# Doubt the Messenger:
# The reputation cost of fact-checking

August 17, 2024

**Abstract**

This article presents a novel experiment measuring the reputation and ideological costs fact-checkers face when informing voters about the accuracy of messages. The study examines how exposure to the counter and pro-attitudinal fact-checking messages impact the perceived quality and ideological leaning of fact-checking organizations. In a well-powered, pre-registered survey experiment conducted during the 2021 mid-term election in Argentina when COVID-19 was a polarizing issue, we exposed 5,757 respondents to real tweets about the number of COVID-19 cases in Argentina, followed by fact-checking corrections. Results show that pro-attitudinal messages increased the quality rating of the fact-checker, *Chequeado*, and made respondents perceive the organization ideologically closer to their own views. Counter-attitudinal fact-checking also increases the perceived quality but has no significant effect on ideological contrast of the fact-checker. Results from this experiment are important to devise fact-checking interventions that are reputation-improving and support the organization's long-term mission.

**Word Count: 3988 words**

# 1 The Reputation Dilemma

"No happier to be here than you are to have me:

Nobody likes the man who brings bad news."

Antigone, Sophocles

Coordinated misinformation campaigns on social media have sought to influence elections, voter turnout, and voting decisions (Bovet and Makse, 2019; Recuero, Soares and Gruzd, 2020; Mello, 2020), increased vaccine hesitancy during the COVID-19 pandemic (Loomba et al., 2021), and fostered distrust in scientific information about climate change (Van der Linden et al., 2017). In response, social media companies, media organizations, and policymakers have developed strategies to pre-empt and debunk online falsehoods, working in collaboration with fact-checkers—professional, independent organizations that verify dubious online claims and report on their accuracy. This strategy is supported by extensive research highlighting the substantial positive impact of fact-checking corrections on individuals' ability to discern true from false information (Walter et al., 2020; Nyhan, 2021; Porter and Wood, 2021; Brashier et al., 2021; Bode and Vraga, 2018).

Reputation is critical to the mission of the fact-checkers. To succeed, fact-checking organizations require readers to accept their TRUE and FALSE adjudications [1] even when the candidates and parties they support may be politically affected by this decision. An important condition for a successful adjudication is that users recognize the fact-checker as a quality and unbiased source. Indeed, there is a wealth of research showing that counter-attitudinal fact-checking adjudications of TRUE or FALSE information are often met with disbelief by users, who often distrust the accuracy (quality) or intentions (unbiasedness) of the correction (Brandtzaeg, Følstad and Chaparro Domínguez, 2018; Brashier et al., 2021). Furthermore, reputation capital may be lost when adjudicating content. This raises the possibility that future interventions will

---

[1]Throughout the article, we use the term *adjudication*, rather than the more frequent term *correction* to allude to the fact that Fact-Checking organizations use several distinct labels to verify the accuracy of rumors and claims circulating online. FALSE ratings are more frequently used, particularly when addressing online misinformation. However, these organizations also routinely use TRUE ratings to confirm that a particular rumor is indeed true, most likely in answering public requests for adjudicating public statements from elected officials. In the case of *Chequeado*, the label "true" is used in approximately one-third of their publications.

