# Individual misinformation tagging reinforces echo chambers; Collective tagging does not.

Junsol Kim[1], Zhao Wang[2], Haohan Shi[3], Hsin-Keng Ling[4], James Evans[1,2,5]*

[1] Department of Sociology, University of Chicago, Chicago, IL, USA
[2] Computational Social Science, University of Chicago, Chicago, IL, USA
[3] School of Communication, Northwestern University, Evanston, IL, USA
[4] Department of Sociology, University of Michigan, Ann Arbor, MI, USA
[5] Santa Fe Institute, Santa Fe, NM, USA
*Corresponding author. E-mail: jevans@uchicago.edu

**Fears about the destabilizing impact of misinformation online have motivated individuals and platforms to respond. Individuals have become empowered to challenge others' online claims with fact-checks in pursuit of a healthier information ecosystem and to break down echo chambers of self-reinforcing opinion. Using Twitter data, here we show the consequences of individual misinformation tagging: tagged posters had explored novel political information and expanded topical interests immediately prior, but being tagged caused posters to retreat into information bubbles. These unintended consequences were softened by a collective verification system for misinformation moderation. In Twitter's new platform, Community Notes, misinformation tagging was peer-reviewed by other fact-checkers before exposure to the poster. With collective misinformation tagging, posters were less likely to retreat from diverse information engagement. Detailed comparison suggests differences in toxicity, sentiment, readability, and delay in individual versus collective misinformation tagging messages. These findings provide evidence for differential impacts from individual versus collective moderation strategies on the diversity of information engagement and mobility across the information ecosystem.**

# Main

The visibility of mis- and disinformation online have attracted substantial attention around the world with demonstrations of their direct influence on major collective action in the world[1–5]. These actions range from buying and selling stocks[2] and avoidance of vaccines[3] to the attempted coup and occupation of the U.S. Capitol by rioters[4]. Legitimate fears about the destabilizing influence of false online information have inspired and put pressure on both individuals and platforms to respond. Individuals proactively correct others' claims by deploying links to fact-checking websites, such as PolitiFact and Snopes[6–10]. With the potential for amplifying misinformation through filter bubbles[11,12], social media platforms like Twitter and Facebook have come under public and political pressure to implement misinformation moderation strategies[13–15].

Individuals have become empowered to challenge others' online claims with misinformation tags (or fact-checks) as "vigilantes" in pursuit of a healthy information ecosystem and to break down ideological echo chambers[6–8]. These in-the-wild, vigilante misinformation tags tend to target political outgroups[6,7,9], exposing tagged posters to opposing ideological perspectives. It is less clear, however, whether their misinformation tagging motivates targeted posters to explore diverse political contents afterward. Earlier research on motivated reasoning suggests that misinformation tags contradicting targeted poster's beliefs could "backfire" and reinforce preexisting beliefs[16,17], which could discourage people from exploring diverse information[18]. By contrast, a growing body of research argues that misinformation tagging does not backfire, but reduces engagement with misinformation and expands it with diverse information[13,14,19,20]. These mixed findings suggest that the effects of misinformation tagging could depend on the method of correcting misinformation. Individual misinformation tagging by other users often involves toxic and intolerant messages that dehumanize targeted posters[9,21], potentially hindering their willingness to explore diverse information[22].

Platforms have experimented with institutionalized systems that verify the accuracy of content through collective inputs from a wider distribution of users. Notably, on Twitter's new platform, Community Notes (formerly Birdwatch), misinformation tags undergo a formal peer-review process by diverse users before being revealed to the original posters and broader Twitter user community[8,13,14]. Rather than indiscriminately exposing users to misinformation tags, Community Notes selectively exposes misinformation tags that receive votes from heterogeneous user groups, ensuring that they are verified across a broad spectrum of perspectives[13] to activate the wisdom of crowds[23,24]. The platform also assesses the alignment of users' prior contributions with the crowd's decisions, filtering out voters who frequently oppose and backlash against valid fact-checks on misinformation. Although individual tags may be noisy and less effective, aggregating them collectively could lead to high-quality crowd judgments that align with expert fact-checks across a range of topics, from COVID-19 to politics[14,25–27]. Furthermore, the Community Notes platform has specifically instituted norms that deter toxic and intolerant misinformation tagging messages[28], potentially enhancing the efficacy of misinformation

moderations and gently encouraging posters to leave their echo chambers and explore a broader world of diverse information.

In this study, we explore the impacts of "individual" and "collective" misinformation tagging on tagged posters' echo chambers. Echo chambers refer to "bounded, enclosed media spaces that have the potential to both magnify messages delivered within them and insulate them from rebuttal"[29,30], which could increase susceptibility to misinformation[11,31,32]. This definition guides our operationalization of echo chambers in two distinct ways.

One indicator of echo chambers is their lack of political diversity and limited interaction with politically diverse, cross-cutting sources of information. Prior research has measured echo chambers by selective engagement with like-minded news sources, which insulate people from opposing perspectives that could empower rebuttal[33,34]. This measure strongly correlates with other echo chamber indicators, such as intensive interactions with like-minded users (i.e., homophily)[35,36]. Literature suggests that lack of exposure to and cross-verification through opposing perspectives could erode the ability to find, evaluate, and use information effectively[11,37,38]. It could provide users with the illusion that their views are publicly supported[39,40], weakening their overall immunity against misinformation.

The other key indicator of echo chambers is their absence of content diversity resulting from limited engagement with diverse, unfamiliar topics. Emerging literature has documented the rise of socio-political endogamy, noting that both left and right increasingly develop distinct topical interests, encompassing knowledge bases, cultural tastes, and lifestyles[41–43]. For example, left-leaning individuals are more likely to engage with basic science books about physics, astronomy, and zoology, while right-leaning individuals prefer those about applied and commercial sciences like criminology, medicine, and geophysics[43]. In this way, political polarization spills over into a variety of other topics, leading to multi-dimensional segregation where opposing political groups share progressively less common ground and inhabit different realities even in topics apparently unrelated to politics[41,44]. Topical echo chambers, which magnify topics prevalent within one political group and insulate them from others, can problematize intergroup communication and interaction.

Does exposure to each type of misinformation tagging encourage or discourage posters from exploring diverse information? To answer this question, we use large-scale digital traces from the platform formerly known as Twitter (*X* as of July, 2023) to identify posters exposed to each approach to misinformation tagging. First, we identify posters targeted by individual misinformation tags shared by online "vigilantes". These posters' tweets received other individuals' voluntary replies, citing fact-checking articles from PolitiFact, one of the largest and most studied professional fact-checking organizations in the United States[7,10]. Second, we examine posters targeted by collective misinformation tags resulting in "verified" tags. These posters' tweets received notes that contain collectively verified fact-checks through Twitter's

Community Notes platform. Fig. 1a visualizes the mechanism of each type of misinformation tagging, which represent the most prevalent misinformation moderation strategies on Twitter[6–10,13–15]. Extended Data Fig. 1 presents an example of individual and collective tags that correct topically identical, COVID-19 misinformation.

Using approximately 700,000 tweets that cite news sources—including posts, retweets, and quotes—posted by 8,000 users before and after they were targeted by misinformation tags, we estimate the effects of these tags on the posters' ideological echo chambers. Specifically, we measure echo chambers using political and content diversity in their posting and sharing behavior (see Fig. 1b). Political diversity measures whether a poster's tweet cites a source with opposing political stance (e.g., a right-leaning poster references left-leaning articles)[5,45]. Content diversity measures whether a tweet discusses novel topics unfamiliar in the poster's historical tweets. We apply a transformer-based sentence embedding model (SentenceBERT) to extract a high-dimensional, semantic vector representation for each tweet, and aggregate the vectors of each author's historical tweets to produce an average semantic vector for each poster. We then measure the distance between a particular tweet and the poster to assess the degree to which this tweet expands the poster's content diversity. As our data focus on tweets citing news sources, we assume that the increase of content diversity indicates the exploration of novel political news topics. For example, consider a user who regularly consumes and shares news about COVID-19 but begins to discuss U.S. tax and labor issues as well. This shift indicates an increase in the user's content diversity, as detailed in Supplementary Table 1. We consider both types of diversity because they represent different dimensions that could reinforce one another in limiting exposure to information and exacerbating filter bubbles and echo chambers on social media[17,46].
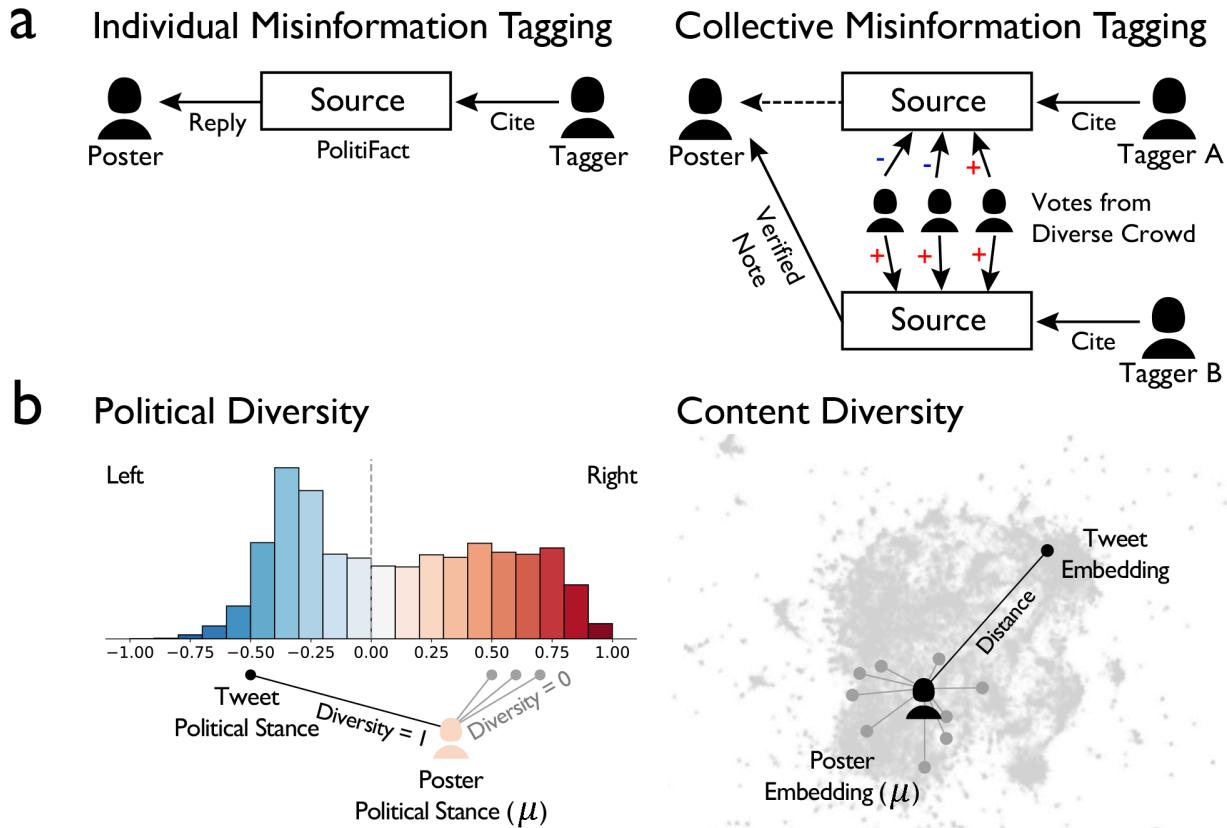
**Fig 1. Misinformation Tagging and Outcomes Measurement. a**, Individual misinformation tagging in which individuals cite PolitiFact fact-checking articles. Collective misinformation tagging through the Community Notes platform, which selectively exposes verified misinformation tags that receive diverse votes as helpful. **b**, Operationalization of tweet political and content diversity. Political diversity captures whether a poster cites a source with opposing political stance (binary 0/1), assessed from the aggregate stances of referenced sources. Content diversity captures whether a post discusses topics unfamiliar to the author's historical tweets (continuous), assessed with the distance between the poster's average tweet and a particular tweet within a contextual embedding (sentenceBERT pre-trained on Twitter)[47].

