*Research Article*

# State media tagging does not affect perceived tweet accuracy: Evidence from a U.S. Twitter experiment in 2022

*State media outlets spread propaganda disguised as news online, prompting social media platforms to attach state-affiliated media tags to their accounts. Do these tags reduce belief in state media misinformation? Previous studies suggest the tags reduce misperceptions but focus on Russia, and current research does not compare these tags with other interventions. Contrary to expectations, a preregistered U.S. experiment found no effect of Twitter-style tags on belief in false state media claims, seemingly because they were rarely noticed. By contrast, fact-check labels decreased belief in false information from state outlets. We recommend platforms design state media tags that are more visible to users.*

Authors: Claire Betzer (1), Montgomery Booth (1), Beatrice Cappio (1), Alice Cook (1), Madeline Gochee (1), Benjamin Grayzel (1), Leyla Jacoby (1), Sharanya Majumder (1), Michael Manda (1), Jennifer Qian (1), Mitchell Ransden (1), Miles Rubens (1), Mihir Sardesai (1), Eleanor Sullivan (1), Harish Tekriwal (1), Ryan Waaland (1), Brendan Nyhan (1)
Affiliation: (1) Department of Government, Dartmouth College, USA

## Research questions

- Do state-affiliated media tags reduce beliefs in tweets containing false information?
- Do fact-checking tags decrease misperceptions more than state-affiliated media tags?
- Does tagging tweets as state-affiliated media reduce trust in the outlet publishing them?
- Do the effects of tagging tweets as state-affiliated media vary between a negatively perceived country and a neutrally perceived country?
- Does tagging true tweets as state-affiliated media reduce their perceived accuracy?

## Essay summary

- In an experiment conducted on Amazon Mechanical Turk (May 7–14, 2022; *N* = 2,555), we tested the effects of exposure to the state-affiliated media tags used on Twitter in 2022 on U.S. respondents' belief in false information from state media outlets.

---

- We found no evidence that state media tags changed the perceived accuracy of false claims from state media and confirmed that fact-check labels reduce the perceived accuracy of misinformation.
- These null effects suggest that users may not have noticed Twitter's state-affiliated media tags, indicating a weakness in their design.

## Implications

Ownership of media firms around the world primarily falls into two categories: media owned by the state and privately owned media (Djankov et al., 2003). Some media outlets that are state-owned or receive government support remain independent from the government. Other outlets are directly controlled by the state and are used by the government to shape political narratives (Djankov et al., 2003; Dragomir & Söderström, 2021; Walker & Orttung, 2014). Although state-affiliated media outlets can exist in democratic countries, this model of state control is in its "most potent" form under authoritarian regimes, where state media outlets are used to spread propaganda in support of the government (Dragomir & Söderström, 2021; Walker & Orttung, 2014).

While television has historically been a prominent form of state-controlled media, social media platforms also provide a mechanism for authoritarian regimes to manipulate information flows and propagate misinformation (Arnold et al., 2021; Bastos & Farkas, 2019; Bradshaw & Howard 2018). State media outlets from authoritarian countries like Russia and China have turned to sites like Facebook and Twitter (now X) to conduct "influence operations" that challenge the existing global order (DiResta et al., 2019; Kinetz, 2021; Xu & Wang, 2022). These government-controlled outlets have sought to broaden their appeal over time by adapting source names and appearing more mainstream to news consumers (e.g., Tiffany, 2022). Russia, in particular, has long sought to use propaganda to exacerbate divisions in the West (Osnos et al., 2017). Chinese state media outlets are also very active in trying to shape perceptions of the regime, including promoting misinformation about the COVID-19 pandemic (Cook, 2020; Molter & DiResta, 2020). Even authoritarian states that are not direct U.S. competitors, like Serbia, use state media platforms to spread propaganda (Mujanović, 2022).

Because the names of state media accounts may be unfamiliar, social media platforms have introduced labels and tags identifying them to users. Twitter, for instance, introduced labels in August 2020 (Robertson, 2020). As CNN wrote at the time, state media accounts after the change "show a podium with a microphone and a label that says 'state-affiliated media'" in gray text under their username and in their bio (Gold, 2020). By contrast, Twitter's fact-check labels at the time the study was conducted appeared in blue underneath the tweet in question and said, for example, "False information: Checked by independent fact-checkers."[2] At the time data was collected for this study in May 2022, YouTube, Facebook, Twitter, and Instagram all provided warnings that certain accounts or posts were from state-affiliated media (Finnegan & Thorbecke, 2021; Gold, 2018; Jackson, 2022). Elon Musk subsequently completed his takeover of Twitter in October 2022. In April 2023, Twitter first labeled National Public Radio as "state-affiliated media" (Folkenflik, 2023) and then dropped all state media labels (Reuters, 2023), seemingly causing views and engagement with state media content to increase (Klepper, 2023; Sadeghi et al., 2023). Most recently, in August 2023, Meta added "state-controlled media" labels to Threads, its text-based social network (Rosen, 2023).[3]

---

[2] Figure 2 below provides an example of how the state media tags and fact-check labels appeared as implemented in our study.

[3] Hundley et al. (2023) provide an instructive overview of how Meta defines state-controlled media and implements its policies.

Source-level labels of untrustworthy sources like authoritarian state media outlets offer a scalable alternative to attaching fact-checking tags to individual claims. Compared to fact checks, however, state media tags remain understudied. At the time this study was conducted, only two published studies had considered the effect of state media tags on the perceived accuracy of the content in social media posts; both found that tags tend to reduce users' belief in false content (Arnold et al., 2021; Nassetta & Gross, 2020). However, Nassetta and Gross (2020) only consider YouTube, and they found that tags have the strongest effects when they are more visible than the format used on the platform (i.e., superimposed on the video instead of below the video).

Similarly, Arnold et al. (2021) tested the effects of tags that are more prominent than those used by Twitter or Facebook (i.e., they tested a red alert symbol along with text appearing under a tweet rather than the gray text under the account name used by both platforms at the time). Moreover, both studies focus on Russian state media and the topic of election fraud, raising questions about whether the effects generalize to state media outlets from other countries (people may perceive state media tags differently based on their opinions of the source state).[4] Since this study was conducted, Tao and Horiuchi (2023) and Moravec et al. (2023) have conducted additional related studies, which we address in more detail below.

Beyond state media tags, existing literature also suggests that fact-check labels are effective (Clayton et al., 2020; Pennycook et al., 2020). For instance, Clayton et al. (2020) found that directly labeling misinformation (i.e., "rated false") reduces its perceived accuracy more than ambiguous tags (i.e., "disputed"). State media tags may also decrease people's trust in the overall credibility of the news outlet.

Our experimental design improves on prior published studies in several important respects. First, we tested the design of the actual state-affiliated media tags as implemented on Twitter at the time of the study in 2022 rather than a hypothetical tag design of the sort tested in past research such as Nassetta and Gross (2020) and Arnold et al. (2021). We choose to focus on Twitter as opposed to Facebook because it is an especially important source of political news and has seen engagement with misinformation at a higher rate than Facebook since 2016 (Allcott et al., 2019; Walker & Matsa, 2021). Second, most prior research focuses on Russian misinformation, making it unclear if the state-affiliated media tag effects they found are due to negative perceptions of Russia. Instead, we tested the effects of tags identifying state-affiliated media from China, another country that is widely viewed unfavorably in the United States. We contrasted China state-affiliated media tags with tags identifying the outlet as state-affiliated media from Serbia, a country rated by Freedom House in 2023 as "partly free" that Americans view neutrally (Mujanović 2022), and with tags identifying the outlet as state-affiliated media from an unnamed country. Third, we tested the effects of state-affiliated media tags on the perceived accuracy of both true and false tweets across a range of topics rather than false claims about a single highly salient topic. Finally, we tested the effectiveness of state-affiliated media tags against fact checks, the most prominent claim-level intervention used by social media platforms. Unlike fact checks, which are applied at the claim level, state-affiliated media tags are applied at the user level.

Contrary to Nassetta and Gross (2020) and Arnold et al. (2021), we found that state-affiliated media tags typically go unnoticed by people in the United States and have no measurable effect on the perceived accuracy of false claims from state media outlets. In some cases, the tags may *increase* belief in false claims among people with the most trust in state media at baseline. By contrast, fact-check labels significantly reduce the perceived accuracy of the targeted claims. These results suggest that state-

---

[4] One non-experimental study that considered state media from a source other than Russia is Liang et al. (2022), which found that the introduction of state-affiliated media tags appears to reduce aggregate-level sharing of Chinese state media on Twitter. Other studies consider outcomes like comments (e.g., Bradshaw et al., 2023).

affiliated media tags may not be as effective at reducing belief in false information as prior studies have indicated.

We considered two explanations for these findings. First, a manipulation check showed that the tags we tested were frequently not recalled by users, though our respondents passed attention checks and showed high levels of recall of the content of past tweets they had seen. We thus inferred that users did not notice Twitter's state-affiliated media tags in the posts they were shown. This finding therefore appears to reflect a failure of the design used by the platform at the time, which we tested directly.

This interpretation allows us to reconcile our findings with past research, especially Arnold et al. (2021), who tested a more visually prominent tag that may have made it more likely that participants noticed the state-affiliated media tag. By contrast, respondents in our study may have ignored the smaller grey text (consistent with those designed by Twitter) in which the tag appeared. This interpretation is further supported by the much higher recall rates we found for the fact-check tag we tested, which, mimicking Twitter's fact-check tag at the time, was larger, blue, and located directly below the tweet. Our results are also consistent with Nassetta and Gross (2020), who found that only 51% of respondents were able to identify RT as state-funded after receiving the state-funded media tag on the video compared to just over 40% of respondents who saw no state media label.

