	(3)		
Literacy Intervention	0.006 (0.006)	0.058 (0.022)	0.045 (0.020)		
Education		0.003 (0.002)			
Male			0.020 (0.017)		
Literacy Intervention Education		0.007 (0.002)			
Literacy Intervention Male			0.044 (0.021)		
Constant	0.937 (0.005)	0.875 (0.018)	0.924 (0.016)		
Observations R ²	1,224 0.001	1,224 0.070	1,224 0.066		
Adjusted R ² Residual Std. Error	0.00004 0.103 (df = 1222)	0.066 0.100 (df = 1218)	0.062 0.100 (df = 1218)		
F Statistic	0.954	18.340	17.181		

Table D.2: ATE and HTE for Confidence in Story Veracity

Note:

=

self-reported confidence in being able to tell the difference between true and false stories, that is, this measures the efficacy of the treatment. Confidence was measured on a three point scale where higher values indicate a greater level of confidence. In Column 2, the dependent variable is self-reported scrutiny of messages; respondents were asked whether they check if messages are true before forwarding them. This is a binary variable. In Column 3, respondents' turnout in the general election is measured. In the final column, I measure trust in WhatsApp on a four-point scale where higher values indicate more trust in the medium.

Table D.3 is the average treatment effect on the four dependent variables described above. Table D.4 is the heterogeneous effect of party identity on the four dependent variables described above.

	Dependent variable:			
	Confidence Message Checking Turno		Turnout	WhatsApp Trust
	(1)	(2)	(3)	(4)
Literacy Intervention	0.001	0.015	0.013	0.041
	(0.023)	(0.026)	(0.030)	(0.040)
Constant	0.170	0.246	0.478	2.539
	(0.019)	(0.021)	(0.025)	(0.033)
Observations	1,224	1,224	1,224	1,224
R2	0.00000	0.0003	0.0002	0.001
Adjusted R ²	0.001	0.001	0.001	0.00004
Residual Std. Error (df = 1222)	0.377	0.425	0.499	0.663
<u>F Statistic (df = 1; 1222)</u>	0.003	0.350	0.192	1.051

Table D.3: Average Treatment Effect on Non-Identification DVs

Note:

	Dependent variable:			
	Confidence	Message Checking	Turnout	WhatsApp Trust
	(1)	(2)	(3)	(4)
Literacy Intervention	0.025 (0.041)	0.016 (0.046)	0.038 (0.054)	0.009 (0.071)
BJP Supporter	0.012 (0.040)	0.022 (0.045)	0.035 (0.053)	0.103 (0.070)
Literacy Intervention x BJP Supporter	0.039 (0.049)	0.002 (0.055)	0.035 (0.065)	0.075 (0.086)
Constant	0.162 (0.033)	0.262 (0.037)	0.454 (0.044)	2.469 (0.058)
Observations R ^Z	1,224 0.003	1,224 0.001	1,224 0.003	1,224 0.003
Adjusted R ²	0.0003	0.002	0.001	0.0004
Residual Std. Error (df = 1220)	0.376	0.425	0.499	0.663
F Statistic (df = 3; 1220)	1.111	0.335	1.377	1.1/5

Table D.4: Heterogeneous Effect of Party on Non-Identification DVs

Note:

E Enumerator Fixed Effects

The endline survey to measure the dependent variable was conducted offline (as a paper survey) for field safety reasons. The main dependent variable consisted of 14 stories, but because the survey was conducted offline, the order of appearance of these stories was pre-determined and limited to 3 random orders. A single enumerator only had access to one of the three random orders. Hence as a robustness check, I replicate the main results with enumerator fixed effects.

Table E.1 replicates results for the main effect of the intervention on the outcome. Results are robust to enumerator fixed effects.

	Dependent variable: Number of Stories Identified as False				
	Pro-BJP Stories	Anti-BJP Stories	Pro-BJP Stories	Anti-BJP Stories	
	(1)	(2)	(3)	(4)	
Literacy Intervention	0.007 (0.058)	0.004 (0.053)			
Literacy + Pro-BJP Fact-Check			0.003 (0.067)	0.001 (0.061)	
Literacy + Anti-BJP Fact-Check			0.017 (0.067)	0.008 (0.061)	
Constant	4.789 (0.060)	5.741 (0.054)	4.789 (0.060)	5.741 (0.054)	
Observations R ² Adjusted R ²	1,224 0.252 0.250	1,224 0.123 0.120	1,224 0.252 0.249	1,224 0.123 0.120	
Residual Std. Error	0.961 (df = 1220)	0.868 (df = 1220)	0.962 (df = 1219)	0.868 (df = 1219)	

Table E.1: Effect of Treatment with Enumerator Fixed Effects

Note:

Table E.2 replicates results with enumerator fixed effects for the heterogeneous effect of party identity. Results are robust to enumerator fixed effects.