be less effective and that, over time, the fact checker's perceived quality and ideological integrity will decline, threatening its viability as an effective strategy to counter beliefs for misinformation.
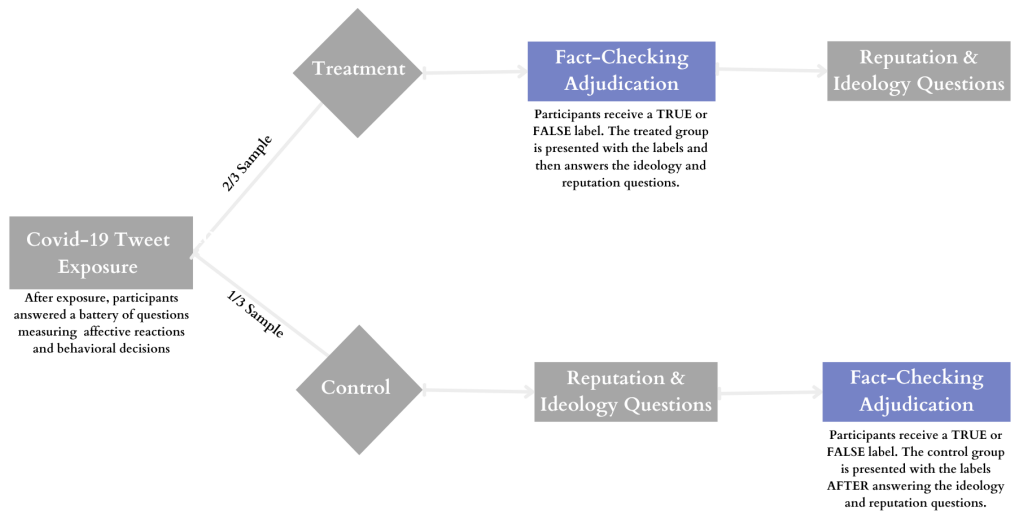
This article describes results from a well-powered and pre-registered survey experiment [2] measuring the reputation costs fact-checkers face when informing voters about the accuracy of online rumors. Our experiment exposes respondents to one of two real social media posts posted by leading Argentine news organizations, which reported that Argentina had the highest number of COVID-19 cases per day on the previous day. Both tweets were published six months apart by *Infobae* with an identical text. On October 2, 2020, Infobae's publication was adjudicated as FALSE by the leading Argentine fact checker, *Chequeado*. The correction was posted on their website and their social media accounts. By contrast, an identically worded article published on May 20, 2021, was factually correct. Depending on the day, content that may be accurate or inaccurate allows researchers to assess the effect of fact checks without requiring deception at the time of implementation. After respondents were treated to one of the two original Infobae tweets, but before exposure to a Fact-checking *Chequeado*'s adjudication, the control group answered questions about the perceived reputation and ideological position of a variety of news organizations and politicians, among which we include the fact checker *Chequeado*. The treated groups, on the other hand, were presented with either a TRUE or a FALSE adjudication before answering the reputation and ideology questions. The experiment assesses differences in the fact checker's perceived reputation and ideological placement between the untreated group and those treated with the TRUE or FALSE adjudication. The publication by *Infobae* aligns well with charges frequently leveled against Argentine President Alberto Fernandez by the opposition, who accused his administration of underperforming during the pandemic. Therefore, our design allows us to measure how partisan-motivated reasoning (Lodge and Taber, 2013) moderates the reputation and ideological cost of fact checks in a polarized political environment. We are the first to explore how assimilation and contrast effects moderate the perceived ideological leaning and the reputation assessment of fact-checking organizations.

---

[2]Our pre-registration is available at https://osf.io/g2mut/?view_only=2b7f93ea495c430b8fd9a19dda79403a

Our results show that respondents reported *Chequeado* as a higher quality organization when treated to a pro-attitudinal adjudication (i.e. the TRUE adjudication for opposition supporters and the FALSE adjudication for government supporters). We identify a large and precise increase (*Cohen's d* of 0.35 SD) (Cohen, 2013) on the reputation scores of *Chequeado*. When treated with counter-attitudinal adjudications, the treatment effects are smaller (*Cohen's d* of 0.1 SD) but still statistically significant, indicating the effects have a positive valence shock, in some sense overcoming voters' partisan motivations. However, the effect of pro- or counter-attitudinal adjudications had a more nuanced effect on perceived ideological bias. On average, respondents perceived *Chequeado* as ideologically closer (assimilation effect) only when exposed to a pro-attitudinal adjudication. However, the effect on ideology is mostly determined by ideological reactions from leftists and pro-government supporters, while right-wing voters do not react ideologically to pro-attitudinal corrections. Results are null when participants are exposed to counter-attitudinal adjudication (contrast effect).

## 2 Experimental Design

Figure 1 summarizes our design. We start our design by exposing respondents to a Tweet that reports the number of COVID-19 deaths in Argentina. We asked respondents if they would share the Tweet and how this Tweet made them feel. We then distract respondents with other questions and split the sample into our treatment and control groups. For our treated group, we proceeded to adjudicate whether the initial Tweet was TRUE or FALSE. We asked if they would share the fact checker's adjudication and whether they believed the original Tweet was true or false. Most importantly, one-third of the individuals are asked questions about the perceived quality and ideological standing of different politicians and news organizations, including *Chequeado*. In contrast, 2/3 of the respondents are asked these same questions after the fact check. Therefore, the design measures the reputation and perceived ideological location of the fact checker before and after the intervention. As a result, we can assess whether pro- or counter-attitudinal fact checks alter the fact checker's perceived quality and ideological position.

a) Survey Flow



b) Original Tweets and Fact-Checking Adjudication

**Figure 1** Figure A depicts the survey design. Figure B presents the original tweets and adjudications.

# 3 Main Hypotheses

The experimental design yields the following hypothesis that addresses questions of **social media sharing**, **reputation**, and **ideological distance**.