# Results

We aim to investigate the effects of individual and collective misinformation tagging on political and content diversity using large-scale Twitter data. In our observational data, treatments (i.e., exposure to misinformation tagging), however, are not randomly assigned to misinformation posters, which pose challenges for identifying the causal effects of misinformation tagging. To address these concerns, we apply interrupted time series (ITS) and delayed feedback (DF) analysis, which help eliminate non-causal explanations under certain assumptions.

## Interrupted Time Series (ITS) Analysis

Interrupted Time Series (ITS) analysis investigates whether the trend in political and content diversity shifts after misinformation tagging. ITS assumes that without the intervention of misinformation tagging, the pre-treatment trend (i.e., before misinformation tagging) would persist, and the immediate change in trend after misinformation tagging is attributed to effects from tagging. We control for user-level fixed effects to correct for time-invariant user characteristics.

As shown in Fig. 2a and Extended Data Table 1, posters manifest an increasing tendency to explore novel political information before being fact-checked by misinformation tags. Specifically, before individual and collective misinformation tagging, posters increase the political diversity ($\beta$=.237, $p$<.001; $\beta$=.309, $p$=.024) and content diversity ($\beta$=.007, $p$<.001; $\beta$=.003, $p$=.461) of their information engagement over time.

Having their posts criticized by individual misinformation tags, however, causes posters to retreat within an "information bubble". Immediately after tagging, posters significantly decrease the political diversity ($\beta$=-1.009, $p$<.001) and content diversity ($\beta$=-.030, $p$<.001) of their posts. After tagging, the slope becomes nearly flat, indicating that posters' future posts continue to collapse in both political diversity ($\beta$=-.150, $p$=.059) and content diversity ($\beta$=-.010, $p$<.001).

By contrast, collective misinformation tagging does not cause individuals to retreat within their previous information bubble. The data even reveals a slight increase in political diversity ($\beta$=.270, $p$=.629) and a significant increase in content diversity ($\beta$=.040, $p$=.006) immediately after tagging. Nevertheless, collective misinformation tagging has only a temporary effect on individual posters. Specifically, increased diversity does not persist, with political diversity ($\beta$=-.358, $p$=.072) and content diversity ($\beta$=-.014, $p$=.009) eventually converging to levels experienced before the initial misinformation tags occur. Despite the steepness of the slope following collective tagging, our analysis indicates that content diversity does not significantly drop below the pre-tagged period (see Supplementary Method 1).

Additional analyses reveal the effects of misinformation tagging on the proximity between posters and misinformation taggers. This suggests that Twitter navigation likely makes posters more visible to fact-checkers as they venture into foreign territory (see Fig. 2b). Exposure to fact-checks causes them to retreat back into their information bubbles, distancing them from the foreign stances that fact-checked them (see Supplementary Method 2).
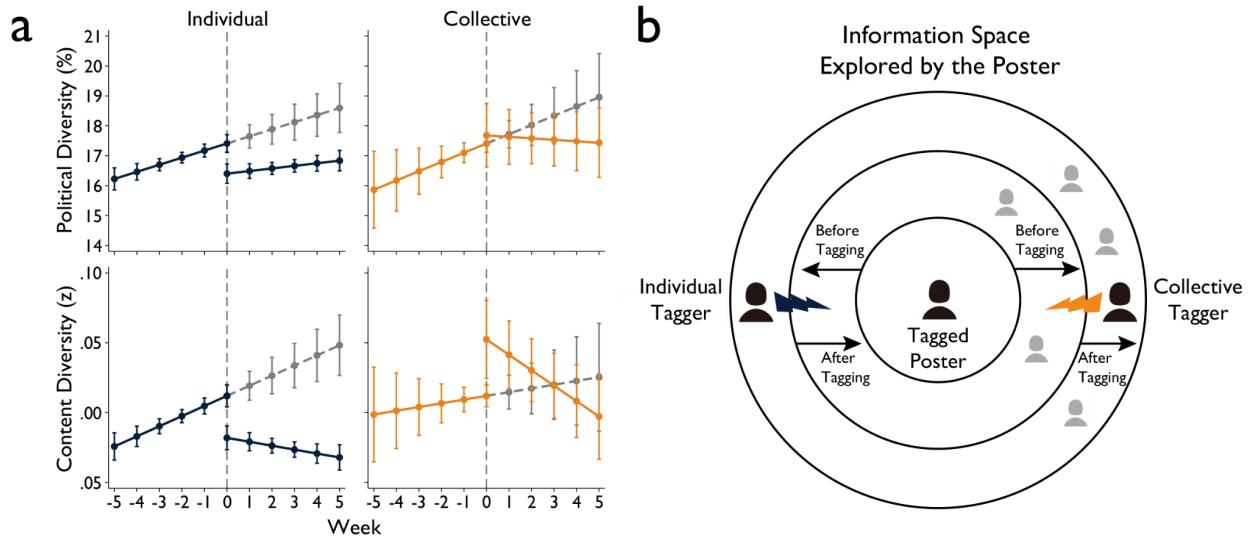


**Fig. 2. Political and Content Diversity Change with the Intervention of Individual and Collective Misinformation Tagging. a**, Results from Interrupted Time Series (ITS) analysis. The *x*-axis denotes the timeline of tweets posted before and after tagging, with negative values representing the number of weeks before posting tagged tweets and positive values the number of weeks after. The *y*-axis denotes political and content diversity, with dots capturing the average diversity score of the corresponding week, and error bars indicating 95% confidence intervals. Solid lines connect the dots revealing trends of political and content diversity before and after tagging, with gray dotted lines tracing the counterfactual trend if fact-checks had not occurred. **b**, Illustration of political and content diversity dynamics before and after tagging. Before individual and collective tagging, posters exhibit increased political and content diversity, which increases the likelihood of encountering a fact-checker. After individual tagging, posters retreat into information bubbles; after collective tagging, they venture further beyond them.

To better understand what happens when posters retreat to their information bubbles, we conduct a series of descriptive analyses (see Supplementary Table 4). When posters reduce their political and content diversity, the number of tweets (comprising posts, retweets, and quotes) posted per day significantly increases, indicating that users are more active within their information bubbles. Specifically, the number of tweets per day is negatively correlated with political diversity ($r=-.107$, $p<.001$) and content diversity ($r=-.052$, $p<.001$). Similarly, we find that the type of posting is different; the proportion of "retweets (i.e., tweets simply sharing other users' tweets)" out of the entire tweets per day is negatively correlated with political diversity ($r=-.046$, $p<.001$) but positively correlated with content diversity ($r=.012$, $p=.001$). This indicates that users actively post tweets rather than passively retweet other users' tweets when they exhibit low

political diversity. To demonstrate the significant effects of misinformation tagging on political and content diversity, irrespective of these factors, we have adjusted for the number of tweets posted per day in both interrupted time series (ITS) and delayed feedback (DF) analyses. We have also controlled for the proportion of links to political sources per day, which did not meaningfully change our results (see Supplementary Table 5).

## Delayed Feedback (DF) Analysis

We employ delayed feedback (DF) analysis to further strengthen our causal inference[48]. In our DF analysis, we estimate baseline changes (i.e., changes in outcomes that occur without tags) to answer the question: "Are shifts in political and content diversity attributable to tagging, or do similar changes occur even without tagging?" Pairs of tweets containing similar misinformation, targeted by misinformation tagging at different times, are matched to construct a control group, consisting of posters whose problematic tweets have not yet been tagged due to "delayed feedback," and a treatment group of posters who have. For instance, Extended Data Fig. 2 presents an illustrative example involving a pair of matched tweets and tags.

In Fig. 3a, post-treatment ($t_1$) represents the time window when treatment tweets are tagged but control tweets are not, and pre-treatment ($t_0$) represents the time window with equal duration $t_1$ when both treatment and control tweets are untagged. Changes in the outcomes between $t_0$ and $t_1$ in the control group reflect "baseline changes," which indicate changes without tags. Changes between $t_0$ and $t_1$ in the treatment group reflect "treated changes," which indicate changes with tags. We compare the difference in pre-post change between control and treatment groups (i.e., baseline vs. treated changes) to identify the effects of misinformation tagging on political and content diversity. DF analysis assumes that, in the absence of treatment, both control and treatment groups would exhibit parallel trends. We control for user-level fixed effects to control for time-invariant, user-specific characteristics.