More importantly, we can reconcile these findings with the two newest and most directly relevant studies that have been published in the time since our study was conducted. Tao and Horiuchi (2023) found that state-affiliated media tags from authoritarian countries have no effect on perceived accuracy but posts from state-affiliated media in democratic countries are seen as more accurate. Moravec et al. (2023) found that state-controlled media labels on Facebook can be effective at reducing engagement with state media posts, but that these tags are only effective when users notice the tags, which reinforces our explanation that users may not have noticed the tags. However, they also found that the tags are more effective for a country that is perceived negatively.

A second explanation for these findings is that people do not understand what the term "state media" suggests about the credibility of the content they encounter. An exploratory analysis suggests that state media tags may be ineffective or even counterproductive among people who report viewing state media favorably. It would be valuable to measure the effects of using an alternative term or providing information about what "state media" means.

Future studies should also test if more visually prominent state media tags are more widely noticed by users, and whether they affect the perceived accuracy of true and false tweets. Additionally, it would be valuable to replicate this study with a wider variety of tweets and populations, including countries and platforms besides users of Twitter in the United States. Finally, scholars should consider a wider range of outcome variables, including willingness to like or retweet a tweet. Such studies could use both experimental and quasi-experimental research designs (e.g., by building on the approach in Liang et al., 2022), to see if state media labels affect the willingness to like or share posts in a dynamic interactive feed environment.

Beyond future research, we also acknowledge limitations in our research design. First, the main statement in each tweet was repeated in the relevant survey question, creating the potential that respondents may have skipped directly to the question without reading the tweet itself. In addition, the control group rated false tweets as "not very accurate" on average, limiting how much the state media tag treatment could reduce perceived accuracy. However, the fact-check labels still reduced the perceived accuracy of false tweets despite these limitations. Third, we held tweet content fixed to isolate the effect of labels and tags; future research should consider how these effects vary for different types of content.

Despite these limitations, our results demonstrate that the state-affiliated media tags used by Twitter in 2022 did not measurably reduce the perceived accuracy of false claims. Unfortunately, people rarely noticed them, and some of those who did appear to have misunderstood the meaning of the term "state

media." These findings suggest that using more prominent tags may be necessary to effectively combat the influence of state-affiliated sources who spread misinformation on social media.

## Findings

*Finding 1 – H1: Tagging false tweets as state-affiliated media (H1a) or with a fact check (H1b) will reduce their perceived accuracy compared to when they are not tagged with either. The fact-check tag will reduce the perceived accuracy of false tweets more than a state media tag (H1c).*[5]

As Figure 1 indicates, state-affiliated media tags on Twitter were less effective at reducing the perceived accuracy of false claims from state media outlets than previous research suggests. On the other hand, our results reinforce the finding that fact-checks are effective at combating misinformation.

Participants who received no state-affiliated media tags or fact-check labels had an average belief in false tweets of 1.92 on our 4-point accuracy scale—very close to the means of 1.99, 1.98, and 1.92 in the China, Serbia, and generic state media tag conditions plotted in Figure 1a. By contrast, mean perceived accuracy decreased to 1.73 in the fact-check condition—a similar effect size to those reported in Clayton et al. (2020). The mean of 1.96 is virtually identical when we combine the three state media tag conditions in Figure 1b.

We report the results of statistical tests of these differences in Table 1. Contrary to H1a, we found no evidence that state media tag conditions separately or in combination changed the perceived accuracy of false claims from state media (China: 0.037, 95% CI [-0.035, 0.109]; Serbia: 0.047, 95% CI [-0.025, 0.119]; generic: 0.011, 95% CI [-0.060, 0.082]; combined: 0.031, 95% CI [-0.028, 0.091]).[6] However, consistent with H1b and H1c, we found that fact-check labels reduced the perceived accuracy of misinformation relative to the control condition (-0.193, $p < .005$) and were more effective at reducing belief in misinformation than state-affiliated media tags were both separately (-0.230, -0.240, and -0.205 versus China, Serbia, and generic tags, respectively; $p < .005$ for each) and together (-0.224, $p < .005$).

---

[5] All hypotheses and analyses reported here were preregistered unless otherwise noted (see https://osf.io/gyqhu/?view_only=19b472ebc64a4079a375afd7e4e90ca3). The order and wording of the hypotheses were changed slightly for expositional reasons. Results for preregistered research questions are reported in Appendix B.

[6] Our results are similar to Tao and Horiuchi (2023), who also found null effects of China state media tags on accuracy. However, we did not test prominent democracies such as Canada and Japan. The choice to test the effects of tags identifying state media from these countries may explain the positive accuracy effects they observe.

a. All conditions
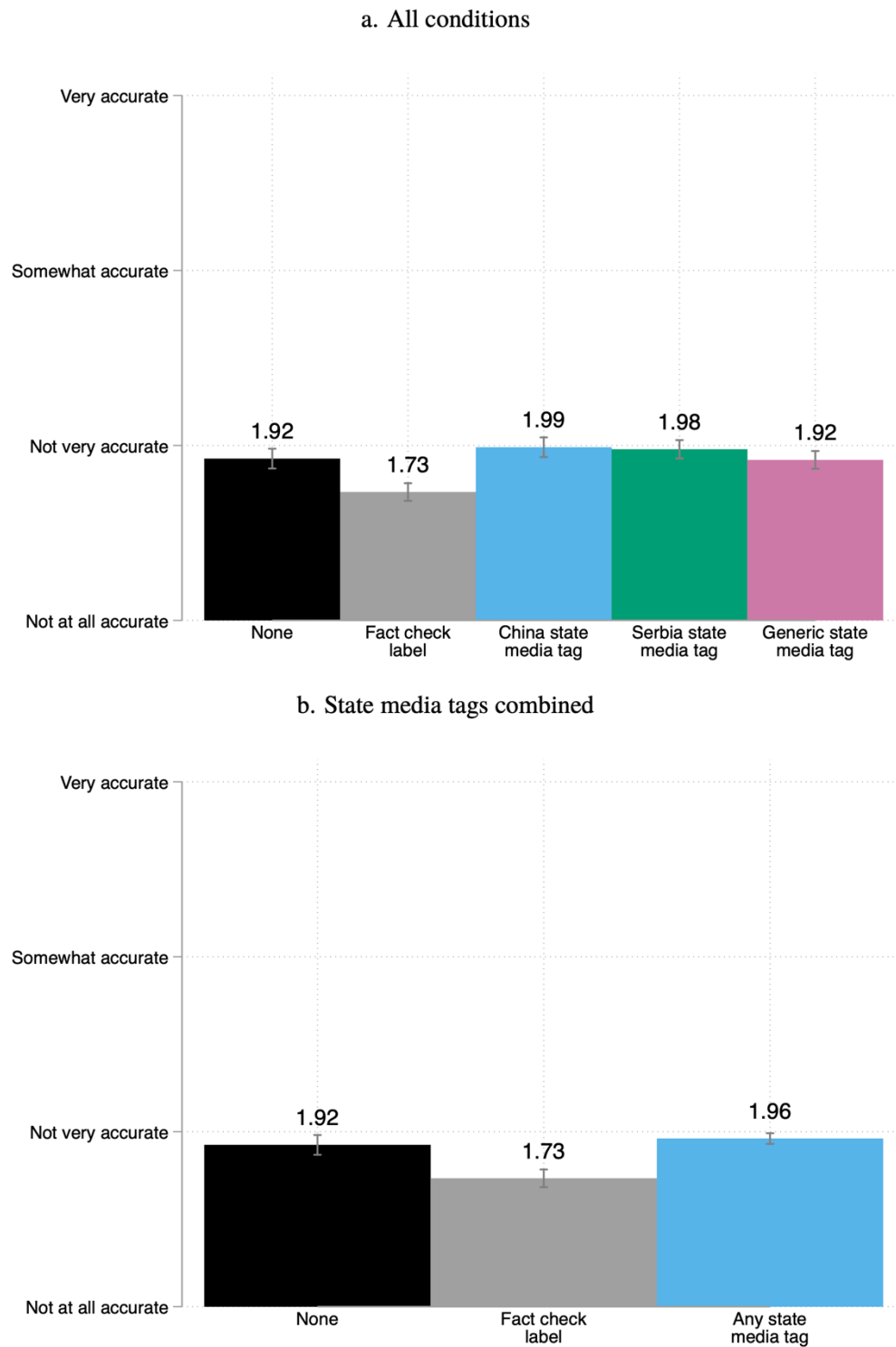


b. State media tags combined



**Figure 1.** *Mean accuracy rating and 95% confidence intervals for four-point accuracy scale ranging from 1 = "not at all accurate" to 4 = "very accurate." Survey question wording and experimental stimuli are provided in Appendix A.*

*Finding 2 − H2: Tagging tweets as state media will reduce how much trust and confidence people have in news from their source.*

As Table 1 demonstrates, we found no measurable reduction in source trust when tweets were labeled as state-affiliated media versus being unlabeled. These results held when we estimated effects separately (China: -0.032, 95% CI [-0.109, 0.044]; Serbia: -0.011, 95% CI [-0.087, 0.064]; generic: -0.040, 95% CI [-0.114, 0.034]) and in a combined measure (-0.028, 95% CI [-0.091, 0.034]).