The first hypothesis of our study expects "liking" and "sharing" to be more frequent among respondents treated with pro-attitudinal content. Therefore, we expect fewer "likes" and fewer "shares" of the Infobae post by government supporters than opposition supporters. This is consistent with social media behavior sharing behavior that is driven by cognitive congruence and attention (Aruguete, Calvo and Ventura, 2021).

$HT_1$: Pro-attitudinal confirmations and refutations messages will be shared at higher rates than counter-attitudinal messages.

The second hypothesis of our study evaluates how partisan-motivated reasoning moderates the perceived quality of the fact checker upon participants' exposure to an adjudication. The effect of motivated reasoning on news organizations' perceived quality has been well documented in the communications literature (Ardèvol-Abreu and Gil de Zúñiga, 2017; Lee, 2005, 2012), and to a lesser extent in the political science literature that is often more interested in political parties valence advantages (Adams, Merrill III and Grofman, 2005). While different in scope and size from traditional news organizations, fact-checkers are integral to today's news media environment. Fact-checkers are frequently cited as authorities to adjudicate misinformation intent to partisans during elections and often face significant backlash from the affected groups that amplified misinformation content. We expect:

$HT_2a$: Pro-attitudinal confirmations and refutations will increase the perceived quality of the fact checker.

$HT_2b$: Counter-attitudinal confirmations and refutations will decrease the perceived quality of the fact checker.

The third hypothesis of our study evaluates the consequences of pro-attitudinal and counter-attitudinal adjudications on the perceived ideological distance between the fact checker and the respondent. This hypothesis connects the experiment to the literature on *assimilation* and

*contrast*, as described in Banks et. al (2021). We discuss our measure of ideological distance, assimilation, and contrast in the supplemental information file and in section 3 of this article:

$HT_3a$: Pro-attitudinal confirmations and refutations will decrease the perceived ideological distance between the respondent and the fact checker (assimilation effect).

$HT_3b$: Counter-attitudinal confirmations and refutations will increase the perceived ideological distance between the respondent and the fact checker (contrast effect).

## 4 Estimation

We estimate treatment effects for all the hypotheses using OLS regression with robust standard errors and covariate adjustment. We selected covariates by fitting a LASSO regression in the entire sample for every outcome variable, with the penalty term selected by cross-validation. These variables include standard demographics, self-reported perceptions about the country's economy and personal economic conditions, self-reported social media usage, ideological placement, partisan preferences, and measures of institution trust, asked before treatment exposure. SI Section 4 presents full regression tables, and models with and without covariate adjustment.

**Dependent Variables:** To test hypotheses 1 (*pro-attitudinal sharing*), we use participants self-reported reaction (*ignore, reply, comment, or like* the tweet published by the Fact-Checking organization. For hypotheses 2 (*reputation models*), our primary models use participants' assessment of *Chequeado* quality using a five-star scale, in which we explain to participants that zero stars mean "poor-quality" news organization and five stars mean "high-quality" news organization. Lastly, the primary models testing hypothesis 3 (*ideogical distance*) uses the absolute distance between the self-reported ideology of the respondent and the reported ideology of *Chequeado*. The ideology variable ranges from 1 (very progressive) to 7 (very conservative).

**Independent Variables:** Our primary right-hand variable measures if participants were exposed to a pro-attitudinal or counter-attitudinal fact-checking adjudication. To build this measure, we take into consideration the interaction between the framing of the fact-checking correction (TRUE or FALSE) and the vote choice of the respondent (*Cambiemos*/right-wing or Frente de

Todos ($FdT$)/left-wing) [3]. When respondents from the left are randomly assigned to a TRUE adjudication, confirming the accuracy of the tweet reporting world record COVID-19 deaths in Argentina during the leftist administration, we consider these voters exposed to a counter-attitudinal adjudication. Meanwhile, when respondents from the same party are exposed to a fact-checking adjudication saying it is FALSE the tweet reporting world record COVID-19 deaths in Argentina, we consider those as exposed to pro-attitudinal adjudication. We use the same logic, but of course, switch the effects of TRUE and FALSE to build the pro and counter attitudinal measure for right-wing voters in our sample [4]. For hypothesis 1, we estimate the models using the entire sample of participants to measure the effects of pro-attitudinal reactions to the Fact-Checking tweet. For hypotheses 2 and 3, we used participants assigned to the control group, who answered the reputation and ideology questions before receiving the pro and counter-attitudinal adjudication as the baseline group in the regression model.