Fig. 3b and Extended Data Table 2 present results from the DF analysis. Our DF analysis demonstrates that changes are indeed due to tagging, showing that "treated changes" are significant above and beyond "baseline changes." Consistent with the ITS findings, DF analysis indicates that individual "vigilante" misinformation tags lead to a notable decrease in political diversity ($\beta$=-5.886, $p$=.002). Nevertheless, individual misinformation tagging does not significantly affect content diversity ($\beta$=.018, p=.652). Although ITS analyses show that content diversity decreases after tagging, DF analyses suggest that content diversity does not decrease beyond baseline changes observed without tags. Collective, "verified" misinformation tags, by contrast, do not produce a significant decrease in political diversity ($\beta$=1.219, p=.690) and even increase content diversity following tagging ($\beta$=.274, $p$<.001).
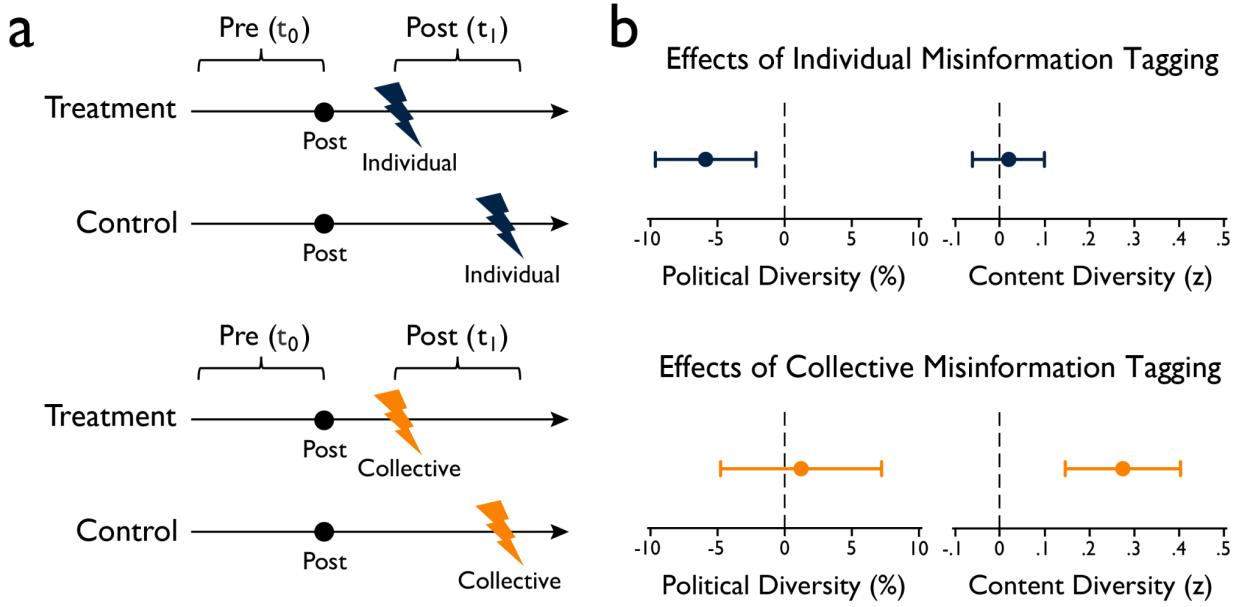
**Fig. 3. Delayed Feedback (DF) Analysis. a**, Pre- and post-treatment periods. Post-treatment ($t_1$) represents the time window when "treated" tweets are tagged but control tweets are not. Pre-treatment ($t_0$) represents the time window with equal duration $t_1$ when both treatment and control tweets remain untagged. **b**, The effects of individual and collective misinformation tagging on political and content diversity, which are estimated by the difference in pre-post changes among outcomes between treatment and control.

## Linguistic Characteristics of Misinformation Tags

Individual and collective misinformation tagging messages manifest different linguistic characteristics. As shown in Fig. 4 and Extended Data Table 3, we find that individual misinformation tags exhibit twice the toxic content ($t(7496)=9.86$, $p<.001$) and convey more negative sentiment ($t(7731)=-2.14$, $p=.033$) compared to collective misinformation tags. Collective tags express slightly higher positive sentiment and produce messages with more neutral sentiment than individual tags. Furthermore, individual tag messages are much shorter ($t(7731)=-26.95$, $p<.001$) and more readable ($\chi2(7)=155.32$, $p<.001$) than collective tags. While 53.53% of individual tags necessitate a college-level reading comprehension or higher, 75.77% of collective tags demand this level. Moreover, the delay between posting misinformation and fact-checks is shorter for individual than collective tagging ($t(7731)=-2.13$, $p=.033$). These findings demonstrate that individual tags convey their messages quickly through messages that are succinct, straightforward, emotive, and sometimes toxic. In contrast, collective tags are more slowly communicated through lengthy, complex messages, devoid of emotional undertone or toxicity.
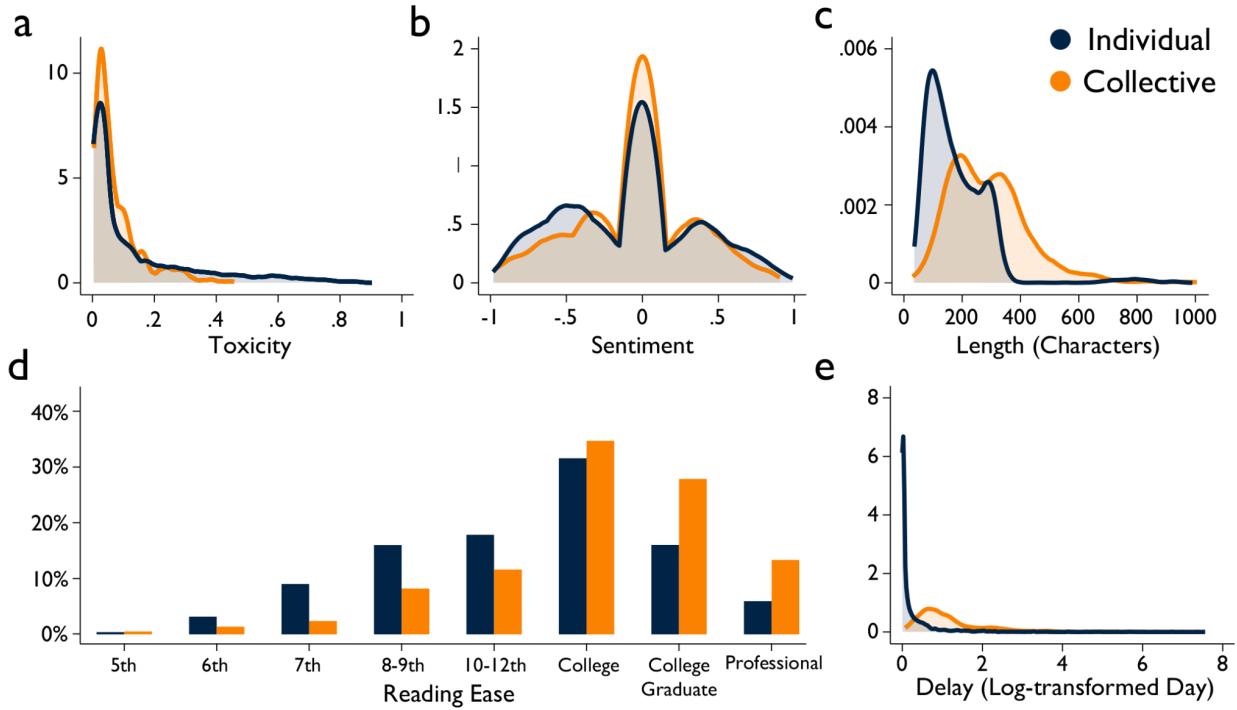
**Fig. 4. Linguistic Characteristics of Fact-checking Messages.** We present a univariate kernel density function for continuous variables (toxicity, sentiment, length, delay) and a histogram for the categorical variable (reading ease). The purple line represents the distribution within individual misinformation tags; the yellow line represents the distribution within collective tags.

Based on linguistic differences between individual and collective tags, we question whether gaps in the effects of individual versus collective tags persist even after controlling for these differences. First, we control for toxicity by excluding tags with a toxicity level higher than 0.4 and retaining only non-toxic tags. Second, we control for sentiment by removing tags with either positive (>0.2) or negative (<-0.2) sentiments, keeping only neutral tags. Third, we control for length by excluding tags longer than 400 characters and retaining short tags. Fourth, we control for readability by excluding tags that require college-level or higher readability and selecting tags that are relatively easy to read. Fifth, we control for delay by omitting any tags associated with delays longer than 48 hours (log(delay+1) > 1.10) and focusing on quick tags.

We find that these controls do not significantly impact our reported outcomes except for length. As shown in Extended Data Table 4, the gap between individual and collective tagging in political diversity—but not content diversity—diminishes from 1.279 to 1.071 when controlling for length, rendering the gap insignificant. Nevertheless, controlling for length accounts for 16.26% of the gap between individual and collective tagging in political diversity. This indicates that linguistic characteristics explain a modest but nontrivial portion of the differential impacts between individual and collective tagging. Nevertheless, these measured qualities do not account for the vast majority of the difference.

## Control Analyses

We find systematic differences in misinformation that receives individual and collective tagging, as well as in the posters who are corrected by each type of tagging. Therefore, we conduct a series of analyses to address these potential confounders with Interrupted Time Series (ITS) models.

First, we observe that individual taggers focus more on political topics, while collective taggers correct a more diverse range of topics (see Supplementary Table 2). As shown in Extended Data Table 5, the nine most frequent topics in our dataset include political topics known to trigger divisive, polarized reactions in US politics (see Methods: Topic modeling). These topics account for 87.25% of the corrections made through individual tagging but only 59.49% of the corrections made through collective tagging. Still, there is a considerable overlap between the topics corrected by individual and collective tagging. Even when we limit the sample to these nine topics, our results do not meaningfully change (see Extended Data Table 6).

We also find that gaps between individual and collective tags are significant and even slightly larger when they correct identical topics of misinformation, employing propensity score weighting (PSW) (see Supplementary Method 3 and Extended Data Table 6). These results demonstrate that impacts from individual and collective tagging differ, even when they correct topically identical messages. We note that collective tagging is less likely to correct political topics than individual tagging, but more effective in causing original posters to explore diverse content when successfully deployed on political topics.

Second, we find that right-leaning posters are more likely to be corrected by individual tagging. The proportion of right-leaning posters corrected by individual tags is 53.17% while right-leaning posters corrected by collective tags is 44.14% (Difference=9.03%; $z$=4.89; $p$<.001). Therefore, we control for political stances of misinformation posters and taggers in our sample. Prior research found that right-leaning individuals are more likely to post misinformation, and left-leaning individuals are more likely to correct it[10], which is reproduced in our data as shown in Supplementary Method 4 and Supplementary Table 3. We compare the effects of individual and collective tagging in this common scenario where right-leaning individuals are corrected by left-leaning ones, either through individual or collective tagging. Specifically, we limit the sample to cases where left-leaning taggers or voters correct right-leaning posters. Even in this analysis, we find a significant difference between individual and collective tagging (see Extended Data Table 7).