	Dependent variable: Number of Stories Identified as False		
	Pro-BJP Stories	Anti-BJP Stories	
	(1)	(2)	
Literacy Intervention	0.254	0.077	
	(0.103)	(0.093)	
BJP Supporter	0.265	0.327	
	(0.102)	(0.092)	
Literacy Intervention x	0.384	0.120	
BJP Supporter	(0.125)	(0.112)	
Constant	4.608	5.521	
	(0.092)	(0.082)	
Observations	1,224	1,224	
R ²	0.258	0.139	
Adjusted R ²	0.255	0.135	
Residual Std. Error (df = 1218)	0.958	0.860	
F Statistic (df = 5; 1218)	84.543	39.252	

Table E.2: Effect of Treatment x Party with Enumerator Fixed Effects

Note:

F All Stories as DV

Below I replicate results where the dependent variable is the number of stories correctly identified given all fourteen stories, true and false. Results hold.

	Dependent variable: Number of Stories Accurately Identified		
	(1)	(2)	
Literacy Intervention Pooled	0.005 (0.097)		
Literacy + Pro-BJP Fact-Check		0.014 (0.112)	
Literacy + Anti-BJP Fact-Check		0.024 (0.113)	
Constant	11.638 (0.080)	11.638 (0.080)	
Observations R ² Adjusted R ² Residual Std. Error F Statistic	1,224 0.00000 0.001 1.604 (df = 1222) 0.002 (df = 1; 1222)	1,224 0.0001 0.002 1.605 (df = 1221) 0.058 (df = 2; 1221)	

Table F.1: Effect of Treatment on Identification of Stories

Note:

Table F.2: Effect of Treatment

Party on Identification of Stories

	Dependent variable: Number of Stories Identified as False	
	(1)	
Literacy Intervention	0.400	
	(0.172)	
BJP Supporter	0.497	
	(0.170)	
Literacy Intervention	0.595	
BJP Supporter	(0.208)	
Constant	11.300	
	(0.140)	
Observations	1,224	
R ²	0.007	
Adjusted R ²	0.005	
Residual Std. Error	1.599 (df = 1220)	
F Statistic	3.067 (df = 3; 1220)	

Note:

G Correlates of Misinformation

Independent of the literacy intervention, it is descriptively interesting for the understudied context of India to understand who is more likely to consume misinformation and more likely to be able to identify news as false. I consider the main effect of several demographic and pre-treatment variables on ability to identify misinformation. The results are presented in Table G.1. For all 12 dependent variable stories taken together, BJP partisans are significantly better at identifying false stories as compared to their non-BJP partisan counterparts. Further, as expected, accurate prior beliefs are more likely to aid in identifying misinformation. Higher levels of digital literacy were negatively associated with identification, underscoring that greater knowledge of WhatsApp leads to more vulnerability to misinformation in this context. However, those who report using WhatsApp more often are more likely to be able to identify misinformation. Interestingly, higher levels of trust in WhatsApp do not correlate with identification of false stories, suggesting that familiarity with the medium itself can make people more vulnerable to misinformation and consequently more likely to share false stories.

With respect to demographic variables, increase in age is associated with a higher capacity to identify misinformation. On the other hand, education has a positive effect on ability to identify false stories.

	Dependent variable: Number of Stories Identified as False		
	Pooled DV : All Stories		
Literacy Intervention	0.060 (0.095)		
BJP Supporter	0.234 (0.113)		
Accurate Priors (Higher = more accurate)	0.480 (0.231)		
Digital Literacy (Higher = more literate)	1.168 (0.252)		
Political Knowedge (Higher = more knowledge)	0.070 (0.046)		
WhatsApp Use Frequency (Higher = more usage)	0.150 (0.047)		
Trust in WhatsApp (Higher = more trust)	0.071 (0.057)		
Education	0.045 (0.018)		
Age	0.022 (0.005)		
Male	0.164 (0.164)		
Hindu	0.185 (0.144)		
Constant	8.987 (0.437)		
Observations	1,158		
R ⁴	0.066		
Adjusted K ⁻	0.057		
Residual Sid. Effor	1.509 (df = 1146) 7.335 (df = 11: 1146)		
	1.000 (u = 11, 1140)		