**Data:** Participants were recruited from Netquest's online panel of Argentine respondents. Participants were at least 18 years old of age and nationals from Argentina. The survey sample included 5,757 respondents in Argentina. The number of participants met national representative samples for each country and guaranteed a well-powered study, capable of identifying with 80% of power effect as small as 0.1 standardized effects with no controls, which is a benchmark considered a small effect in the literature (Cohen, 2013). The survey was conducted between November 4 and December 7.

## 5 Results

**Pro-Attitudinal Sharing:** The first hypothesis of this study, $HT_1$, expected pro-attitudinal messages to be shared at a higher rate than counter-attitudinal messages. We expected supporters

---

[3]We ask participants "if the general presidential election were to take place next week." In a multi-party system, presidential vote choice is often a more appropriate measure for partisan preferences than classic partisanship questions (Calvo and Ventura, 2021; Samuels and Zucco, 2018). To show that our results are robust to different measurement choices, we test our hypothesis in the supplemental materials (SI Section 6) using self-reported partisanship as our measure of partisan identity. Results go in the same direction as those presented in the manuscript.

[4]Since our hypotheses focus on directional effects, we remove from the models for hypotheses 2 and 3 voters who self-reported voting blank
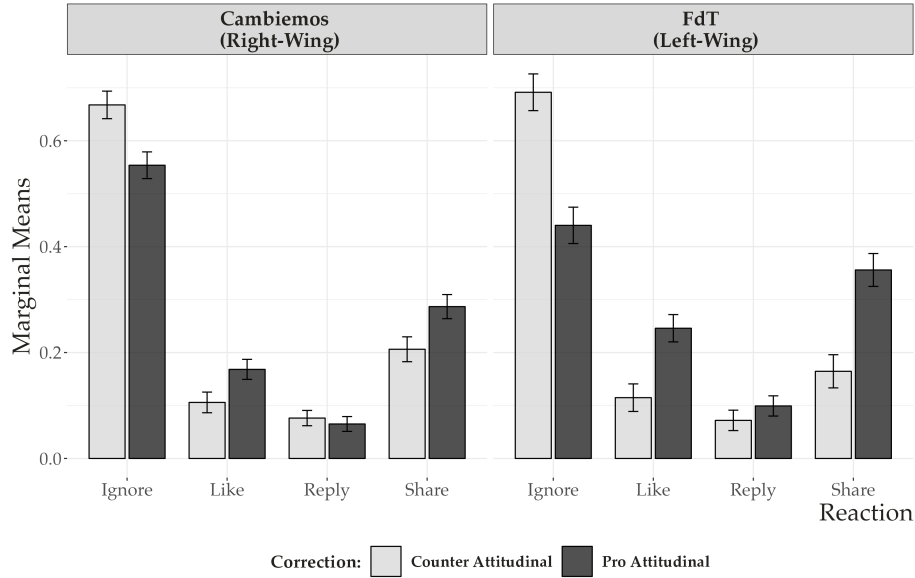
**Figure 2** The vertical axis describes the Marginal Means of respondents that "Ignore", "Like", "Reply" or "Share" the adjudication. Participants were allowed to select multiple responses

of *Cambiemos* (right-wing voters) to share the TRUE adjudication at higher rates than the supporters of President Fernandez, the *FdT*. Meanwhile, we expected the *FdT* voters (left-wing) to share the adjudication of FALSE at a higher rate than *Cambiemos*. Pro- and counter-attitudinal sharing are also crucial preconditions for any testing of the reputation and ideological placement effects. The decision to share content that aligns with the voters preferences is an important marker of the positive and negative content conveyed by the treatments.

Figure 2 confirms the pro- and counter-attitudinal response to the adjudication by *Chequeado* and is expected by H1. Across both partisan groups, voters are more likely to share and like pro-attitudinal adjudication and more likely to ignore counter-attitudinal. We present these results as marginal means, and in the appendix, we present marginal effects for pro-attitudinal sharing (SI Figure 1). The effects are considerably larger among *FdT* voters, with pro-attitudinal behavior for "like" and "share" being two times the marginal means of counter-attitudinal. These findings are consistent with other studies looking at pro-attitudinal sharing of fact-checking corrections Aruguete et al. (2021); Walter et al. (2020) and serve as a validity check, showing that respondents understood the publication and reacted according to expectations.