Third, we find that popular users are more likely to receive collective tags than individual tags, which is consistent with prior literature (see Supplementary Fig. 1)[8]. To examine the differences between individual and collective tags when focusing on less popular, everyday users, we exclude those whose number of followers exceeds 2,967, the average number of followers

among users corrected by individual tags. We find the results are consistent overall (see Extended Data Table 8), but suggest that collective tagging of low popularity posters is even more effective, relative to individual tagging, than with high popularity users. This may indicate the "inoculation" of popular users to critique, an increased sensitivity among unpopular users to collective "nudges" [49], or both.

## Robustness Checks

We verify our findings with a battery of robustness checks. First, we seek to avert concerns over the presence of bots on Twitter by reanalyzing our data excluding identified bot accounts [2,5]. Second, we reanalyze our relationship controlling for potentially insincere informational activities, such as citing sources of low credibility and intentionally spreading fake news. Third, we attempt to avoid situations in which posters simply criticize distant information without honest consideration by filtering out posts with negative sentiment. Fourth, considering the low visibility of individual tags in Twitter's message-reply interface [6,8], we restrict the sample to original posters who replied to (and thereby read) the individual tags and remove non-responders. Fifth, we identify all tweets within the sample that mention keywords related to receiving community notes broadly and remove them, as they could confound our measure of content diversity. To address concern regarding the effect of replying directly to individual taggers, which could confound the measure of political diversity, we also identify and remove all tweets that reply directly to individual taggers. Sixth, to strictly identify PolitiFact links that correct the original posters, we prompt ChatGPT to annotate whether the links are used to correct the original poster rather than support them. Then, we limit the sample to links that correct the original posters. These alterations do not significantly impact our reported outcomes (see Method: Robustness Checks).

# Discussion

This study provides empirical evidence regarding the impact of individual and collective misinformation tagging on echo chambers. Before misinformation tagging, posters show an increased curiosity in diverse political and topical content. This challenges the conception that misinformation is generated and corrected when people retreat into echo chambers [11,31]. On the contrary, posters become fact-checked when they venture outside those bubbles. Why is exploration followed by misinformation tagging? First, posters could misinterpret unfamiliar and diverse information from a lack of information literacy [50], increasing the chance of posting the misinformation being tagged. Second, news feed algorithms may increase the probability that posters' tweets become visible to people from political outgroups, who are highly motivated to fact-check foreign posters [6,7,14]. Our analysis shows that posters increase the "closeness" to misinformation taggers before fact-checks, which could increase the chance of appearing in fact-checkers' news feeds.

Individual misinformation tagging discourages posters from exploring diverse information. Posters tagged by vigilantes manifest an immediate drop in political diversity, as evidenced by both interrupted time series (ITS) and delayed feedback (DF) analyses. Content diversity also decreases in ITS analyses, although DF analyses do not reveal a significant drop. This suggests that while content diversity decreases after tagging, it does not fall below the baseline change expected without tags. These unintended consequences are mitigated by collective misinformation tagging. Unlike individual tagging, collective tagging does not diminish political and content diversity in both ITS and DF analyses; it even results in a short-term rise in content diversity.

Our analyses show that individual tagging involves short, toxic, and emotion-driven messages. Collective tagging, on the other hand, involves longer, less toxic, emotionally neutral, and deliberative messages revealed to posters longer after their offending posts. These results suggest the trade-off between the effectiveness of established systems for promoting openness and mobility across the information ecosystem, but the efficiency of vigilantes in cleaning it. Low visibility of individual misinformation tagging in Twitter's message-reply interface[6,8] may motivate taggers to use short and potentially toxic messages. Community Notes responded by implementing a more visible interface for collaborative tagging, which reduces the tendency to terseness, facilitating long and deliberate discussion. Also, norms and values underlying participation in Community Notes could prevent taggers from disseminating succinct yet inflammatory messages viewed as unhelpful and instead source diverse perspectives[13].

What mechanisms drive differences in the effects of individual and collective misinformation tagging on echo chambers? We find that linguistic characteristics, such as toxicity, sentiments, and length only partially explain differential impacts between individual and collective tagging. This implies that differences in quality other than linguistic characteristics also exert a direct influence. Literature on the "wisdom of crowds" suggests that while individual tags are susceptible to biases and noise, aggregating tags collectively could correct individual bias, increasing the quality of nonexpert fact-checks[26,51,52]. For example, compared to individual tags, collective tags are more closely aligned with professional fact-checks from experts on a variety of topics, ranging from COVID-19 to politics[14,25–27]. Even though we focus on individual tags that cite professional fact-checks (i.e., PolitiFact), it is possible that interpretations within individual tags might be less effective when not cross-validated like collective tags. For example, individual tags might fail to convey the key points of PolitiFact articles or clearly articulate the relevance of these articles to the original post. Additionally, when multiple fact-checkers co-validate collective tags, these decisions may be perceived as more legitimate and less susceptible to biases, encouraging the original posters to seek out more diverse and cross-validating information[26].

Overall, our findings suggest that misinformation is posted and fact-checked when original posters who were accustomed to like-minded sources associated with low credibility (see Supplementary Table 4) suddenly increase their political and content diversity. In the short term, some might believe that pushing them back into their echo chambers with individual tags seems like an effective way to curb misinformation. Nevertheless, over the long term, this approach could expand the cluster of users immersed in misinformation, depriving them of opportunities to educate themselves with opposing perspectives. The ethical and normative aspects of our research remain open questions, but we argue that collective tagging encouraging exploration could be better for the long-term health of the information ecosystem.

Our analyses have several notable limitations. First, our method for assessing posters' political stances is indirect, through their posting behavior[5]. This approach has been successfully applied to predict political party affiliation and self-described ideology in previous literature[52], but using a direct measure of political ideology or affiliation with social media and survey data would strengthen our assessments. Second, our quasi-experimental methodologies (ITS and DF) depend on assumptions for causal inference. We employ topic modeling and matching to enhance tweet comparability within treatment and control groups, but acknowledge that unobserved time-variant confounders may influence posters' responses. Third, although we have employed a popular bot detection algorithm, recent studies have suggested that algorithmic removal of bots is challenging and may introduce additional bias[53]. Therefore, we report the full results with and without the algorithmic removal of bots, demonstrating that our results are consistent. To thoroughly remove bots, future research could match social media data with survey or administrative data (e.g., voter records) to ensure the authenticity of participants[54]. Fourth, Twitter (X as of July, 2023) closed access to the Academic Research API, which had been freely available to eligible researchers until May 2023. This could limit other researchers' ability to reproduce our findings with recent data after May 2023[55]. Despite these limitations, our study uncovers a significant and substantial relationship between fact-checks and reduced information diversity. We also demonstrate the power of designed institutions, like collective fact-checking on Twitter, to moderate the negative, narrowing effects of fact-checking on information exploration.

# Methods

## Data

Using Twitter API v2.0 with academic research access, we collected Twitter data to explore the effects of individual and collective misinformation tagging. First, we identified 9,372 users targeted by individual misinformation tagging from 2021/10/1 to 2022/3/25. We selected users whose tweets received fact-checking replies that contain URLs to fact-checking articles from "politifact.com." Second, we identified 1,465 users targeted by collective tagging from 2022/12/19 to 2023/3/31, when Community Notes were made public to Twitter users globally[56]. In Community Notes, users can flag any tweets as misinformation with notes, and other members vote for the helpfulness of the notes. (Users also have the option to flag tweets they believe are free from misinformation; however, these instances have been excluded from our analysis.) Collectively verified notes that received the above-threshold helpfulness votes from a diverse set of users are then made public to the original user (who posted the misinformation) and the broad Twitter audience[13]. In our work, we only considered notes with above-threshold helpfulness votes. Note that the platform also assesses the alignment of users' prior contributions with the crowd's decisions, filtering out voters who frequently oppose and backlash against valid fact-checks on misinformation (see Supplementary Method 5).

Due to the rate limit of Twitter API, we only collected data from regular Twitter users, excluding organizations' and celebrities' accounts with 50,000 or more followers. Additionally, to focus on individual users, rather than organizational accounts (e.g., CNN, Fox News, etc), we removed 1,659 users identified as organization accounts by the M3Inference library[57,58]. We further removed 1,445 users who were fact-checked more than once within the period of data collection to avoid the potential for them to become desensitized for repeated fact-checks. After filtering the data, our final dataset included 7,733 users, where 6,760 users were targeted by individual misinformation tagging and 973 users were targeted by collective misinformation tagging. We found that individual tagging is more frequent than collective tagging in our dataset due to the cross-validation process required to expose collective tags. This leads to an imbalance in group size between users corrected by individual and collective tags. Nevertheless, our statistical models (interrupted time series and delayed feedback models) do not assume equal group size for comparison between the effects of individual and collective tagging. Also, we found that 16.33% of tweets that received individual tags and 15.60% of tweets that received collective tags were removed by Twitter or by the original poster. The probability of removal is similar between individual and collective tags (Difference=0.73%, $z=.629$, $p=.529$).

Finally, we collected users' historical tweets—including posts, retweets, and quotes—which span two months before posting tagged tweets and two months after exposure to misinformation tagging, resulting in 1,409,845 tweets in total. Posts typically indicate active engagement with diverse political sources and topics, allowing users to express their opinions. In contrast, retweets

and quotes—which involve sharing others' tweets—suggest more passive engagement, not necessarily reflecting personal views. We utilize these three types of behaviors for a more comprehensive measurement of users' information engagement[59,60]. We assume that individual misinformation taggings are exposed to users when they are posted, and collective misinformation taggings are exposed to users when they are made public following the above-threshold helpfulness votes. For our statistical analyses, we included 712,948 tweets with observed political and content diversity scores. We received a determination from the Institutional Review Board that the study is not considered human subjects research and did not require review.

## Political Diversity

Political diversity measures whether a user posted a tweet that referenced sources having an opposite political stance. Specifically, we determine the political stance of the referenced source by extracting the domain (e.g., cnn.com) of the source and check it from MediaBias/FactCheck database (MBFC; https://mediabiasfactcheck.com/), as suggested in prior works[5,45]. MBFC provides a continuous score for 4,874 websites to indicate each source's political stance, ranging from -1 (extreme left) to 1 (extreme right). Our additional analysis shows that political stance scores from MBFC show significant inter-rater reliability with another database of the political stance of news media, AllSides.com (see Supplementary Method 6).