Table G.1: Main Effect of Demographic and Pre-Treatment Variables

Note:

=

H Age and Digital Literacy

I explore further the relationship between age, misinformation, and digital literacy. The tables below look at age as variable. In Table H.1, I demonstrate that older respondents are better at identification. However in Table H.2, I find that older respondents have lower levels of digital literacy, demonstrating that despite having better digital literacy skills, younger respondents are worse are identifying false stories.

Dependent variable: Number of Stories Identified As False	
(1)	
0.024	
(0.005)	
9.276	
(0.136)	
1,224	
0.019	
0.019	
1.553 (df = 1222)	
24.246 (df = 1; 1222)	

Table H.1: Effect of Age on Identification of Stories

p<0.1; p<0.05; p<0.01

Table H.2: Effect of Age on Digital Literacy

	Dependent variable: Digital Literacy (Higher = More Literate)	
	(1)	
Age (Continuous)	0.001	
	(0.001)	
Constant	0.796	
	(0.017)	
Observations	1,224	
R ²	0.005	
Adjusted R ²	0.004	
Residual Std. Error	0.194 (df = 1222)	
F Statistic	5.716 (df = 1; 1222)	

Note:

I now consider whether the literacy intervention worked better depending on age or digital literacy. In Table H.3 I interact the treatment with age and digital literacy, and do not find an interaction effect.

	Dependent variable: Number of Stories Identified as False	
	(1)	(2)
Literacy Intervention	0.386 (0.306)	0.349 (0.382)
Age (Continuous)	0.032 (0.009)	
Literacy Intervention x Age	0.015 (0.011)	
Digital Literacy (Higher = more literate)		0.984 (0.394)
Literacy Intervention x Digital Literacy		0.454 (0.489)
Constant	10.797 (0.257)	12.378 (0.306)
Observations R ²	1,224 0.016	1,224 0.025
Residual Std. Error (df = 1220) <u>F Statistic (df = 3; 1220)</u>	0.014 1.592 6.673	0.022 1.585 10.270

Table H.3: Effect of Treatment x Age and Digital Literacy

Note:

I True Stories

The outcome measure for this study comprised of more false stories than true (rather than a 50-50 split between true and false stories). This was done to maximize reducing belief in as many false stories as possible. However, several steps were taken to ensure that the imbalance of true vs. false stories did not affect the efficacy of the treatment. Before measuring the outcomes, respondents were told that some of the stories were false and some true, likely reducing the urge to default to the stories being false. Further, with the comprehension check, respondents were not only asked whether stories were true or false but were also asked how they identified the veracity of these stories. Importantly, a majority of respondents in the treatment groups said that their responses were motivated by enumerators teaching them about these stories during the household visit, rather than having learnt about the stories on the news or through a friend. Further, enumerators were instructed for this question to not read out response options aloud, but to allow respondents to organically speak about their views on the false stories in a way that minimized the ability of respondents to provide socially desirable answers. Thus, much care was taken in the experiment to ensure that the skew towards false stories would not impact respondents' answers.

I now analyze whether the treatment worked for the two true stories alone. Results are in Table I.1. I find that the perceptions of veracity of these stories did not depend on the treatment. However, respondents accurately classified a high proportion of the true stories, 76 percent and 95 percent respectively.

	Dependent variable: Accurate Identification	
	1st True Story	2nd True Story
Literacy Intervention	0.009 (0.026)	0.005 (0.013)
Constant	0.776 (0.021)	0.951 (0.010)
Observations R ² Adjusted R ² Residual Std. Error (df = 1222) F Statistic (df = 1; 1222)	1,224 0.0001 0.001 0.421 0.134	1,224 0.0001 0.001 0.209 0.171

Table I.1: Identification of True Stories

Note:

J Predicted Story Identification

Here I visualize the interaction heterogeneous effect from partisan identity. I graph the predicted values from the interaction model in the bar plots below. In Figure J.1 I plot the predicted number of stories identified among pro-BJP stories; in Figure J.2 I do the same for anti-BJP stories.



Figure J.1: Predicted Identification of Pro-BJP Stories



Figure J.2: Predicted Identification of Anti-BJP Stories