**Reputation Effects:** The second hypothesis of this study, $HT_2a$, stated that pro-attitudinal confirmations and refutations would increase the perceived quality of the fact checker. Figure 3 shows the average treatment effect when respondents are treated with the pro-attitudinal and counter-attitudinal adjudication. Dark blue points indicate the average effect across the parties, light blue indicates the effects for right-wing voters, and red dots for left-wing voters.

We recover strong support for positive pro-attitudinal reputation shock. Upon being exposed to a pro-attitudinal adjudication, participants increase their perceived reputation of Chequeado by 0.35 SD ($t = 8.68$, $p$-value $< 0.01$). Interestingly, contrary to $HT_2b$, results also show more moderate but still positive gains vis-á-vis participants in the control group, with respondents reporting a higher reputation when reading the counter-attitudinal post. Given that, no information is provided to the respondents about *Chequeado*, this increase may reflect an issue-change effect rather than a net gain, as fact-checking on COVID-19 may already carry an independent positive charge when compared to fact-checking statements by partisans or elected officials. While partisan-motivated reasoning moderates the reputation costs of fact-checking adjudications, we show that at least on the COVID-19 issue, even counter-attitudinal corrections bring positive reputation gains for fact-checking organizations.
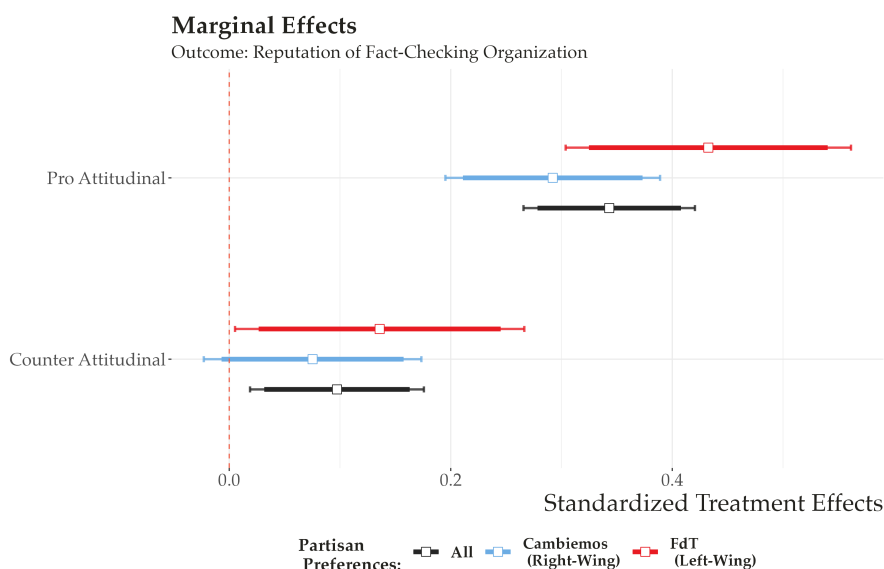


**Figure 3** Standardized Average treatment effect (ATE) for the reputation models with no controls.

Lastly, figure 4 provides a different view of the reputation effects of pro- and counter-attitudinal adjudications, this time conditional on both party and ideology. As it is possible to observe, counter-attitudinal adjudications have a more detrimental effect on the conservative and very conservative subgroup of supporters of *Cambiemos*. The difference between the pro-attitudinal and counter-attitudinal adjudication for supporters of the *FdT*, on the other hand, is not affected by the self-reported ideological location of the respondent. Next we turn our attention to the estimates of assimilation and contrast effects.
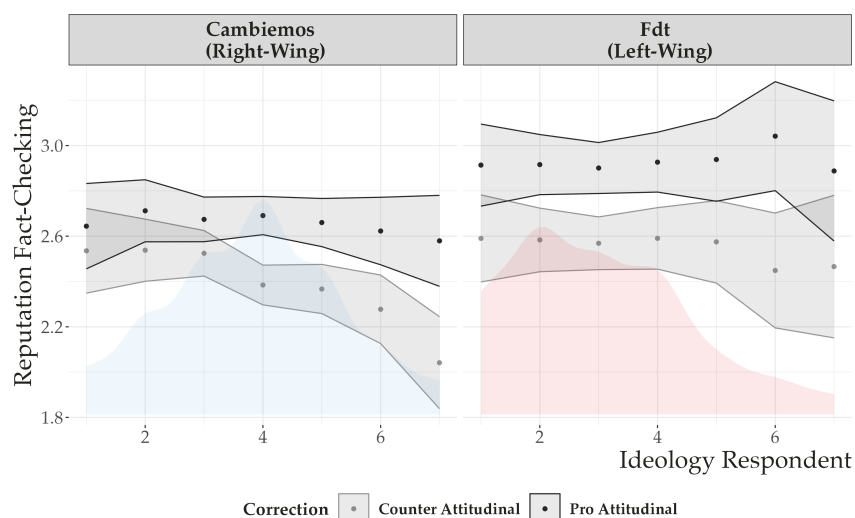