We then calculate a user's political stance by averaging the political stance scores of sources referenced in their historical tweets which span two months before posting tagged tweets and two months after misinformation tagging (see Supplementary Fig. 2). Users who predominantly cite left-leaning media are considered "left," and those who cite right-leaning media are considered "right". Specifically, users with negative average political stance scores are categorized as left, while those with positive scores are categorized as right. Finally, we assign a binary value to represent a user's political diversity: 1 (diverse) if a user cited a source that has an opposite political stance from the user's own political stance, 0 (not-diverse) if a user cited a source with the same political stance.

The mean political diversity score is .166, and the standard deviation is .372 (N=712,948). Political diversity is negatively correlated with the number of tweets posted per day (r=-.107, p<.001) and the proportion of retweets (r=-.046, p<.001). This indicates that users are more active within information bubbles, actively posting tweets rather than passively retweeting other users' tweets within these bubbles (see Supplementary Table 4).

## Content Diversity

Content diversity measures whether a user posted a tweet with a topic that is rarely discussed in the user's historical tweets. We apply the Twitter4SSE model, a transformer-based sentence embedding model (SentenceBERT) that was initialized from BERTweet (a RoBERTa model trained on 850 million tweets from 2012/1 to 2019/8 and 5 million tweets related to COVID-19 pandemic), to encode the meaning of a tweet into a 768-dimensional vector[61,62]. The model was further optimized based on recent data (75 million tweets from 2020/11 to 2020/12) using Multiple Negatives Ranking Loss (MNRL) to identify semantic similarity based on the principle that tweets quoting or replying to the same original tweet are likely discussing related ideas[61]. If a pair of tweets quoted or replied to the same tweet, the semantic similarity between them is assumed to be high.

To apply the Twitter4SSE model, we first conduct the identical data preprocessing steps to clean the tweets, which includes: eliminate URLs and mentions and transform the text to lowercase to reduce the presence of generic texts[61]. Next, we represent each tweet with a 768-dimensional semantic embedding (Supplementary Fig. 3 shows the visualization). Finally, we measure the cosine distance between the user embedding and tweet embedding (see Fig. 1b) to represent the content diversity of the current tweet. The user embedding is the average embedding of the user's historical tweets (see Fig. 1b). Estimating the distance in the embedding space has been frequently used to quantify the diversity of user activities in the online platform[46,63]. The distance ranges from 0 to .835, with 0 representing homogeneous content and .835 representing extremely diverse content. The mean content diversity score is .357, and the standard deviation is .109 (N=712,948). We find that political and content diversity are slightly correlated (r=.020, p<.001), assessing conceptually distinct aspects of diversity.

Extended Data Table 1 shows an example of how content diversity scores are assigned. In this example, the user primarily shows interests in COVID-19 related misinformation. However, as the user explores diverse topics—tax, LGBTQ+, international issues, and labor—the content diversity score increases.

Content diversity is negatively correlated with the number of tweets posted per day (r=-.052, p<.001) but positively correlated with the proportion of retweets (r=.012, p<.001). In other words, users tend to retweet others' tweets rather than posting their own tweets when increasing content diversity (see Supplementary Table 4).

## Interrupted Time Series (ITS) Analysis

We apply Interrupted Time Series (ITS) analysis to examine how individual and collective misinformation tagging affect the trend of political and content diversity in posting behavior. We fit the ITS model to the time series around fact-checking events, spanning five weeks (35 days)

before posting the fact-checked tweet and five weeks after fact-checking. To compare the differential impacts of individual and collective misinformation tagging, we formulate the following multi-group ITS model. We control for user fixed effects to eliminate the user-related unobserved time-invariant heterogeneity that could possibly affect the outcomes. Additionally, the number of tweets posted per day is negatively correlated with political diversity (r=-.107, p<.001) and content diversity (r=-.052, p<.001), indicating that users are more active within information bubbles. Therefore, we control for the number of tweets posted per day to ensure that our analysis focuses on variations in diversity rather than engagement volume.

For each tweet, let $Y$ be the outcome variable (i.e., political or content diversity of a specific tweet), $W$ is the weeks before posting the tweet with misinformation (negative values) or after misinformation tagging (positive values). Note that we measure $W$ by dividing the days by 7. For example, if a particular tweet is posted 3 days before posting the tweet, $W$ is -3/7. $T$ is an indicator of the treatment status where 0 represents a tweet posted before misinformation tagging and 1 represents after tagging. $C$ is an indicator of the type of misinformation tagging where 0 represents individual tagging and 1 represents collective tagging. $N$ corresponds to the number of tweets per day (control variable), $\alpha$ corresponds to the user fixed effect, and and $\epsilon$ is the error term. Then the ITS model is defined:

$$Y = \beta_0 + \beta_1 W + \beta_2 T + \beta_3 WT + \beta_4 WC + \beta_5 TC + \beta_6 WTC + \beta_7 N + \alpha + \epsilon \qquad (1)$$

Here, $\beta_0$ is the intercept, $\beta_1$ is the slope before individual misinformation tagging. $\beta_2$ is the change in the outcome immediately after the individual misinformation tagging. $\beta_3$ is the slope change before and after individual misinformation tagging. $\beta_1 + \beta_4$ is the slope before collective misinformation tagging. $\beta_2 + \beta_5$ is the change in the outcome immediately after the collective misinformation tagging. $\beta_3 + \beta_6$ is the slope change before and after collective misinformation tagging. Thus, $\beta_4$, $\beta_5$, $\beta_6$ are the terms that estimate the "differences" between the effects of individual and collective misinformation tagging. Supplementary Table 6 shows how these estimates correspond to each cell in Extended Data Table 2 for each outcome.

Before estimating the model, political diversity (binary variable) has been multiplied by 100 so that the coefficients are interpretable as absolute percentage point changes. Content diversity has been normalized to $z$-scores (i.e., the number of standard deviations from the mean). When estimating the statistical significance of the estimates, all $p$-values are two-sided. The thresholds for statistical significance is set at $p < .05$, and marginal significance is set at $p < .1$.

## Delayed Feedback (DF) Analysis

In addition to the interrupted time series (ITS) analysis, we conduct a delayed feedback (DF) analysis to estimate the causal impacts. We begin by establishing control and treatment groups: each tweet is paired with another tweet that was subject to misinformation tagging at an earlier

time. Specifically, for every tweet in a control group, we search for a corresponding treatment tweet using the following criteria: (1) They must have been fact-checked using the same approach, either individual or collective misinformation tagging. (2) They should have been fact-checked prior to the control tweet. (3) They should have the same topic, considering that distinct topics of misinformation could lead to different levels of political and content diversity (see Method: Topic Modeling for a detailed explanation of the topic modeling process). (4) They should have been posted no more than seven days apart from the control tweet. In cases where we have multiple tweets that meet these criteria, we choose the one with the closest posting time to the control tweet. This results in 476 pairs of tweets in control and treatment groups.

For each pair of tweets, we identify two time windows: pre-treatment ($t_0$) and post-treatment ($t_1$). $t_1$ represents the time window when the treatment tweets are fact-checked but the control tweets are not. If the duration of $t_1$ exceeds a seven-day window, we use the data within the seven-day window after receiving the tags, considering that the timing of the fact-check could affect the outcome. $t_0$ represents the time window (with equal duration of $t_1$) when both the treatment and control tweets are not fact-checked. Then we design the following difference-in-differences model to assess the impacts of misinformation tagging.

For each tweet, let $Y$ be the outcome variable (i.e., political or content diversity of a specific tweet). $T$ is a binary variable indicating whether the treatment tweet, but not control tweet, receives the treatment (i.e., misinformation tagging). $G$ is a binary variable indicating whether the tweet is assigned in the treatment (i.e., 1) or control group (i.e., 0). $C$ is an indicator of the type of misinformation tagging where 0 represents individual tagging and 1 represents collective tagging. $N$ corresponds to the number of tweets per day (control variable), $\alpha$ corresponds to the user fixed effect, and and $\epsilon$ is the error term.

$$Y = \beta_0 + \beta_1 TG + \beta_2 TGC + \beta_3 T + \beta_4 TC + \beta_5 N + \alpha + \epsilon \qquad (2)$$

$\beta_0$ is the intercept. $\beta_1$ is the difference in pre-post change in the outcome between the control and treatment group for "individual misinformation tagging". $\beta_1 + \beta_2$ is the difference in pre-post change for "collective misinformation tagging". Thus, $\beta_2$ is the term that estimates the "difference" between the effects of individual and collective misinformation tagging. $\beta_3$ and $\beta_4$ account for the baseline changes in the outcomes. Supplementary Table 7 shows how these estimates correspond to each cell in Extended Data Table 3 for each outcome.

Like ITS models, political diversity (binary variable) has been multiplied by 100 so that the coefficients are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (i.e., the number of standard deviations from the mean). When estimating the statistical significance of these coefficients (e.g., as in Table 1), all p-values are two-sided. The thresholds for statistical significance is set at $p < .05$, and marginal significance is set at $p < .1$.

Supplementary Fig. 4 illustrates average baseline changes of political and content diversity obtained from DF analysis. According to the baseline changes, we find that political diversity significantly increases ($\beta$=1.956, $p$=.042), but content diversity marginally decreases ($\beta$=-.039, $p$=.060) if the problematic tweet is not tagged. In other words, we find that political diversity consistently increases over time without tagging in the control group. On the other hand, we find that content diversity may decrease over time. Our results show that the effects of individual and collective tagging are above and beyond these baseline changes (See Results: Delayed Feedback Analysis).

## Topic modeling

We apply BERTopic to extract latent topics from tweets that received misinformation tags[64]. Specifically, we first represent each tweet with a 768-dimensional semantic embedding using Twitter4SSE. Then, we map the embeddings to a 5-dimensional space via UMAP (Uniform Manifold Approximation and Projection) to mitigate the curse of dimensionality[65,66]. Next, we apply HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise) to identify clusters of topics[67]. Unlike k-means algorithms, HDBSCAN does not require the user to pre-specify the number of clusters, and HDBSCAN is adept at identifying and handling noise, distinguishing between topics and outliers, which is crucial for maintaining the integrity of the clustered topics.