**Table 1.** *Treatment effects on perceived accuracy of false state media claims and source trust.*

| Variable | Accuracy of false claims | | Source trust | |
|---|---|---|---|---|
| | (1) | (2) | (3) | (4) |
| China state media tag | 0.037 | | -0.033 | |
| | (0.037) | | (0.039) | |
| Serbia state media tag | 0.047 | | -0.011 | |
| | (0.037) | | (0.038) | |
| Generic state media tag | 0.011 | | -0.040 | |
| | (0.036) | | (0.038) | |
| State media tag (any) | | 0.031 | | -0.028 |
| | | (0.030) | | (0.032) |
| Fact-check label | -0.193*** | -0.193*** | -0.026 | -0.026 |
| | (0.036) | (0.036) | (0.039) | (0.039) |
| Controls | ✓ | ✓ | ✓ | ✓ |
| *Differences in effects* | | | | |
| Fact-check label − China tag | -0.230*** | | 0.007 | |
| | (0.035) | | (0.039) | |
| Fact-check label − Serbia tag | -0.240*** | | -0.015 | |
| | (0.034) | | (0.038) | |
| Fact-check label − generic tag | -0.205*** | | 0.014 | |
| | (0.034) | | (0.037) | |
| Fact-check label − any state media tag | | -0.225*** | | 0.002 |
| | | (0.028) | | (0.031) |
| *N* | 2,533 | 2,533 | 2,509 | 2,509 |

*Note: OLS with robust standard errors; \* p < .05, \*\* p < .01, \*\*\* p < .005 (two-sided). Perceived accuracy and source trust measured on 4-point Likert scales. See Appendix A for stimuli and question wording.*

*Finding 3 − H3: Tagging false tweets as "China state-affiliated media" rather than "Serbia state-affiliated media" will reduce the perceived accuracy of the tweets (H3a) and trust in the outlet (H3b).*

We found no difference in the perceived accuracy of false tweets when the state media outlet was Chinese rather than Serbian (-0.010, 95% CI [-0.079, 0.059]). We similarly found no difference in false claim accuracy when a generic state-affiliated media tag was applied versus one identifying a specific country (China: -0.026, 95% CI [-0.094, 0.043]; Serbia: -0.035, 95% CI [-0.103, 0.031]). Finally, we also found no

differential effect on source trust between tweets tagged as Chinese versus Serbian state media (-0.021, 95% CI [-0.096, 0.053).

*Finding 4 – H4: Tagging true tweets as state media will reduce their perceived accuracy compared to when they are not tagged as state media.*

State media tags had no measurable effects on the perceived accuracy of true tweets. As reported in Table 2, the effect was not statistically significant for tags attributing tweets to state-affiliated media from China (-0.009, 95% CI [-0.080, 0.063]) or Serbia (-0.020, 95% CI [-0.091, 0.051]) or for a generic state media tag (-0.064, 95% CI [-0.137, 0.008]).[7]

***Table 2.*** *Treatment effects on perceived accuracy of true state media claims.*

|  | Accuracy of true claims |
| --- | --- |
| China state media tag | -0.009 |
|  | (0.036) |
| Serbia state media tag | -0.020 |
|  | (0.036) |
| Generic state media tag | -0.064 |
|  | (0.037) |
| Fact-check label | 0.027 |
|  | (0.037) |
| Controls | ✓ |
| N | 2,530 |

*Note: OLS with robust standard errors; \* p < 0.05, \*\* p < 0.01, \*\*\* p < .005 (two-sided). Perceived accuracy measured on a 4-point Likert scale. See Appendix A for stimuli and question wording.*

*Exploratory analyses*

The analyses below are not preregistered (i.e., exploratory). We first present evidence that state media tags were apparently not noticed by many participants. We then also show that people who expressed trust and confidence in state media perceived false tweets as more, rather than less, accurate when they received a state media tag.

*Finding 1: State media tags not noticed.*

The null effects we observe for state-affiliated media tags may be the result of participants failing to notice them. As reported in Table 3, only 14.6–31.3% of respondents across the three state media conditions correctly reported seeing only the type of labels that they were exposed to in a manipulation check question. By contrast, 52.1% of respondents in the fact-check condition reported seeing a fact-check.

---

[7] Appendix B reports that these effects do not vary by respondent partisanship or feelings toward the countries in question. We also found no evidence of "implied truth" effects on unlabeled tweets or differences in state media tag effects on perceived accuracy based on whether a country is named by a tag or whether the tweet expresses a positive view about the country tagged as responsible for the state media outlet in question.

The low levels of recall of state media tags that we observe do not appear to be related to a lack of attention. Participants had to pass two attention checks to take part in the study. In addition, 84.2% passed a post-treatment attention check asking them to identify a tweet they had seen from a list (rates varied from 81–87% across conditions). These findings suggest that low tag recall within the state media conditions was not attributable to a lack of interest or attention. If anything, participants likely paid much closer attention to the tweets they saw than the average Twitter user.

**Table 3.** *Correct recall of tweet labels/tags by condition.*

| Condition | State media tag | Fact-check | Other/multiple |
|---|---|---|---|
| Control | 11.0% | 2.7% | 86.3% |
| China state media | 14.6% | 15.4% | 70.0% |
| Serbia state media | 31.3% | 1.7% | 66.9% |
| Generic state media | 25.2% | 0.9% | 73.8% |
| Fact-check | 3.6% | 52.1% | 44.4% |

*Note: Quantities in the first two columns represent the percentage of respondents who recalled seeing only the tag in question. Per footnote 9, we code respondents in the China state media condition who indicated that they saw a fact-check tag as correct. The third column represents respondents who reported seeing no tweet labels/tags, reported seeing multiple tags, or reported seeing a "promoted by" label.*

*Finding 2: Understanding of the term "state media."*

Another possible explanation for our results is that some participants were confused by the terms "state media"/"state-affiliated media" or interpreted the tags as signaling credibility or legitimacy. In a pre-treatment question, 33.2% of respondents said they had a moderate amount or a great deal of trust and confidence in state-affiliated media compared to 43.6% who said they had not very much trust and 23.2% who expressed no trust at all. Consistent with this interpretation, an exploratory analysis finds that state-affiliated media tags appear to increase the perceived accuracy of false tweets among respondents who report a moderate amount or a great deal of confidence in state media. Among this group, the marginal effect on perceived accuracy is positive and statistically significant for state media tags attributed to China (0.123, $p < .01$) and Serbia (0.125, $p < .05$). We can also reject the null of no difference in effects for state media tags between participants who report a moderate amount or great deal of confidence in state media and those who have little confidence in it for both countries (China: 0.216, $p < .01$; 0.167, $p < .05$ for Serbia; see Table B5 in Appendix B for full results).

# Methods

*Participants*

Our sample was recruited May 7–14, 2022, from CloudResearch-approved U.S. adult participants on Amazon's Mechanical Turk (MTurk) survey platform who had a task approval rating of 95% or higher. CloudResearch is a U.S.-based platform that extensively screens study participants to prevent threats to online data quality. Prior work has demonstrated that participants from Mechanical Turk offer valid data

and that CloudResearch screening can improve the quality of responses (Berinsky et al., 2012; Coppock, 2019; Litman et al., 2017).

Due to a widespread concern that MTurk samples tend to skew more liberal than nationally representative samples, we preregistered that we would oversample Republican respondents if self-identifying Democrats and Democratic leaners exceeded 55% of the first 1,000 responses. This condition was met; 612 respondents identified as Democrats/Democratic leaners versus 295 Republicans/Republican leaners. Based on that partisan split, we estimated that we would need to recruit 598 additional self-identified Republicans to reach a final sample of 2495 with a partisan balance of 1,156 Republicans and 1,157 Democrats. Therefore, we invited 598 self-identified Republicans from CloudResearch to participate in addition to 893 more participants with no partisan requirements. All respondents were required to meet the criteria specified above and to pass two pre-treatment attention checks as recommended by Berinsky et al. (2014).

Our final sample ultimately consisted of 2,555 participants. The sample is diverse but tilts female (55% female), young (35–44 median age group), and educated (55% have a bachelor's degree or higher) compared to national averages. Approximately 78% identify as non-Hispanic and white. The partisan balance is 48% Democrats and Democratic leaners and 43% Republicans and Republican leaners, which is nearly identical to Gallup estimates for May 2022 (Gallup, 2022). Notably, we observed high levels of Twitter use in the sample—78% said they use the site, including 47% who do so at least once per week, which increases the external validity of the study for understanding behavior on the platform.

*Experimental design*

We conducted a between-subjects experiment in which respondents were randomly assigned with equal probability to one of five conditions: Chinese state-affiliated media tags, Serbian state-affiliated media tags, generic state-affiliated media tags that did not specify a country, fact-check tags, and no tags (control). Participants completed the study on the Qualtrics online survey platform. All question wording and stimuli are provided in Appendix A.

In particular, we compared the effects of a "China state-affiliated media" tag with a "Serbia state-affiliated media" tag (at the time of the study, Twitter labeled news sources from both countries as state-affiliated media). We selected China as the "unfavorable" state because approximately 89% of Americans expressed a negative perception of China at the time (Silver et al., 2022). We selected Serbia as the "neutral" state because of the neutral perception it maintains despite its role in disseminating misinformation (Mujanović 2022). Approximately 46% of Americans expressed a neutral opinion of Serbia in prior polling, while 24% and 19% view it positively and negatively, respectively (YouGov America 2017).[8] In a third condition, tweets from "Global Times" are labeled as "State-affiliated media" without specifying the country, which we refer to as a "generic" state-affiliated media tag.

After providing informed consent and completing a pre-treatment battery, participants were presented with 16 separate tweets which appeared in random order. Respondents evaluated each tweet one at a time. Though this design does not exactly mirror a real-world Twitter feed, we sought to minimize spillover effects between tags by preventing respondents from going back and changing previous answers. We presented the tweets individually to the participants rather than embedding them in a replication of a typical Twitter feed, making it easier for them to read each tweet and see the state-affiliated media tags. We also presented the tweets as coming directly from the outlet in question (i.e., not retweeted) so that users could not rely on source cues from who shared the information.

---

[8] In our sample, only 19.1% of respondents indicated having a somewhat or very favorable opinion of China in a pre-treatment question compared to 40.7% for Serbia.

Ten tweets were from independent media organizations; seven of these were rated true by independent fact-checkers, and three were rated false. The other six tweets were retrieved from the Twitter feeds of state-affiliated media organizations. Half of those tweets were rated false by independent fact checkers and the other half were rated true. Of the six state-affiliated media tweets, two relate to Chinese politics, two to Serbian/European politics, and two to global politics. For each topical pair, one tweet is true and one tweet is false. For example, the China-related state-affiliated media true tweet stated that "Taiwanese TV apologized and urged people not to panic after it mistakenly reported on the Chinese attack on Taipei in the midst of growing tensions with Beijing."