**Figure 4** Conditional effect of ideology by Treatment Condition on Reputation.

**Ideological Bias: Assimilation and Contrast**    The third hypothesis of our study evaluates the consequences of pro-attitudinal and counter-attitudinal adjudications on the perceived ideology of the fact checker. This hypothesis, $HT_3a$, expects pro-attitudinal adjudications to decrease the perceived ideological distance between the respondent and the fact checker (assimilation effect). On the other hand $HT_3b$, we expect counter-attitudinal messages to increase the perceived ideological distance between the respondent and the fact checker (contrast effect).

Figures 5 present the results. As pre-registered, we find that exposure to pro-attitudinal messages decreases the perceived distance between the respondent and the fact checker (-0.022

SD, ($t = -2.35$, $p$-value $< 0.05$). When looking at the effects by party, we see the effects are mostly driven by left-wing voters and null for right-wing voters. Contrary to hypothesis 3b, we do not find statistically significant effects for exposure to counter-attitudinal fact-checking messages. In other words, when exposed to a pro-attitudinal correction, voters perceive the fact-checking organization as moving closer to their own ideological position. However, this effect only exists among left-wing voters; conservatives are inelastic to ideological assimilation even when reading a pro-attitudinal correction.
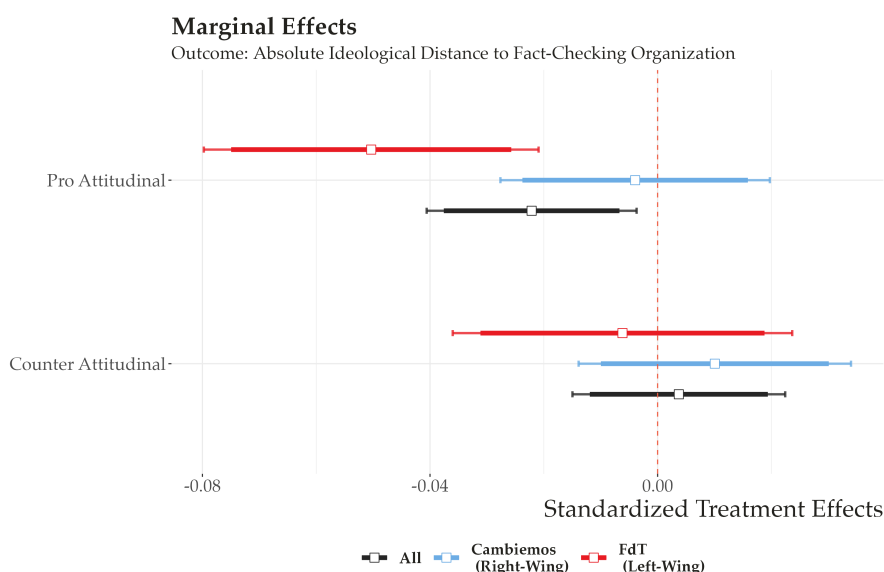


**Figure 5** Standardized Average treatment effect (ATE) for the directional effect of the treatment on ideological distance.