Traditional methods such as LDA extract topics based on bag-of-words and often fall short when applied to short texts like tweets[68]. BERTopic emerges as particularly advantageous for analyzing data from Twitter and it preserves the semantic structure of the text[62], thus enhancing its ability for short-text analysis compared to traditional models.

We generate 23 topics for 6,660 fact-checked tweets, and 1,073 tweets are not assigned any topic and thus considered as outliers. These outliers are excluded from the process of assigning tweets into control and treatment groups in the delayed feedback (DF) analysis. Most frequent topics with the keywords are shown in the Extended Data Table 5. As shown in Extended Data Table 5, the nine most frequent topics in our dataset include political topics that are known to trigger divisive, polarized reactions in US politics, such as COVID-19 vaccine-related misinformation (Topic 1), election- and politician-related misinformation (Topic 4, 5, 7, 8), policy-related misinformation (Topic 3, 6), and environment and disaster-related misinformation (Topic 9). These topics account for 87.25% of the corrections made through individual tagging and 59.49% of the corrections made through collective tagging. Given the time period of collection, the most frequent topic is about COVID-19 pandemic and vaccination.

## Linguistic Characteristics of Misinformation Tagging Messages

For each misinformation tagging, we analyze the message's toxicity, sentiment, length, reading ease, and delayed response time to provide insights into the qualitative differences between individual and collective misinformation tagging. Extended Data Table 4 shows the descriptive statistics of the following variables.

- *Toxicity*: We apply Google Jigsaw Perspective API to measure the probability that a particular message is toxic (range from 0 to 1[10,69]).
- *Sentiment*: We conduct Vader sentiment analysis to estimate sentiment scores of messages (on a [-1, 1] scale[14,70]). The scale spans from -1, denoting negative sentiment, to 1, denoting positive sentiment.
- *Length*: We measure the length of messages based on the number of characters[14].
- *Reading ease score*: We evaluate the readability of messages with the Flesch-Kincaid Reading Ease score (on a [1,100] scale, where large value indicates easier readability[14,71]). The Flesch–Kincaid reading ease score was transformed into an 8-level categorical variable: "5th grade" for scores 100–90, "6th grade" for 90–80, "7th grade" for 80–70, "8th & 9th grade" for 70–60, "10th to 12th grade" for 60–50, "College" for 50–30, "College graduate" for 30–10, and "Professional" for 10–0.
- *Delayed response time*: We calculate it as the number of days between original tweets and misinformation tagging.

## Robustness Checks

We verify our findings with a battery of robustness checks (see Supplementary Tables 8 and 9). First, there might be concerns that our conclusions about the effects of misinformation tagging on human users may be biased by the presence of bots on Twitter. Many studies have utilized bot detection algorithms to exclude users who are likely to be bots to address this concern[2,5], but others argue these algorithms lead to false negatives (i.e., bots misclassified as humans) and positives (i.e., humans misclassified as bots) that could further bias analyses, even when used cautiously[53]. To mitigate concerns of bot prevalence, we reanalyze our data excluding accounts identified as bots using BotometerLite API. Specifically, using BotometerLite API, we evaluate the likelihood of users in our dataset being bot accounts and remove 360 accounts that have a likelihood higher than 50%[72], which do not meaningfully change our results.

In terms of applying BotometerLite API, some features are missing in our dataset collected with Twitter API 2.0: (1) *default_profile* (whether the user altered the theme or background of their user profile); (2) *profile_use_background_image* (whether the user has a background image or not); and (3) *favourites_count* (number of likes posted by the user, which were only available in Twitter API 1.1). To address this issue, we conduct missing data imputation with the IterativeImputer in Sklearn. We train an imputation model with 90,000 tweets randomly selected in August 2021 from the Twitter Stream Grab (https://archive.org/details/twitterstream). We then

evaluate the model with a held-out sample of 10,000 tweets. The model performance for predicting the missing features is as follows: default_profile at .95 *F-1* score, profile_use_background_image at .90 *F-1* score, favourites_count was .10 $R^2$ value. Finally, we apply this imputation model to recover the missing features in our dataset.

Second, we control for potentially insincere informational activities, such as citing sources of low credibility and intentionally spreading fake news. Some might question whether the increase in political and content diversity is associated with these insincere activities. Put simply, users might be engaging with diverse information that includes misleading claims and conspiracy theories. For each tweet posted by each poster, we measure the credibility of the referenced source. Specifically, we use the binary credibility scores (1=low credibility; 0=medium or high credibility) from the MediaBias/FactCheck database. Our analysis indicates a strong negative correlation between the engagement of low-credibility sources and measures of political diversity ($r$=-.227, $p$<.001) and content diversity ($r$=-.030, $p$<.001). Furthermore, we reassess our data while controlling for credibility of sources, and find that our results remain unaffected. This implies that the increase of diversity in information engagement reflects engagement with a healthier information ecosystem, rather than the reverse.

Third, we attempt to avoid situations in which posters simply criticize distant information without honest consideration by filtering out tweets with negative sentiment. For each tweet posted by each poster, we conduct Vader sentiment analysis to estimate sentiment scores (on a [-1, 1] scale[14,70]). Then we exclude tweets that have a sentimental score lower than 0, which do not meaningfully change our results, except for making the immediate change of content diversity after collective misinformation tagging less significant.

Fourth, we restrict the sample to original posters who have replied to (and thereby read) the individual tags (i.e., fact-checking replies) and remove non-responders. Specifically, out of 6,760 original posters who received individual tags, we remove 4,288 posters who did not reply to the tags, resulting in 2,472 posters. After that, we compare these 2,472 posters with 973 posters who received collective tags. Even after removing the non-responders, we find that results regarding tagging's effects remain consistent. Specifically, as with the complete sample, we identically find that individual tagging causes immediate decrease in political and content diversity in ITS. The gap between individual and collective tagging's effects on the immediate change of political and content diversity did not change in both ITS and DF analyses.

Fifth, we address concerns regarding the possibility of miscoding mentions of "community note." Specifically, we identify all tweets that mention keywords about receiving community notes broadly (i.e., community note, birdwatch, fact-check, factcheck, politifact) within the sample and remove those tweets. To address the concern regarding the effect of replying back to individual taggers, we identify all tweets that reply directly to the individual taggers and remove

them. As shown in Supplementary Tables 10 and 11, we find that the effects on political and content diversity do not meaningfully change in both ITS and DF analyses.

Sixth, to strictly identify PolitiFact links that correct the original posters, we submit original posts, replies containing PolitiFact links, and the cited PolitiFact fact-checking articles to ChatGPT (gpt-4o-2024-05-13). We prompt the model to annotate whether the PolitiFact link was used to correct the original poster rather than support them (see Supplementary Method 7). Consequently, we identify 5,592 PolitiFact links out of 6,760 links (82.72%) as corrective. Subsequently, we limit the sample to the 5,592 links identified by ChatGPT from the individual tagging data, which does not meaningfully alter the results (see Supplementary Method 7).

**References**

1. Lazer, D. M. J. *et al.* The science of fake news. *Science* **359**, 1094–1096 (2018).

2. Vosoughi, S., Roy, D. & Aral, S. The spread of true and false news online. *Science* **359**, 1146–1151 (2018).

3. Loomba, S., de Figueiredo, A., Piatek, S. J., de Graaf, K. & Larson, H. J. Measuring the impact of COVID-19 vaccine misinformation on vaccination intent in the UK and USA. *Nat Hum Behav* **5**, 337–348 (2021).

4. Green, J., Hobbs, W., McCabe, S. & Lazer, D. Online engagement with 2020 election misinformation and turnout in the 2021 Georgia runoff election. *Proc. Natl. Acad. Sci. U. S. A.* **119**, e2115900119 (2022).

5. Bovet, A. & Makse, H. A. Influence of fake news in Twitter during the 2016 US presidential election. *Nat. Commun.* **10**, 7 (2019).

6. Hannak, A., Margolin, D., Keegan, B. & Weber, I. Get back! You don't know me like that: The social mediation of fact checking interventions in Twitter conversations. *Proceedings of the International AAAI Conference on Web and Social Media* **8**, 187–196 (2014).

7. Shin, J. & Thorson, K. Partisan Selective Sharing: The Biased Diffusion of Fact-Checking Messages on Social Media. *J. Commun.* **67**, 233–255 (2017).

8. Pilarski, M., Solovev, K. & Pröllochs, N. Community Notes vs. Snoping: How the Crowd Selects Fact-Checking Targets on Social Media. *arXiv [cs.SI]* (2023).

9. Micallef, N., He, B., Kumar, S., Ahamad, M. & Memon, N. The role of the crowd in countering misinformation: A case study of the COVID-19 infodemic. in *2020 IEEE International Conference on Big Data (Big Data)* (IEEE, 2020). doi:10.1109/bigdata50022.2020.9377956.

10. Mosleh, M. & Rand, D. G. Measuring exposure to misinformation from political elites on Twitter. *Nat. Commun.* **13**, 7144 (2022).

11. Rhodes, S. C. Filter Bubbles, Echo Chambers, and Fake News: How Social Media Conditions Individuals to Be Less Critical of Political Misinformation. *Political Communication* **39**, 1–22 (2022).

12. Bhadani, S. *et al.* Political audience diversity and news reliability in algorithmic ranking. *Nat Hum Behav* **6**, 495–505 (2022).

13. Wojcik, S. *et al.* Birdwatch: Crowd Wisdom and Bridging Algorithms can Inform Understanding and Reduce the Spread of Misinformation. *arXiv [cs.SI]* (2022).

14. Allen, J., Martel, C. & Rand, D. G. Birds of a feather don't fact-check each other: Partisanship and the evaluation of news in Twitter's Birdwatch crowdsourced fact-checking program. in *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* 1–19 (Association for Computing Machinery, New York, NY, USA, 2022).

15. Clayton, K. *et al.* Real solutions for fake news? Measuring the effectiveness of general warnings and fact-check tags in reducing belief in false stories on social media. *Polit. Behav.* **42**, 1073–1095 (2020).

16. Nyhan, B. & Reifler, J. When Corrections Fail: The Persistence of Political Misperceptions. *Political Behavior* **32**, 303–330 (2010).

17. Bail, C. A. *et al.* Exposure to opposing views on social media can increase political polarization. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 9216–9221 (2018).