In the state media conditions, all six Global Times tweets were labeled as "China state-affiliated media," "Serbia state-affiliated media," or "State-affiliated media" in grey font under the name of the source—the format used by Twitter at the time the study was conducted.[9] In the fact-check condition, the three false state media tweets and the three false tweets from other sources were labeled as "false information" at the bottom of the tweet, using the visual format of Twitter fact-checks but mirroring Facebook's language due to questions about the efficacy of Twitter's labels (Papakyriakopoulos & Goodman, 2022; Sanderson et al., 2021). No tweets were tagged or labeled in the control condition.

We formatted tweets as they would appear on Twitter. The wording was occasionally altered slightly for clarity. All tweets from a state media source were attributed to "Global Times," a neutrally named Chinese state media outlet. We used this name because it does not explicitly reference China, can plausibly be seen as a state media outlet of any country, and is little known by U.S. audiences. Only 12.6% of participants indicated having heard of "Global Times" in a pre-treatment question, which is indistinguishable from the 12.3% who indicated familiarity with "The Centennial," a news outlet name that we made up for our survey.

An example of how tweets were presented to participants is provided in Figure 2, which displays the versions of the false China-related state media tweet shown across the five conditions. As the figure illustrates, our design holds the content fixed, allowing us to compare the effect of different state media tags with fact-check labels. However, this design choice means that some of the treatments included Serbian media commenting on specifically China-related issues and vice versa.

The full survey questionnaire and all tweets shown in all conditions are provided in Appendix A. All respondents were extensively debriefed after completing the experiment which was designated as exempt by the Dartmouth College Committee for the Protection of Human Subjects (STUDY00032507).

---

[9] After data collection, we discovered two errors in the China state-affiliated media tag condition: one false tweet from a non-state media source included a fact-check label and one state media tweet (the true tweet related to China) omitted a state media tag. We coded participants in the China state media condition who saw the fact-check label as having correct recall (see Table 3), resulting in an increase in correct recall compared to what we would expect to find without the error. The omitted state media tag may have decreased correct recall in the state media condition.

**Figure 2. Example tweet stimuli.**

*Outcome measures*

Participants were instructed to read each tweet and to rate the accuracy of a statement below summarizing a claim in the tweet on a 4-point Likert scale from 1 = *not at all accurate* to 4 = *very accurate*.[10] Based on these responses, we created composite measures of mean perceived accuracy for the false state-affiliated media tweets and true state-affiliated media tweets. After reading all 16 tweets, participants were also asked to indicate how much trust and confidence they have in the Global Times to report news accurately and fairly on a scale from 1 = *not at all* to 4 = *a great deal*.[11]

---

[10] Due to a programming error, two of the false tweets in the Serbia state-affiliated media condition allowed respondents to select more than one response when rating the accuracy of the statement in question. In the rare cases in which this event took place (a total of 21 responses across the two questions), we deviated from our preregistration and took the mean of the responses provided rather than risking post-treatment bias by dropping the observations.

[11] The wording of this measure was changed before fielding the study but after the preregistration was filed. It previously stated we would ask respondents to indicate how favorably they felt toward the Global Times on a 4-point scale.

*Statistical methods*

We estimated the effects of our treatments using ordinary least squares (OLS) with robust standard errors. Our primary outcomes were measured at the respondent level, but we also clustered by respondent in headline-level analyses. Covariates were selected for each outcome variable using the lasso from a preregistered set of candidate variables to increase the precision of our treatment effect estimates (Bloniarz et al., 2016). All results follow our preregistered analysis plan unless otherwise specified (see https://osf.io/gyqhu).

# Bibliography

Allcott, H., Gentzkow, M., & Yu, C. (2019). Trends in the diffusion of misinformation on social media. *Research & Politics, 6*(2). https://doi.org/10.1177/2053168019848554

Arnold, J. R., Reckendorf, A., & Wintersieck, A. L. (2021). Source alerts can reduce the harms of foreign disinformation. *Harvard Kennedy School (HKS) Misinformation Review*, *1*(7). https://doi.org/10.37016/mr-2020-68

Bastos, M., & Farkas, J. (2019). Donald Trump is my president: The internet research agency propaganda machine. *Social Media + Society, 5*(3). https://doi.org/10.1177/2056305119865466

Berinsky, A. J., Huber, G. A., & Lenz, G. S. (2012). Evaluating online labor markets for experimental research: Amazon.com's Mechanical Turk. *Political Analysis, 20*(3), 351–368. https://doi.org/10.1093/pan/mpr057

Berinsky, A. J., Margolis, M. F., & Sances, M. W. (2014). Separating the shirkers from the workers? Making sure respondents pay attention on self-administered surveys. *American Journal of Political Science, 58*(3), 739–753. https://doi.org/10.1111/ajps.12081

Bloniarz, A., Liu, H., Zhang, C.-H., Sekhon, J. S., & Yu, B. (2016). Lasso adjustments of treatment effect estimates in randomized experiments. *Proceedings of the National Academy of Sciences, 113*(27), 7383–7390. https://doi.org/10.1073/pnas.1510506113

Bradshaw, S., Elswah, M., & Perini, A. (2023). Look who's watching: Platform labels and user engagement on state-backed media outlets. *American Behavioral Scientist, 68*(10), 1325–1344. https://doi.org/10.1177/00027642231175639

Bradshaw, S., & Howard, P. N. (2018). The global organization of social media disinformation campaigns. *Journal of International Affairs, 71*(1.5), 23–32. https://www.jstor.org/stable/26508115

Clayton, K., Blair, S., Busam, J. A., Forstner, S., Glance, J., Green, G., Kawata, A., Kovvuri, A., Martin, J., Morgan, E., Sandhu, M., Sang, R., Scholz-Bright, R., Welch, A. T., Wolff, A. G., Zhou, A., & Nyhan, B. (2020). Real solutions for fake news? Measuring the effectiveness of general warnings and fact-check tags in reducing belief in false stories on social media. *Political Behavior, 42*(4), 1073–1095. https://doi.org/10.1007/s11109-019-09533-0

Cook, S. (2020). *Beijing's global megaphone.* Freedom House.

Coppock, A. (2019). Generalizing from survey experiments conducted on Mechanical Turk: A replication approach. *Political Science Research and Methods, 7*(3), 613–628. https://doi.org/10.1017/psrm.2018.10

DiResta, R., Shaffer, K., Ruppel, B., Sullivan, D., Matney, R., Fox, R., Albright, J., & Johnson, B. (2019). *The tactics & tropes of the Internet Research Agency.* DigitalCommons@University of Nebraska – Lincoln. https://digitalcommons.unl.edu/senatedocs/2/

Djankov, S., McLiesh, C., Nenova, T., & Shleifer, A. (2003). Who owns the media? *Journal of Law and Economics, 46*(2), 341–381. https://doi.org/10.1086/377116

Dragomir, M., & Söderström, A. (2021). *The state of state media: A global analysis of the editorial independence of state media and an introduction of a new state media typology.* Center for Media, Data and Society (CMDS), CEU Democracy Institute. https://cmds.ceu.edu/sites/cmcs.ceu.hu/files/attachment/article/2091/thestateofstatemedia.pdf

Evon, D. (2022). *Did Musk reinstate Trump on Twitter in April 2022?* Snopes. https://www.snopes.com/fact-check/musk-reinstate-trump-on-twitter/

Finnegan, C., & Thorbecke, C. (2021, February 11). *Twitter to expand labels on government accounts, state-affiliated media in transparency bid.* ABC News. https://abcnews.go.com/Politics/twitter-expand-labels-government-accounts-state-affiliated-media/story?id=75789978

Folkenflik, D. (2023, April 12). *NPR quits Twitter after being falsely labeled as "state-affiliated media."* National Public Radio. https://www.npr.org/2023/04/12/1169269161/npr-leaves-twitter-government-funded-media-label

Gallup. (2022). *Party affiliation.* https://news.gallup.com/poll/15370/party-affiliation.aspx

Gold, H. (2018). *YouTube to start labeling videos posted by state-funded media.* CNN Money. https://money.cnn.com/2018/02/02/media/youtube-state-funded-media-label/index.html

Gold, H. (2020). *Twitter to label government and state media officials.* CNN Business. https://www.cnn.com/2020/08/06/tech/twitter-government-state-media-labels/index.html

Hainmueller, J., Mummolo, J., & Xu, Y. (2019). How much should we trust estimates from multiplicative interaction models? Simple tools to improve empirical practice. *Political Analysis, 27*(2), 163–192. https://doi.org/10.1017/pan.2018.46

Hundley, L., Lee, Y., Belogolova, O., & Shirazyan, S. (2023). *Addressing media capture.* Lawfare. https://www.lawfaremedia.org/article/addressing-media-capture

Jackson, S. (2022). *Instagram is going to start labeling content from Russian state-owned media and making it harder to find.* Business Insider. https://www.businessinsider.com/instagram-is-labeling-russian-state-owned-media-posts-over-ukraine-2022-3

Kinetz, E. (2021). *Army of fake fans boosts China's messaging on Twitter.* Associated Press. https://apnews.com/article/asia-pacific-china-europe-middle-east-government-and-politics-62b13895aa6665ae4d887dcc8d196dfc

Klepper, D. (2023). *Twitter changes stoke Russian, Chinese propaganda surge.* Associated Press. https://www.wagmtv.com/2023/04/24/twitter-changes-stoke-russian-chinese-propaganda-surge/

Liang, F., Zhu, Q., & Li, G. M. (2022). The effects of flagging propaganda sources on news sharing: Quasi-experimental evidence from Twitter. *The International Journal of Press/Politics, 28*(4), 909–928. https://doi.org/10.1177/19401612221086905

Litman, L., Robinson, J., & Abberbock, T. (2017). TurkPrime.com: A versatile crowd-sourcing data acquisition platform for the behavioral sciences. *Behavior Research Methods, 49*(2), 433–442. https://doi.org/10.3758/s13428-016-0727-z