Figure 6 clarifies the results using a standard description of assimilation (positive slope) and contrast (negative slope) (See SI section 7 for a more extensive discussion about assimilation and contrast). Results show a statistically significant change among the $FdT$ (left-wing) respondents, with the aggregate effect moving from a "contrasted" relationship after observing the counter-attitudinal adjudication to an assimilated relationship after the pro-attitudinal adjudication. As depicted in the standardized treatment effects, assimilation and contrast effects are not identified for treated respondents that support *Cambiemos*. To ensure the robustness of our findings, we conducted a set of placebo tests presented at SI Section 5. We re-estimate
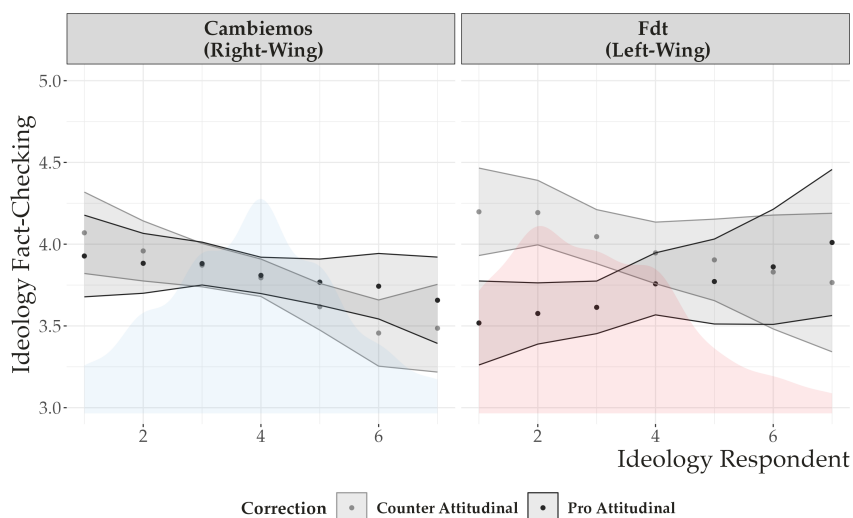
**Figure 6** Marginal Means of *Chequeado*'s ideology conditional on the self-reported ideology of the respondent.

all the primary models of the paper using other media organizations that are also included in our survey instrument as a placebo. We do not find effects for most of our placebo estimates.

## 6  Discussion

This article implements a survey experiment to measure the reputation cost of publishing pro- and counter-attitudinal fact checks in a polarized political environment. Our pre-registered hypotheses expected pro-attitudinal fact checks to increase the perceived reputation of the fact-checker and reduce the ideological distance to readers. Concurrently, we expected counter-attitudinal fact checks to reduce the fact-checkers reputation and increase the perceived ideological distance to readers. Results confirm the expected pro-attitudinal effects on reputation but are less robust when measuring ideological assimilation and contrast effects of the adjudication. In particular, the results of ideological distance are statistically significant only among supporters of the *FdT*.

Some limitations of this study are worthy of notice. First, while our study is internally valid, exposure and attention to the treatments are higher in surveys than in corrections that

organically circulate on social media. Our study describes the short-term effect of pro- and counter-attitudinal corrections, but the expected effects could be lower in the wild, given that social media users have competing information to attend to. A second limitation of this study is that treated respondents acquire information about *Chequeado* when they read the correction, while the control group does not. Given that we cannot provide information about *Chequeado* to the control group without biasing the study, the treated group may not only alter their prior perception of the organization but may also be exposed to it for the first time. Given that *Chequeado* is Argentina's most important fact-checking organization, it is highly likely that most respondents have been exposed to one of their corrections.

What do our results tell about the overall dynamics of misinformation corrections and the long-run mission of fact-checkers? Our findings validate that fact-checking organizations can maintain long-term reputation stocks, particularly when publishing adjudication across the political divide. However, context matters. Our results are encouraging in a hypothetical environment of balanced production and consumption of rumors. However, in a more realistic context where the production of misinformation is uneven, coming more from one political side than another (González-Bailón et al., 2023), fact-checkers live in a reputational dilemma; those continuously exposed to pro-attitudinal messages will update their priors about the ideological position of fact-checkers (assimilation effect), affecting the general equilibrium ideological position of these organizations. In the long run, these updates may lead to increased perceptions of ideological bias against the fundamental tools in mitigating beliefs for misinformation.

# References

Adams, James F, Samuel Merrill III and Bernard Grofman. 2005. *A unified theory of party competition: A cross-national analysis integrating spatial and behavioral factors.* Cambridge University Press.

Ardèvol-Abreu, Alberto and Homero Gil de Zúñiga. 2017. "Effects of editorial media bias perception and media trust on the use of traditional, citizen, and social media news." *Journalism & mass communication quarterly* 94(3):703–724.