18. Mosleh, M., Martel, C., Eckles, D. & Rand, D. Perverse Downstream Consequences of Debunking: Being Corrected by Another User for Posting False Political News Increases Subsequent Sharing of Low Quality, Partisan, and Toxic Content in a Twitter Field

Experiment. in *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* 1–13 (Association for Computing Machinery, New York, NY, USA, 2021).

19. Swire-Thompson, B., DeGutis, J. & Lazer, D. Searching for the Backfire Effect: Measurement and Design Considerations. *J. Appl. Res. Mem. Cogn.* **9**, 286–299 (2020).

20. Wood, T. & Porter, E. The Elusive Backfire Effect: Mass Attitudes' Steadfast Factual Adherence. *Political Behavior* **41**, 135–163 (2019).

21. Jiang, S. & Wilson, C. Linguistic Signals under Misinformation and Fact-Checking: Evidence from User Comments on Social Media. *Proc. ACM Hum.-Comput. Interact.* **2**, 1–23 (2018).

22. Masullo, G. M. & Kim, J. Exploring 'Angry' and 'Like' Reactions on Uncivil Facebook Comments That Correct Misinformation in the News. *Digital Journalism* **9**, 1103–1122 (2021).

23. Page, S. E. *The Difference: How the Power of Diversity Creates Better Groups, Firms, Schools, and Societies - New Edition*. (Princeton University Press, 2008).

24. Page, S. E. *The Diversity Bonus: How Great Teams Pay Off in the Knowledge Economy*. (Princeton University Press, 2019).

25. Saeed, M., Traub, N., Nicolas, M., Demartini, G. & Papotti, P. Crowdsourced Fact-Checking at Twitter: How Does the Crowd Compare With Experts? in *Proceedings of the 31st ACM International Conference on Information & Knowledge Management* 1736–1746 (Association for Computing Machinery, New York, NY, USA, 2022).

26. Allen, J., Arechar, A. A., Pennycook, G. & Rand, D. G. Scaling up fact-checking using the wisdom of crowds. *Sci Adv* **7**, eabf4393 (2021).

27. Martel, C., Allen, J., Pennycook, G. & Rand, D. G. Crowds Can Effectively Identify

Misinformation at Scale. *Perspect. Psychol. Sci.* 17456916231190388 (2023).

28. Twitter. Values. *Community Notes Guide*

    https://communitynotes.twitter.com/guide/en/contributing/values (2023).

29. Jamieson, K. H. & Cappella, J. N. *Echo Chamber: Rush Limbaugh and the Conservative Media Establishment*. (Oxford University Press, 2008).

30. Muller, D. Democracy under strain. (2021) doi:10.1007/978-3-030-76761-7_2.

31. Törnberg, P. Echo chambers and viral misinformation: Modeling fake news as complex contagion. *PLoS One* **13**, e0203958 (2018).

32. Del Vicario, M. *et al.* The spreading of misinformation online. *Proceedings of the National Academy of Sciences* **113**, 554–559 (2016).

33. Flamino, J. *et al.* Political polarization of news media and influencers on Twitter in the 2016 and 2020 US presidential elections. *Nat Hum Behav* **7**, 904–916 (2023).

34. Samory, M., Abnousi, V. K. & Mitra, T. Characterizing the Social Media News Sphere through User Co-Sharing Practices. *ICWSM* **14**, 602–613 (2020).

35. Bessi, A. *et al.* Homophily and polarization in the age of misinformation. *Eur. Phys. J. Spec. Top.* **225**, 2047–2059 (2016).

36. Cinelli, M., De Francisci Morales, G., Galeazzi, A., Quattrociocchi, W. & Starnini, M. The echo chamber effect on social media. *Proc. Natl. Acad. Sci. U. S. A.* **118**, (2021).

37. Bode, L. Political News in the News Feed: Learning Politics from Social Media. *Mass Communication and Society* **19**, 24–48 (2016).

38. Conover, M. D., Goncalves, B., Ratkiewicz, J., Flammini, A. & Menczer, F. Predicting the Political Alignment of Twitter Users. in *2011 IEEE Third International Conference on Privacy, Security, Risk and Trust and 2011 IEEE Third International Conference on Social*

*Computing* 192–199 (IEEE, 2011).

39. Luzsa, R. & Mayr, S. False consensus in the echo chamber: Exposure to favorably biased social media news feeds leads to increased perception of public support for own opinions. *Cyberpsychology (Brno)* **15**, (2021).

40. Cooke, N. A. *Fake News and Alternative Facts: Information Literacy in a Post-Truth Era*. (American Library Association, 2018).

41. DellaPosta, D., Shi, Y. & Macy, M. Why Do Liberals Drink Lattes? *AJS* **120**, 1473–1511 (2015).

42. Mutz, D. C. & Rao, J. S. The Real Reason Liberals Drink Lattes. *PS Polit. Sci. Polit.* **51**, 762–767 (2018).

43. Shi, F., Shi, Y., Dokshin, F. A., Evans, J. A. & Macy, M. W. Millions of online book co-purchases reveal partisan differences in the consumption of science. *Nature Human Behaviour* **1**, 1–9 (2017).

44. Milbauer, J., Mathew, A. & Evans, J. A. Aligning multidimensional worldviews and discovering ideological differences. *Empir Method Nat Lang Process* 4832–4845 (2021).

45. Gallotti, R., Valle, F., Castaldo, N., Sacco, P. & De Domenico, M. Assessing the risks of 'infodemics' in response to COVID-19 epidemics. *Nat Hum Behav* **4**, 1285–1293 (2020).

46. Anderson, A., Maystre, L., Anderson, I., Mehrotra, R. & Lalmas, M. Algorithmic Effects on the Diversity of Consumption on Spotify. in *Proceedings of The Web Conference 2020* 2155–2165 (Association for Computing Machinery, New York, NY, USA, 2020).

47. Reimers, N. & Gurevych, I. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *arXiv [cs.CL]* (2019).

48. Srinivasan, K. B., Danescu-Niculescu-Mizil, C., Lee, L. & Tan, C. Content Removal as a

Moderation Strategy: Compliance and Other Outcomes in the ChangeMyView Community. *Proc. ACM Hum.-Comput. Interact.* **3**, 1–21 (2019).

49. Thaler, R. H. & Sunstein, C. R. *Nudge: The Final Edition*. (Penguin, 2021).

50. Jones-Jang, S. M., Mortensen, T. & Liu, J. Does Media Literacy Help Identification of Fake News? Information Literacy Helps, but Other Literacies Don't. *Am. Behav. Sci.* **65**, 371–388 (2021).

51. Becker, J., Porter, E. & Centola, D. The wisdom of partisan crowds. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 10717–10722 (2019).

52. Shi, F., Teplitskiy, M., Duede, E. & Evans, J. A. The wisdom of polarized crowds. *Nat Hum Behav* **3**, 329–336 (2019).

53. Rauchfleisch, A. & Kaiser, J. The False positive problem of automatic bot detection in social science research. *PLoS One* **15**, e0241045 (2020).

54. Grinberg, N., Joseph, K., Friedland, L., Swire-Thompson, B. & Lazer, D. Fake news on Twitter during the 2016 U.S. presidential election. *Science* **363**, 374–378 (2019).

55. Ledford, H. Researchers scramble as Twitter plans to end free data access. *Nature Publishing Group UK* http://dx.doi.org/10.1038/d41586-023-00460-z (2023) doi:10.1038/d41586-023-00460-z.

56. Community Notes. Beginning today, Community Notes are visible around the world 🌍🌏🌎. *Twitter* https://twitter.com/CommunityNotes/status/1601753552476438528 (2022).

57. Wang, Z. *et al.* Demographic Inference and Representative Population Estimates from Multilingual Social Media Data. in *The World Wide Web Conference* 2056–2067 (Association for Computing Machinery, New York, NY, USA, 2019).

58. Bagrow, J. P., Liu, X. & Mitchell, L. Information flow reveals prediction limits in online social activity. *Nat Hum Behav* **3**, 122–128 (2019).

59. Rao, A., Morstatter, F. & Lerman, K. Retweets Amplify the Echo Chamber Effect. in *Proceedings of the 2023 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining* 30–37 (Association for Computing Machinery, New York, NY, USA, 2024).

60. Barberá, P., Jost, J. T., Nagler, J., Tucker, J. A. & Bonneau, R. Tweeting From Left to Right: Is Online Political Communication More Than an Echo Chamber? *Psychol. Sci.* **26**, 1531–1542 (2015).

61. Di Giovanni, M. & Brambilla, M. Exploiting Twitter as Source of Large Corpora of Weakly Similar Pairs for Semantic Sentence Embeddings. in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing* 9902–9910 (Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 2021).

62. Nguyen, D. Q., Vu, T. & Nguyen, A. T. BERTweet: A pre-trained language model for English Tweets. *arXiv [cs.CL]* (2020).

63. Waller, I. & Anderson, A. Generalists and Specialists: Using Community Embeddings to Quantify Activity Diversity in Online Platforms. in *The World Wide Web Conference* 1954–1964 (Association for Computing Machinery, New York, NY, USA, 2019).

64. Grootendorst, M. BERTopic: Neural topic modeling with a class-based TF-IDF procedure. *arXiv [cs.CL]* (2022).

65. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv [stat.ML]* (2018).

66. Houle, M. E., Kriegel, H.-P., Kröger, P., Schubert, E. & Zimek, A. Can Shared-Neighbor

Distances Defeat the Curse of Dimensionality? in *Scientific and Statistical Database Management* 482–500 (Springer Berlin Heidelberg, 2010).

67. Ester, M., Kriegel, H.-P., Sander, J. & Xu, X. A density-based algorithm for discovering clusters in large spatial databases with noise. in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining* 226–231 (AAAI Press, 1996).

68. Tierney, G., Bail, C. & Volfovsky, A. Author Clustering and Topic Estimation for Short Texts. *arXiv [cs.IR]* (2021).

69. Lees, A. *et al.* A New Generation of Perspective API: Efficient Multilingual Character-level Transformers. in *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* 3197–3207 (Association for Computing Machinery, New York, NY, USA, 2022).

70. Hutto, C. & Gilbert, E. VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text. *ICWSM* **8**, 216–225 (2014).

71. Peter Kincaid, J., Fishburne, R. P., Jr, Rogers, R. L. & Chissom, B. S. Derivation Of New Readability Formulas (Automated Readability Index, Fog Count And Flesch Reading Ease Formula) For Navy Enlisted Personnel. (1975).

72. Yang, K.-C., Ferrara, E. & Menczer, F. Botometer 101: social bot practicum for computational social scientists. *SIAM J. Sci. Comput.* **5**, 1511–1528 (2022).