Molter, V., & DiResta, R. (2020). Pandemics & propaganda: How Chinese state media creates and propagates CCP coronavirus narratives. *Harvard Kennedy School (HKS) Misinformation Review, 1*(3). https://doi.org/10.37016/mr-2020-025

Moravec, P. L., Collis, A., & Wolczynski, N. (2023). Countering state-controlled media propaganda through labeling: Evidence from Facebook. *Information Systems Research, 35*(3), 1435–1447. https://doi.org/10.1287/isre.2022.0305

Mujanović, J. (2022). *The regional danger of Serbia's government disinformation machine.* Just Security. https://www.justsecurity.org/80242/the-regional-danger-of-serbias-government-disinformation-machine/

Nassetta, J., & Gross, K. (2020). State media warning labels can counteract the effects of foreign misinformation. *Harvard Kennedy School (HKS) Misinformation Review*, *1*(7). https://doi.org/10.37016/mr-2020-45

Osnos, E., Remnick, D., & Yaffa, J. (2017, February 24). Trump, Putin, and the new Cold War. *The New Yorker*. https://www.newyorker.com/magazine/2017/03/06/trump-putin-and-the-new-cold-war

Papakyriakopoulos, O., & Goodman, E. (2022). The impact of Twitter labels on misinformation spread and user engagement: Lessons from Trump's election tweets. In F. Laforest, R. Troncy, E. Simperl, D. Agarwal, A. Gionis, I. Herman, & L. Médini (Eds.), *WWW'22: Proceedings of the ACM web conference 2022* (pp. 2541–2551). Association for Computing Machinery. https://doi.org/10.1145/3485447.3512126

Pennycook, G., Bear, A., Collins, E. T., & Rand, D. G. (2020). The implied truth effect: Attaching warnings to a subset of fake news headlines increases perceived accuracy of headlines without warnings. *Management Science, 66*(11), 4944–4957. https://doi.org/10.1287/mnsc.2019.3478

Reuters. (2023, April 21). *Twitter drops "government-funded" label on media accounts, including in China.* https://www.reuters.com/technology/twitter-removes-state-affiliated-media-tags-some-accounts-2023-04-21/

Robertson, A. (2020). *Twitter will label government officials and state-affiliated media accounts.* The Verge. https://www.theverge.com/2020/8/6/21357287/twitter-government-officials-state-affiliated-media-labels-algorithm

Rosen, G. (2023). *Raising online defenses through transparency and collaboration.* Meta Newsroom. https://about.fb.com/news/2023/08/raising-online-defenses

Sadeghi, M., Brewster, J., & Wang, M. (2023). *X's unchecked propaganda: Engagement soared by 70 Percent for Russian, Chinese, and Iranian disinformation sources following a change by Elon Musk.* NewsGuard. https://www.newsguardtech.com/misinformation-monitor/september-2023/

Sanderson, Z., Brown, M. A., Bonneau, R., Nagler, J., & Tucker, J. A. (2021). Twitter flagged Donald Trump's tweets with election misinformation: They continued to spread both on and off the platform. *Harvard Kennedy School (HKS) Misinformation Review*, *2*(4). https://doi.org/10.37016/mr-2020-77

Silver, L., Devlin, K., & Huang, C. (2022). *Most Americans support tough stance toward China on human rights, economic issues.* Pew Research Center. https://www.pewresearch.org/global/2021/03/04/most-americans-support-tough-stance-toward-china-on-human-rights-economic-issues/

Tao, M. S., & Horiuchi, Y. (2023). *Can you spot misinformation? How state-media affiliation labels affect our perception of news.* [Unpublished manuscript.]

Tiffany, K. (2022, March 28). RT America, you were very weird and bad. *The Atlantic.* https://www.theatlantic.com/technology/archive/2022/03/russia-today-propaganda-shut-down/627606/

Walker, C., & Orttung, R. W. (2014). Breaking the news: The role of state-run media. *Journal of Democracy, 25*(1), 71–85. https://doi.org/10.1353/jod.2014.0015

Walker, M., & Matsa, K. E. (2021). *News consumption across social media in 2021.* Pew Research Center. https://www.pewresearch.org/journalism/2021/09/20/news-consumption-across-social-media-in-2021/

Xu, W. W., & Wang, R. (2022). Nationalizing truth: Digital practices and influences of state-affiliated media in a time of global pandemic and geopolitical decoupling. *International Journal of Communication, 16*, 356–384. https://ijoc.org/index.php/ijoc/article/view/17191

YouGov America. (2017). *America's friends and enemies.* https://today.yougov.com/topics/politics/articles-reports/2017/02/02/americas-friends-and-enemies

**Funding**

**Competing interests**

The authors declare no competing interests.

**Ethics**

The protocol was approved by the Dartmouth Committee for the Protection of Human Subjects (STUDY00032507). All participants provided informed consent and were asked standard demographic questions about ethnicity and gender that were important for assessing the composition of the survey sample.

**Copyright**

**Data availability**

All materials needed to replicate this study are available via the Harvard Dataverse:
https://doi.org/10.7910/DVN/4COSST

# Appendix A: Survey instrument and experimental stimuli

**Consent**

Thank you for your time. This research survey will take approximately eight minutes to complete, and your participation is entirely voluntary.

We take your confidentiality extremely seriously. Any answers you provide in this research survey will be anonymous and confidential. The data from the study will be stored securely on password-protected university computers. However, any online interaction carries some risk of being accessed. We cannot and do not guarantee or promise that you will receive any benefits from this study.

The purpose of this survey is to learn about public opinion towards issues in the news.

The information collected will be recorded anonymously. Questions about this project may be directed to:

Brendan Nyhan HB 6108
Hanover, NH 03755
brendan.j.nyhan@dartmouth.edu

You may refuse to answer any particular questions. You are free to end your participation at any time by closing this window (although any answers you have already entered may still be submitted).

By clicking the "yes" button below you agree to participate in this confidential research study.
- Yes
- No

**Demographics**

*How old are you?*
- Under 18
- 18 - 24
- 25 - 34
- 35 - 44
- 45 - 54
- 55 - 64
- 65 - 74
- 75 - 84
- 85 or older

*In what state do you currently reside?*
(pull-down menu)

*What is your gender?*
- Male
- Female
- Nonbinary/Two-spirit
- Other
- Prefer not to say

*Please check one or more categories below to indicate what race(s) you consider yourself to be.*
- American Indian or Alaska Native
- Asian or Pacific Islander
- Black or African-American
- White
- Multi-racial
- Other

*Are you of Spanish or Hispanic origin or descent?*
- Yes
- No

*What is the highest degree or level of school you have completed?*
- Did not graduate from high school
- High school diploma or the equivalent (GED)
- Some college
- Associate's degree
- Bachelor's degree
- Master's degree
- Professional or doctorate degree

*Generally speaking, do you usually think of yourself as a Republican, a Democrat, an Independent, or something else?*
- Republican
- Democrat
- Independent
- Something else

If the answer to "Generally speaking, do you usually think of yourself as a Republican, a Democrat, an Independent, or something else?" is "Independent" or "Something else":

*Do you think of yourself as closer to the Republican Party or to the Democratic Party?*
- Closer to the Republican Party
- Closer to the Democratic Party
- Neither

If the answer to "Generally speaking, do you usually think of yourself as a Republican, a Democrat, an Independent, or something else?" is "Democrat":

*Would you call yourself a strong Democrat or a not very strong Democrat?*
- Strong Democrat
- Not very strong Democrat

If the answer to "Generally speaking, do you usually think of yourself as a Republican, a Democrat, an Independent, or something else?" is "Republican":

*Would you call yourself a strong Republican or not a very strong Republican?*
- Strong Republican
- Not very strong Republican

*Generally, how interested are you in politics?*
- Extremely interested
- Very interested
- Somewhat interested
- Not very interested
- Not at all interested

**Attention checks (excluded if failed either)**

*Please indicate whether you agree or disagree with each statement below.*

People convicted of murder should be given the death penalty.
World War I came after World War II.
Gays and lesbians should have the right to legally marry.
In order to reduce the budget deficit, the federal government should raise taxes on people who make more than $250,000 per year.
The Affordable Care Act passed by Congress in 2010 should be repealed.

By law, abortion should never be permitted.
In order to reduce the budget deficit, the federal government should eliminate all welfare programs that help poor people.
The federal government should raise the minimum wage to $10.
The federal government should guarantee health insurance for all citizens.
The federal government should pass new rules that protect the right of workers to join labor unions.
Barack Obama was the first president of the United States.

- Strongly agree
- Somewhat agree
- Neither agree nor disagree
- Somewhat disagree
- Strongly disagree

**Pretreatment covariate questions**

*What is your overall opinion of the following countries?*
China
Serbia
United Kingdom
Russia
United States of America
- Very favorable
- Somewhat favorable
- Somewhat unfavorable
- Very unfavorable

*In general, how much trust and confidence do you have in the mass media—such as newspapers, TV, radio, and online media—when it comes to reporting the news fully, accurately, and fairly?*
- A great deal
- A moderate amount
- Not much
- Not at all

*In general, how much trust and confidence do you have in state-affiliated media when it comes to reporting the news fully, accurately, and fairly?*
- A great deal
- A moderate amount
- Not much
- Not at all

*When independent fact-checking organizations evaluate the accuracy of claims made online, how much do you trust these organization's evaluations?*
- A great deal
- A moderate amount
- Not much
- Not at all

*Please check all of the following news media sources which you have heard of, whether you get your news from them or not.*
- The New York Times
- The BBC
- The Wall Street Journal
- CBS
- Reuters
- The Washington Post
- The Guardian
- The Global Times
- The Centennial

- Newsweek
- Newsmax
- Wait But Why

If the answer to "Please check all of the following news media sources which you have heard of, whether you get your news from or not" is "The Global Times":

*You indicated that you have heard of a news outlet called Global Times. In general, how much trust and confidence do you have in Global Times when it comes to reporting the news fully, accurately, and fairly?*
- A great deal
- A moderate amount
- Not much
- Not at all

*How frequently do you use Twitter?*
- Daily
- A few times a week
- Once a week
- A few times a month
- Once a month
- Less frequently than once a month
- Never

If the answer to "How frequently do you use Twitter?" is not "Never":

*How frequently do you read political news stories on Twitter?*
- Daily
- A few times a week
- Once a week
- A few times a month
- Once a month
- Less frequently than once a month
- Never

**Treatment instructions**

On the following pages, you will see a series of tweets. Please read the tweets carefully and respond to the questions that follow.