Aruguete, Natalia, Ernesto Calvo and Tiago Ventura. 2021. "News by Popular Demand: Ideological Congruence, Issue Salience, and Media Reputation in News Sharing." *The International Journal of Press/Politics* 4.

Aruguete, Natalia, Ingrid Bachman, Ernesto Calvo, Sebastian Valenzuela and Tiago Ventura. 2021. "Truth be told: Cognitive and affective moderators of selective sharing of fact-checks on social media.".

Banks, Antoine, Ernesto Calvo, David Karol and Shibley Telhami. 2021. "# polarizedfeeds: Three experiments on polarization, framing, and social media." *The International Journal of Press/Politics* 26(3):609–634.

Bode, Leticia and Emily K Vraga. 2018. "See something, say something: Correction of global health misinformation on social media." *Health communication* 33(9):1131–1140.

Bovet, Alexandre and Hernán A. Makse. 2019. "Influence of Fake News in Twitter during the 2016 US Presidential Election." *Nature Communications* 10:7.

Brandtzaeg, Petter Bae, Asbjørn Følstad and Maria Ángeles Chaparro Domínguez. 2018. "How journalists and social media users perceive online fact-checking and verification services." *Journalism practice* 12(9):1109–1129.

Brashier, Nadia M, Gordon Pennycook, Adam J Berinsky and David G Rand. 2021. "Tim-

ing matters when correcting fake news." *Proceedings of the National Academy of Sciences* 118(5):e2020043118.

Calvo, Ernesto and Tiago Ventura. 2021. "Will I get Covid-19? Partisanship, social media frames, and perceptions of health risk in Brazil." *Latin American politics and society* 63(1):1–26.

Cohen, Jacob. 2013. *Statistical power analysis for the behavioral sciences.* Routledge.

González-Bailón, Sandra, David Lazer, Pablo Barberá, Meiqing Zhang, Hunt Allcott, Taylor Brown, Adriana Crespo-Tenorio, Deen Freelon, Matthew Gentzkow, Andrew M Guess et al. 2023. "Asymmetric ideological segregation in exposure to political news on Facebook." *Science* 381(6656):392–398.

Lee, Eun-Ju. 2012. "That's not the way it is: How user-generated comments on the news affect perceived media bias." *Journal of Computer-Mediated Communication* 18(1):32–45.

Lee, Tien-Tsung. 2005. "The liberal media myth revisited: An examination of factors influencing perceptions of media bias." *Journal of Broadcasting & Electronic Media* 49(1):43–64.

Lodge, Milton and Charles S Taber. 2013. *The rationalizing voter.* Cambridge University Press.

Loomba, Sahil, Alexandre de Figueiredo, Simon J Piatek, Kristen de Graaf and Heidi J Larson. 2021. "Measuring the impact of COVID-19 vaccine misinformation on vaccination intent in the UK and USA." *Nature human behaviour* 5(3):337–348.

Mello, Patrícia Campos. 2020. *A máquina do ódio: notas de uma repórter sobre fake news e violência digital.* Companhia das Letras.

Nyhan, Brendan. 2021. "Why the backfire effect does not explain the durability of political misperceptions." *Proceedings of the National Academy of Sciences* 118(15):e1912440117.

Porter, Ethan and Thomas J Wood. 2021. "The global effectiveness of fact-checking: Evidence from simultaneous experiments in Argentina, Nigeria, South Africa, and the United Kingdom." *Proceedings of the National Academy of Sciences* 118(37):e2104235118.

Recuero, Raquel, Felipe Bonow Soares and Anatoliy Gruzd. 2020. "Hyperpartisanship, Disinformation and Political Conversations on Twitter: The Brazilian Presidential Election of 2018." *Proceedings of the International AAAI Conference on Web and Social Media* 14:569–578.

Samuels, David J and Cesar Zucco. 2018. *Partisans, antipartisans, and nonpartisans: Voting behavior in Brazil.* Cambridge University Press.

Van der Linden, Sander, Anthony Leiserowitz, Seth Rosenthal and Edward Maibach. 2017. "Inoculating the public against misinformation about climate change." *Global challenges* 1(2):1600008.

Walter, Nathan, Jonathan Cohen, R Lance Holbert and Yasmin Morag. 2020. "Fact-checking: A meta-analysis of what works and for whom." *Political Communication* 37(3):350–375.