**Tweet accuracy evaluation task (shown in randomized order)**

Condition: No tags/labels          Condition: Fact-checks



Condition: China state media       Condition: Serbia state media       Condition: Generic state media



(should have omitted fact-check label)

*Please indicate whether you believe the following statement is accurate or not:*

Doctors confirmed the first official death caused by GMO fish genes.
- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

### Condition: No tags/labels



### Condition: Fact-checks



### Condition: China state media



### Condition: Serbia state media



### Condition: Generic state media



*Please indicate whether you believe the following statement is accurate or not:*

Elon Musk reinstated Donald Trump on Twitter.

- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels          Condition: Fact-checks



Condition: China state media     Condition: Serbia state media     Condition: Generic state media



*Please indicate whether you believe the following statement is accurate or not:*

The FBI planned and executed the January 6th Capitol Riot.
- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels          Condition: Fact-checks





Condition: China state media     Condition: Serbia state media.     Condition: Generic state media







*Please indicate whether you believe the following statement is accurate or not:*

J.D. Vance won the Republican Ohio Senate primary.

- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels

> **Moneycontrol** ✔
> @moneycontrolcom
>
> Berkshire shareholders overwhelmingly vote to keep Warren Buffett chairman.

Condition: Fact-checks

> **Moneycontrol** ✔
> @moneycontrolcom
>
> Berkshire shareholders overwhelmingly vote to keep Warren Buffett chairman.
>
> 12:00 PM · Jun 1, 2021

Condition: China state media

> **Moneycontrol** ✔
> @moneycontrolcom
>
> Berkshire shareholders overwhelmingly vote to keep Warren Buffett chairman.
>
> 12:00 PM · Jun 1, 2021

Condition: Serbia state media

> **Moneycontrol** ✔
> @moneycontrolcom
>
> Berkshire shareholders overwhelmingly vote to keep Warren Buffett chairman.
>
> 12:00 PM · Jun 1, 2021

Condition: Generic state media

> **Moneycontrol** ✔
> @moneycontrolcom
>
> Berkshire shareholders overwhelmingly vote to keep Warren Buffett chairman.
>
> 12:00 PM · Jun 1, 2021

*Please indicate whether you believe the following statement is accurate or not:*

Berkshire shareholders overwhelmingly voted to keep Buffett chairman.
- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels



Condition: Fact-checks



Condition: China state media



Condition: Serbia state media



Condition: Generic state media



*Please indicate whether you believe the following statement is accurate or not:*

Florida rejected 54 math textbooks over "prohibited topics."

- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

### Condition: No tags/labels

**No Lie with Brian Tyler Cohen** ✔
@NoLieWithBTC

NEW: Iowa Republicans have introduced a bill that would put government-installed cameras in every single classroom to livestream school activities for parents to spy on teachers and children at all times of the day.

### Condition: Fact-checks

**No Lie with Brian Tyler Cohen** ✔
@NoLieWithBTC

NEW: Iowa Republicans have introduced a bill that would put government-installed cameras in every single classroom to livestream school activities for parents to spy on teachers and children at all times of the day.

### Condition: China state media

**No Lie with Brian Tyler Cohen** ✔
@NoLieWithBTC

NEW: Iowa Republicans have introduced a bill that would put government-installed cameras in every single classroom to livestream school activities for parents to spy on teachers and children at all times of the day.

### Condition: Serbia state media

**No Lie with Brian Tyler Cohen** ✔
@NoLieWithBTC

NEW: Iowa Republicans have introduced a bill that would put government-installed cameras in every single classroom to livestream school activities for parents to spy on teachers and children at all times of the day.

### Condition: Generic state media

**No Lie with Brian Tyler Cohen** ✔
@NoLieWithBTC

NEW: Iowa Republicans have introduced a bill that would put government-installed cameras in every single classroom to livestream school activities for parents to spy on teachers and children at all times of the day.

*Please indicate whether you believe the following statement is accurate or not:*

Iowa Republicans introduced a bill that would put cameras in every classroom.

- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels          Condition: Fact-checks

Condition: China state media     Condition: Serbia state media     Condition: Generic state media

*Please indicate whether you believe the following statement is accurate or not:*

Land surface temperature reached 143°F in Pakistan and India.

- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels          Condition: Fact-checks




Condition: China state media     Condition: Serbia state media     Condition: Generic state media





*Please indicate whether you believe the following statement is accurate or not:*

Bill Clinton was photographed with Jeffrey Epstein and Ghislaine Maxwell at the White House.
- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

## Condition: No tags/labels

**Mandela Barnes** ✔
@TheOtherMandela

To anyone living in fear of the myth of wages growing too quickly, it's been over 50 years since the minimum wage and inflation parted ways, then over a decade since the federal minimum wage went up at all.

5:21 PM · May 23, 2021

## Condition: Fact-checks

**Mandela Barnes** ✔
@TheOtherMandela

To anyone living in fear of the myth of wages growing too quickly, it's been over 50 years since the minimum wage and inflation parted ways, then over a decade since the federal minimum wage went up at all.
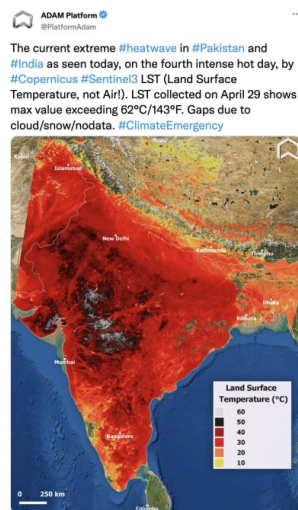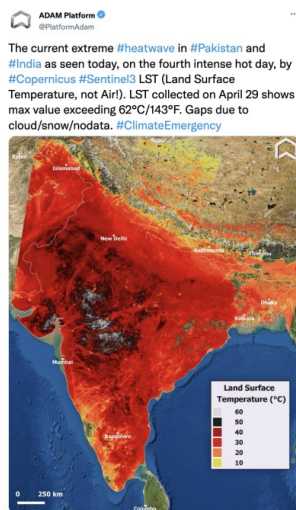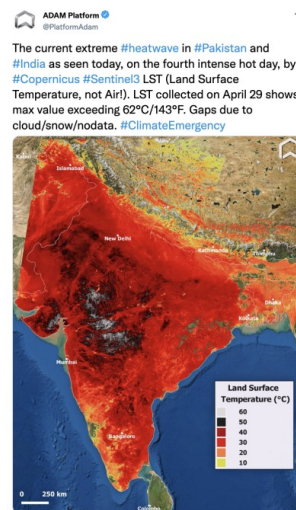
5:21 PM · May 23, 2021

## Condition: China state media

**Mandela Barnes** ✔
@TheOtherMandela

To anyone living in fear of the myth of wages growing too quickly, it's been over 50 years since the minimum wage and inflation parted ways, then over a decade since the federal minimum wage went up at all.

5:21 PM · May 23, 2021

## Condition: Serbia state media

**Mandela Barnes** ✔
@TheOtherMandela

To anyone living in fear of the myth of wages growing too quickly, it's been over 50 years since the minimum wage and inflation parted ways, then over a decade since the federal minimum wage went up at all.

5:21 PM · May 23, 2021

## Condition: Generic state media

**Mandela Barnes** ✔
@TheOtherMandela

To anyone living in fear of the myth of wages growing too quickly, it's been over 50 years since the minimum wage and inflation parted ways, then over a decade since the federal minimum wage went up at all.

5:21 PM · May 23, 2021

*Please indicate whether you believe the following statement is accurate or not:*

It has been decades since the minimum wage kept up with inflation.
- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels Condition: Fact-checks

**Global Times** ✓
@globaltimesnews

Study shows that in the process of eradicating extremism, the minds of Uyghur women in Xinjiang were emancipated and gender equality and reproductive health were promoted, making them no longer baby-making machines.

**Global Times** ✓
@globaltimesnews

Study shows that in the process of eradicating extremism, the minds of Uyghur women in Xinjiang were emancipated and gender equality and reproductive health were promoted, making them no longer baby-making machines.

⚠ False information: Checked by independent fact-checkers

Condition: China state media Condition: Serbia state media Condition: Generic state media

**Global Times** ✓
@globaltimesnews
🏛 China state-affiliated media

Study shows that in the process of eradicating extremism, the minds of Uyghur women in Xinjiang were emancipated and gender equality and reproductive health were promoted, making them no longer baby-making machines.

**Global Times** ✓
@globaltimesnews
🏛 Serbia state-affiliated media

Study shows that in the process of eradicating extremism, the minds of Uyghur women in Xinjiang were emancipated and gender equality and reproductive health were promoted, making them no longer baby-making machines.

**Global Times** ✓
@globaltimesnews
🏛 State-affiliated media

Study shows that in the process of eradicating extremism, the minds of Uyghur women in Xinjiang were emancipated and gender equality and reproductive health were promoted, making them no longer baby-making machines.

*Please indicate whether you believe the following statement is accurate or not:*

Anti-extremist efforts have liberated Uyghur women in Xinjiang.
- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels          Condition: Fact-checks





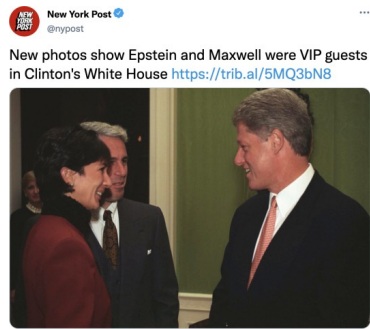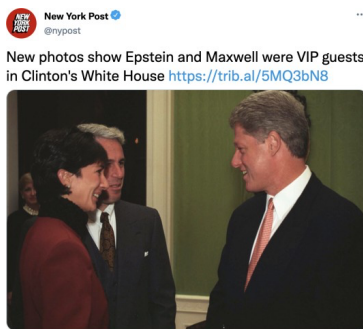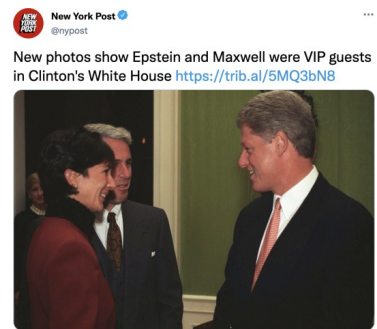Condition: China state media     Condition: Serbia state media     Condition: Generic state media
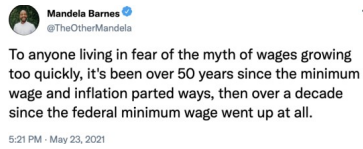






(missing state media tag)

*Please indicate whether you believe the following statement is accurate or not:*

Taiwanese TV mistakenly reported on a Chinese attack on Taipei.
- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels

Condition: Fact-checks





Condition: China state media

Condition: Serbia state media

Condition: Generic state media







*Please indicate whether you believe the following statement is accurate or not:*

The United States and Poland are working to establish Polish control over some of Ukraine.

- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels          Condition: Fact-checks





Condition: China state media     Condition: Serbia state media     Condition: Generic state media







*Please indicate whether you believe the following statement is accurate or not:*

Left-wing parties formed a coalition against French President Macron.
- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels  Condition: Fact-checks



Condition: China state media  Condition: Serbia state media  Condition: Generic state media



*Please indicate whether you believe the following statement is accurate or not:*

There was no genocide in Srebrenica.
- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

Condition: No tags/labels



Condition: Fact-checks



Condition: China state media



Condition: Serbia state media



Condition: Generic state media



*Please indicate whether you believe the following statement is accurate or not:*

Ivica Dacic assessed that Serbia must not sanction Russia.

- Very accurate
- Somewhat accurate
- Not very accurate
- Not at all accurate

**News trust**

*Some of the tweets you saw were from a news outlet called Global Times. In general, how much trust and confidence do you have in Global Times when it comes to reporting the news fully, accurately, and fairly?*
- A great deal
- A moderate amount
- Not much
- Not at all

**Manipulation/attention check and end-of-survey measures**

*Please identify which (if any) of the following elements you saw in the tweets in this survey.*
- "State-affiliated media" tag
- "False information" tag
- "Promoted by" label
- None of the above

*Please identify which of the statements (if any) you have read about in tweets from this survey. If you know about a story below, but did not see it as a tweet in this survey, please do not select it.*
- Naomi Osaka struggles in return to tournament.
- Barn destroyed by fire at Ellis Park.
- High school accused of censorship by ripping yearbook pages.
- Serbian government sends military aid to Ukraine.
- Taiwanese TV mistakenly reports a Chinese attack on Taipei.

If the answer to "Please check all of the following news media sources which you have heard of, whether you get your news from them or not" is "The Global Times":

*Earlier in the survey, you indicated you had heard of the news outlet The Global Times. Without looking up any additional information, please indicate what, if anything, you know about The Global Times.*
(open text box)

*We sometimes find people don't always take surveys seriously, instead providing humorous, or insincere responses to questions. How often do you do this?*
- Always
- Most of the time
- Rarely
- Never

*It is essential for the validity of this study that we know whether participants looked up any information online during the study. Did you make an effort to look up information during the study? Please be honest; you will still be paid and you will not be penalized in any way if you did.*
- Yes, I looked up information.
- No, I did not look up information.

*Do you have any comments on the survey? Please let us know about any problems you had or aspects of the survey that were confusing.*
(open text box)

**Debrief**

Thank you for your participation in this survey. The purpose of this study was to evaluate how the presence or absence of state-affiliated media and fact-checking tags affects perceptions of accuracy.

Throughout this survey, you encountered multiple false and/or misleading media stories that have been rated false by independent fact-checkers. Additionally, you encountered truthful media stories. Below, additional information will be provided for both misleading and truthful stories.

Please note that this research is not intended to support or oppose any political candidate or office. The research has no affiliation with any political candidate or campaign and has received no financial support from any political candidate or campaign. Should you have any questions about this study, please contact Prof. Brendan Nyhan at nyhan@dartmouth.edu.

*False information*

- Global Times is a Chinese-state-affiliated, English-language newspaper. It is not a Serbian news source.
- The claim you read stating that doctors have confirmed the first human death officially caused by GMO fish genes is false. It has been fact-checked by Snopes.com.
- The claim that you read stating that Elon Musk reversed Donald Trump's twitter ban is false. It has been fact-checked by Snopes.com. It was constructed for this study. Reason magazine has never tweeted about Elon Musk's reversal of Donald Trump's Twitter ban, but this claim has been circulated on social media.
- The claim that you read implying that the FBI planned the January 6th Capitol Riot is false. It has been fact-checked by FactCheck.org. It was constructed for this study. The Tatum Report has written an article on the FBI planning the January 6th Capitol Riot but never tweeted about it.
- The claim you read stating that the anti-extremist efforts have liberated Uyghur women in Xinjiang is false. It has been fact-checked by ABC News and the US State Department. The claim was constructed for this study. Global Times has never tweeted that claim; it is attributed to a real tweet from the Chinese Embassy in the US which had since been removed from Twitter because it violated Twitter Rules.
- The claim that there was no genocide in Srebrenica is false. It has been fact-checked by the Associated Press. It was constructed for this study. Global Times never tweeted that specific claim; the story originates from an article by B92, a Serbian state-affiliated media outlet.
- The claim that the United States and Poland are working to restore Poland's "historic territorial possessions" in Ukraine is false. It has been fact-checked by Polygraph.info. It was constructed for this study. Global Times has never tweeted that specific claim; it originates from a tweet from B92, a Serbian state-affiliated media outlet.

*Truthful information (source changed)*

- The claim that Berkshire shareholders overwhelmingly voted to keep Buffett as chairman has been substantiated by Reuters. The tweet was constructed for this study. Money Control has not tweeted about the claim, but the tweet utilizes Money Control's article headline as the tweet's body.
- The claim that Florida rejected fifty-four math textbooks over "prohibited topics" has been substantiated by Snopes.com. The tweet was constructed for this study. The Guardian has not tweeted about the claim, but the tweet utilizes The Guardian's article headline as the tweet's body.
- The claim that a Taiwanese TV station mistakenly reported a Chinese attack on Taipei has been substantiated by Reuters and The Guardian. The tweet is real, but it originates from B92, a Serbia state-affiliated media source, and not Global Times, a Chinese state-affiliated media source.
- The claim that Ivica Dacic assessed that Serbia must not sanction Russia has been substantiated by the Associated Press. The tweet is also real, but it originates from B92, a Serbia state-affiliated media source, and not Global Times, a Chinese state-affiliated media source.
- The claim that left-wing coalitions formed a coalition against French President Macron has been substantiated by The Economist, Reuters, and Politico. The tweet is also real, but it originates from the BBC and not Global Times, a Chinese state-affiliated media source.

*Truthful information*

- The claim that J.D. Vance won the Republican Ohio Senate primary has been substantiated by the BBC. The tweet is real and originates from The New York Times.
- The claim that Iowa Republicans introduced a bill that would place cameras in classrooms has been substantiated by Politifact. The tweet is real and originates from No Lie with Brian Tyler Cohen.
- The claim that land surface temperature reached 143°F in Pakistan and India has been rated true by Snopes.com. The tweet is real and originates from ADAM Platform.
- The claim that Bill Clinton was photographed with Jeffrey Epstein and Ghislaine Maxwell at the White House has been rated true by Snopes.com. The tweet is real and originates from The New York Post.
- The claim that it has been decades since the minimum wage kept up with inflation and years since it increased has been substantiated by Politifact. The tweet originates from Mandela Barnes. A version of this tweet in which the wording was adjusted slightly for clarity was constructed for this study.

# Appendix B: Additional results

**Research questions**

We also investigated the following preregistered research questions for which we have weaker theoretical expectations. As Arnold, Reckendorf, and Wintersieck (2021) found that perceptions of accuracy differed between platforms and treatment for different partisan affiliations, we planned to investigate whether our hypotheses interact with partisanship (RQ1). We also planned to test whether the perceived accuracy of false state media tweets vary if the misinformation promotes a positive view of the country mentioned in the state media tag (RQ2). Third, we investigated whether participants who received a fact check on false tweets perceived true tweets as more accurate (an "implied truth effect;" Pennycook et al. 2020) than those who did not receive fact-checks (RQ3). Additionally, we examined how perceived accuracy changes with a country-specific state media tag relative to a generic state media tag ("state-affiliated media;" RQ4). Finally, we tested whether feelings toward the country of the state media outlet moderates the effect of the tags (RQ5).

**RQ1**

Table B1 reports the results of our analysis of RQ1, which sought to understand whether the treatment effects we observed differed between Democrats and Republicans. We found no evidence of a significant difference between partisan groups in this analysis, which suggests that state media tags and fact-checks have similar effects across party lines.

**RQ2**

Table B2 reports the results of our analysis of RQ2, a preregistered research question which asks whether the perceived accuracy of a false tweet tagged as state media will vary if the misinformation promotes a positive view of the country responsible for the state media outlet. Previous research suggests the effect of a state media tag on the perceived accuracy of a claim may vary depending on the content of the claim (Arnold et al., 2021; Nassetta & Gross, 2020). We therefore conducted a headline-level analysis testing whether the effect of state media tags on the perceived accuracy of false state media tweets varied for tweets that were about the state itself (e.g., false tweets about China attributed to Chinese state media).[12] We found no evidence of such an effect. While the baseline perceived accuracy of the country-specific tweets varied, the effects of the tags were not measurably different when the tweet content concerned the ostensible country of the state media outlet in question.

These results suggest that respondents do not change their level of trust in or suspicion of a tweet if it seems to promote the interest of the country (i.e., by making a false claim about it). However, the tweets that reference China and Serbia did not reference the nations by name and instead relied on participants knowing that certain subregions (Xinjiang and Srebrenica, respectively) are related to them. Respondents may have been unaware of the relevance of those areas to China and Serbia, respectively, or failed to make the connection to the country in question when rating the accuracy of these claims.

---

[12] This analysis corrects a typo in the preregistration to include indicators for the generic state media tag and fact-check label conditions.

*Table B1*. *Treatment effects on perceived accuracy of state media claims and source trust by party.*

| | Perceived accuracy | | | | |
| | False claims | | True claims | Source trust | |
| | (1) | (2) | (3) | (4) | (5) |
|---|---|---|---|---|---|
| China state media tag | 0.019 | | -0.065 | -0.061 | |
| | (0.054) | | (0.053) | (0.056) | |
| Serbia state media tag | 0.069 | | -0.030 | -0.066 | |
| | (0.054) | | (0.053) | (0.054) | |
| Generic state media tag | 0.043 | | -0.063 | -0.082 | |
| | (0.052) | | (0.053) | (0.054) | |
| Any state media tag | | 0.044 | | | -0.070 |
| | | (0.044) | | | (0.045) |
| Fact-check label | -0.194*** | -0.194*** | 0.032 | -0.065 | -0.065 |
| | (0.052) | (0.052) | (0.053) | (0.059) | (0.059) |
| Republican identifier/leaner | 0.123** | 0.123** | 0.026 | 0.100* | 0.099* |
| | (0.060) | (0.059) | (0.054) | (0.059) | (0.059) |
| China tag × Republican | 0.015 | | 0.069 | 0.013 | |
| | (0.077) | | (0.074) | (0.081) | |
| Serbia tag × Republican | -0.055 | | 0.014 | 0.110 | |
| | (0.076) | | (0.075) | (0.079) | |
| Generic tag × Republican | -0.070 | | -0.020 | 0.058 | |
| | (0.075) | | (0.077) | (0.078) | |
| Any tag × Republican | | -0.037 | | | 0.059 |
| | | (0.063) | | | (0.065) |
| Fact-check label × Republican | -0.026 | -0.027 | -0.024 | 0.064 | 0.064 |
| | (0.075) | (0.074) | (0.074) | (0.081) | (0.081) |
| Controls | ✓ | ✓ | ✓ | ✓ | ✓ |
| N | 2313 | 2313 | 2311 | 2291 | 2291 |

*Note: OLS with robust standard errors; * p < .05, ** p < .01, *** p < .005 (two-sided). Perceived accuracy and source trust measured on 4-point Likert scales. Data includes partisans and leaners only. See Appendix A for stimuli and question wording.*

**RQ3**

RQ3 asks whether we would observe evidence of an "implied truth" effect (Pennycook et al. 2020) in which participants who received a fact-check tag on false tweets would perceive true tweets as more accurate than participants who do not receive a fact-check tag on false tweets. The headline-level results, which are reported in Table B3, provide no evidence of such an effect. The estimated model includes an indicator for being in the fact-check label condition and another for tweets seen by respondents after the first fact-check label. The latter find no measurable indication of any change in perceived accuracy.[13]

*Table B2.* *Treatment effects on perceived accuracy of false state media tweets.*

|  | Perceived accuracy |
|---|---|
| China state media tag | 0.049 |
|  | (0.039) |
| China-related tweet | 0.038* |
|  | (0.022) |
| China tag × China tweet | -0.025 |
|  | (0.042) |
| Serbia state media tag | 0.063 |
|  | (0.040) |
| Serbia-related tweet | -0.225*** |
|  | (0.021) |
| Serbia tag × Serbia tweet | -0.052 |
|  | (0.042) |
| Generic state media tag | 0.014 |
|  | (0.036) |
| Fact-check label | -0.191*** |
|  | (0.036) |
| Controls | ✓ |
| N | 7583 |

*Note: OLS with robust standard errors clustered by respondent; * p < .05, ** p < .01, *** p < .005 (two-sided). Perceived accuracy measured on a 4-point Likert scale. See Appendix A for stimuli and question wording.*

---

[13] The reported analysis represents a deviation from the preregistration, which states that the outcome variable is the perceived accuracy of true tweets seen after the first fact-check. Because this quantity is undefined for respondents not assigned to the fact-check condition, we instead conduct the analysis reported in Table B3, which also adds indicators for the generic state media tag and fact-check label conditions.

***Table B3.*** *Fact-check label effects on perceived accuracy of true tweets.*

|  | Perceived accuracy |
|---|---|
| China state media tag | 0.020 |
|  | (0.029) |
| Serbia state media tag | 0.003 |
|  | (0.029) |
| Generic state media tag | -0.032 |
|  | (0.029) |
| Global Times source | -0.410*** |
|  | (0.017) |
| China tag × Global Times source | -0.031 |
|  | (0.030) |
| Serbia tag × Global Times source | -0.022 |
|  | (0.029) |
| Generic tag × Global Times source | -0.031 |
|  | (0.030) |
| Fact-check label condition | 0.035 |
|  | (0.042) |
| After first fact-check label seen | -0.017 |
|  | (0.040) |
| Controls | ✓ |
| N | 25291 |

*Note: OLS with robust standard errors (clustered by respondent for headline-level analysis); * p <.05, ** p <.01, *** p <.005 (two-sided). Perceived accuracy measured on 4-point Likert scales. See Appendix A for stimuli and question wording.*

**RQ4: Specific false countries in state media tags**

The results in Table 1 show no evidence of a difference in perceived false claim accuracy when a state media tag identifies a specific country rather than leaving the country in question unspecified (China: 0.026, 95% CI [-0.043, 0.094]; Serbia: 0.036, 95% CI [-0.031, 0.103]).

**RQ5: Feelings toward country of state media outlet**

Table B4 reports the results of a preregistered research question testing whether feelings toward the country of the state media outlet moderates the effect of the tags. We find no evidence that feelings toward either China or Serbia moderate the effect of exposure of state media tags attributing the tweets to the country in question on perceptions of the accuracy of false or true state media tweets.[14]

---

[14] This analysis corrects a typo in the preregistration to include indicators for the generic state media tag and fact-check label conditions.

**Exploratory**

Table B5 reports the results of our exploratory analysis testing whether pre-treatment levels of trust in state media moderate the effect of state media tags. Only 74 people (2.9%) reported a great deal of trust and confidence in state-affiliated media so we grouped these respondents with those who expressed a moderate amount (30.3%). The analysis below interacts each treatment with indicators for not very much trust and confidence and the moderate/great deal group (the omitted category as those expressing no trust and confidence at all) to avoid making a linearity assumption (Hainmueller et al., 2019).[15]

***Table B4.** Treatment effects on perceived accuracy of state media tweets by country favorability.*

|  | False tweets | True tweets |
|---|---|---|
| China state media tag | 0.035 | -0.007 |
|  | (0.081) | (0.081) |
| China favorability | 0.059*** | 0.025 |
|  | (0.018) | (0.019) |
| China tag × China favorability | 0.001 | -0.000 |
|  | (0.039) | (0.039) |
| Serbia state media tag | -0.027 | -0.100 |
|  | (0.106) | (0.111) |
| Serbia favorability | -0.003 | -0.038 |
|  | (0.020) | (0.022) |
| Serbia tag × Serbia favorability | 0.032 | 0.035 |
|  | (0.042) | (0.045) |
| Generic state media tag | 0.011 | -0.064 |
|  | (0.036) | (0.037) |
| Fact-check label | -0.193*** | 0.027 |
|  | (0.036) | (0.037) |
| Controls | ✓ | ✓ |
| N | 2533 | 2530 |

*Note: OLS with robust standard errors (clustered by respondent for headline-level analysis); \* p < .05, \*\* p < .01, \*\*\* p < .005 (two-sided). Perceived accuracy measured on 4-point Likert scales. See Appendix A for stimuli and question wording.*

---

[15] Results are similar, however, if the state media trust moderator is treated as continuous (available upon request).

**Table B5.** *Treatment effects on perceived accuracy of state media claims by trust in state media.*

| | Perceived accuracy | |
|---|---|---|
| | False claims | True claims |
| China state media tag | -0.094 | -0.066 |
| | (0.065) | (0.075) |
| Serbia state media tag | -0.042 | -0.070 |
| | (0.065) | (0.067) |
| Generic state media tag | -0.057 | -0.076 |
| | (0.069) | (0.076) |
| Fact-check label | -0.191*** | 0.024 |
| | (0.036) | (0.037) |
| Not very much trust in state media | 0.096* | -0.008 |
| | (0.046) | (0.048) |
| Moderate/great deal of trust in state media | 0.126* | 0.066 |
| | (0.054) | (0.056) |
| China tag × not very much | 0.137 | 0.109 |
| | (0.077) | (0.085) |
| Serbia tag × not very much | 0.082 | 0.127 |
| | (0.076) | (0.080) |
| Generic tag × not very much | 0.056 | 0.045 |
| | (0.078) | (0.086) |
| China tag × moderate/great deal | 0.216** | 0.027 |
| | (0.084) | (0.091) |
| Serbia tag × moderate/great deal | 0.167* | -0.017 |
| | (0.083) | (0.085) |
| Generic tag × moderate/great deal | 0.138 | -0.023 |
| | (0.086) | (0.094) |
| Controls | ✓ | ✓ |
| N | 2533 | 2530 |

*Note: OLS with robust standard errors; * p < .05, ** p < .01, *** p < .005 (two-sided). Perceived accuracy on 4-point Likert scales. See Appendix A for stimuli and question wording.*