73. Li, F., Morgan, K. L. & Zaslavsky, A. M. Balancing Covariates via Propensity Score Weighting. *J. Am. Stat. Assoc.* **113**, 390–400 (2018).

74. A. Smith, J. & E. Todd, P. Does matching overcome LaLonde's critique of nonexperimental estimators? *J. Econom.* **125**, 305–353 (2005).

# Extended Data

**Extended Data Table 1. Interrupted Time Series (ITS) Model Results for Political and Content Diversity**

| Outcome | Political diversity (%) | | | Content Diversity (z) | | |
|---|---|---|---|---|---|---|
| Type of misinformation tagging | Individual | Collective | Difference (Collective - Individual) | Individual | Collective | Difference (Collective - Individual) |
| Slope before posting the tweet | .237*** (.057) | .309* (.137) | .072 (.148) | .007*** (.002) | .003 (.004) | -.005 (.004) |
| Immediate intercept change after misinformation tagging | -1.009*** (.223) | .270 (.558) | 1.279* (.601) | -.030*** (.006) | .040** (.015) | .070*** (.016) |
| Slope after misinformation tagging | .087 (.054) | -.049 (.145) | -.136 (.155) | -.003+ (.001) | -.011** (.004) | -.008* (.004) |
| Slope change (After - Before) | -.150+ (.080) | -.358+ (.199) | -.208 (.215) | -.010*** (.002) | -.014** (.005) | -.004 (.006) |
| Observations | 424,969 | | | | | |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. We multiply political diversity by 100 to interpret the estimates as absolute percentage point changes. We normalize content diversity to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests. More details can be found in Methods: Interrupted Time Series (ITS) Analysis.*

**Extended Data Table 2. Delayed feedback (DF) model results for political and content diversity**

| Outcome | Political diversity | | | Content Diversity | | |
|---|---|---|---|---|---|---|
| Type of misinformation tagging | Individual | Collective | Difference (Collective - Individual) | Individual | Collective | Difference (Collective - Individual) |
| Difference in Pre-Post Change (Treatment - Control) | -5.886** (1.911) | 1.219 (3.059) | 7.105* (3.589) | .018 (.041) | .274*** (.066) | .256** (.077) |
| Observations | 8,901 | | | | | |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. We multiply political diversity by 100 to interpret the estimates as absolute percentage point changes. We normalize content diversity to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests. More details can be found in Methods: Delayed Feedback (DF) Analysis.*

**Extended Data Table 3. Descriptive Statistics of Linguistic Characteristics.**

| Type of Misinformation Tagging | Individual | Collective | Significance of Difference (Collective - Individual) | Observations |
|---|---|---|---|---|
| Toxicity (0~1) | .139 (.183) | .076 (.078) | t(7496)=9.86*** | 7,498 |
| Sentiment (-1~1) | -.082 (.439) | -.050 (.382) | t(7731)=-2.14* | 7,733 |
| Length (# Characters) | 179.310 (115.951) | 288.873 (135.276) | t(7731)=-26.95*** | 7,733 |
| Reading Ease | | | | |
| 5th grade | .39% | .51% | | |
| 6th grade | 3.17% | 1.37% | | |
| 7th grade | 9.04% | 2.39% | | |
| 8-9th grade | 16.01% | 8.19% | | |
| 10-12th grade | 17.86% | 11.77% | $\chi 2(7)=155.32$*** | 6,493 |
| College | 31.56% | 34.64% | | |
| College graduate | 16.03% | 27.82% | | |
| Professional | 5.94% | 13.31% | | |
| Delay (# Days) | 3.04 (42.330) | 6.322 (60.315) | t(7731)=-2.13* | 7,733 |

Notes: ***p<.001 **p<.01 *p<.05 +p<.1. The values present the mean value or proportion with the standard deviation in parentheses. Significance of difference is estimated by independent two-sample t-test (continuous variable) or chi-square test (categorical variable). Statistical significance levels (P values) are derived from two-sided tests.

**Extended Data Table 4. Control Analyses for Linguistic Characteristics**

| Outcome | Political diversity (%) | | | | | |
|---|---|---|---|---|---|---|
| Controls | Initial results (N=424,969) | Non-toxic tags (N=381,552) | Neutral tags (N=160,451) | Short tags (N=405,829) | Easy tags (N=167,766) | Quick tags (N=364,024) |
| **Individual** | | | | | | |
| Slope before | .237*** (.057) | .205** (.061) | .312** (.096) | .258*** (.058) | .240*** (.064) | .256*** (.061) |
| Immediate change | -1.009*** (.223) | -.969*** (.237) | -1.530*** (.377) | -1.046*** (.225) | -1.131*** (.252) | -1.176*** (.235) |
| Slope after | .087 (.054) | .076 (.058) | .169+ (.093) | .092+ (.055) | .114+ (.061) | .121* (.058) |
| Slope change | -.150+ (.080) | -.128 (.084) | -.143 (.135) | -.166* (.080) | -.126 (.090) | -.135 (.084) |
| **Collective** | | | | | | |
| Slope before | .309* (.137) | .352* (.140) | .184 (.201) | .421** (.154) | .183 (.157) | .228 (.186) |
| Immediate change | .270 (.558) | .275 (.570) | .280 (.827) | .024 (.625) | .189 (.645) | .399 (.744) |
| Slope after | -.049 (.145) | -.052 (.148) | .221 (.213) | .070 (.161) | .092 (.169) | -.031 (.192) |
| Slope change | -.358+ (.199) | -.405* (.204) | .037 (.292) | -.351 (.223) | -.091 (.230) | -.259 (.268) |
| **Individual vs. collective** | | | | | | |
| Slope before | .072 (.148) | .148 (.153) | -.128 (.223) | .163 (.164) | -.057 (.170) | -.028 (.196) |
| Immediate change | 1.279* (.601) | 1.245* (.617) | 1.809* (.909) | 1.071 (.664) | 1.320+ (.692) | 1.575* (.780) |
| Slope after | -.136 (.155) | -.129 (.159) | .052 (.233) | -.022 (.171) | -.022 (.179) | -.152 (.200) |
| Slope change | -.208 (.215) | -.276 (.220) | .181 (.322) | -.185 (.237) | .035 (.247) | -.124 (.281) |
| Outcome | Content diversity (z) | | | | | |
| Robustness checks | Initial results (N=424,969) | Non-toxic tags (N=381,552) | Neutral tags (N=160,451) | Short tags (N=405,829) | Easy tags (N=167,766) | Quick tags (N=358,201) |
| **Individual** | | | | | | |
| Slope before | .007*** (.002) | .008*** (.002) | .008** (.002) | .008*** (.002) | .006*** (.002) | .008*** (.002) |
| Immediate change | -.030*** (.006) | -.031*** (.006) | -.063*** (.010) | -.030*** (.006) | -.027*** (.007) | -.037*** (.006) |
| Slope after | -.003+ (.001) | -.003* (.002) | .001 (.002) | -.003* (.001) | -.001 (.002) | -.003+ (.002) |
| Slope change | -.010*** (.002) | -.011*** (.002) | -.007* (.003) | -.011*** (.002) | -.007** (.002) | -.011*** (.002) |
| **Collective** | | | | | | |
| Slope before | .003 (.004) | .002 (.004) | .020*** (.005) | .005 (.004) | .002 (.004) | -.007 (.005) |
| Immediate change | .040** (.015) | .044** (.015) | .011 (.021) | .043* (.017) | .040* (.017) | .069** (.020) |
| Slope after | -.011** (.004) | -.011** (.004) | -.021*** (.005) | -.013** (.004) | -.017*** (.004) | -.023*** (.005) |
| Slope change | -.014** (.005) | -.013* (.005) | -.041*** (.007) | -.018** (.006) | -.019** (.006) | -.016* (.007) |
| **Individual vs. collective** | | | | | | |
| Slope before | -.005 (.004) | -.006 (.004) | .011* (.006) | -.003 (.004) | -.004 (.004) | -.015** (.005) |
| Immediate change | .070*** (.016) | .075*** (.016) | .074** (.023) | .072*** (.018) | .067*** (.018) | .106*** (.021) |
| Slope after | -.008* (.004) | -.008* (.004) | -.022*** (.006) | -.010* (.005) | -.016** (.005) | -.020*** (.005) |
| Slope change | -.004 (.006) | -.003 (.006) | -.034*** (.008) | -.007 (.006) | -.012+ (.006) | -.005 (.008) |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Observations refer to the number of tweets that received individual or collective tags. Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests.*

**Extended Data Table 5. Most Frequent Topics Corrected by Individual and Collective Tagging**

| Topic ID | Top keywords | Frequency (N) | Total (%) | Individual (%) | Collective (%) |
|---|---|---|---|---|---|
| 1 | covid, vaccine, vaccines, vaccinated | 1,753 | 26.32 | 28.42 | 10.26 |
| 2 | fact, just, lies, true | 1,295 | 19.44 | 20.76 | 9.35 |
| 3 | gun, state, state lines, lines | 554 | 8.32 | 8.42 | 7.53 |
| 4 | projection, lie republican, monster, fyi | 491 | 7.37 | 7.91 | 3.25 |
| 5 | tax, biden, inflation, bernie | 487 | 7.31 | 7.67 | 4.55 |
| 6 | ukraine, russia, biden, putin | 218 | 3.27 | 2.84 | 6.62 |
| 7 | votes, election, voters, vote | 216 | 3.24 | 3.6 | 0.52 |
| 8 | hillary, thomas, clinton, just | 199 | 2.99 | 3.06 | 2.47 |
| 9 | turkey, earthquake, climate, al gore | 196 | 2.94 | 1.38 | 14.94 |

*Notes*: Top keywords have been identified by counting the most frequent words within tweets corresponding to each topic, following the exclusion of stopwords. Total (%) column indicates the proportion of a particular topic in relation to all corrected misinformation. Individual (%) column indicates the proportion of a particular topic out of the misinformation corrected through individual tagging. Collective (%) column indicates the proportion of a particular topic out of the misinformation corrected through collective tagging.

**Extended Data Table 6. Control Analyses on Topics**

| Outcome | Political diversity (%) | | |
|---|---|---|---|
| Controls | Initial results (N = 424,969) | Limiting topics (N = 301,380) | Propensity score weighting (N = 301,380) |
| **Individual** | | | |
|   Slope before | .237*** (.057) | .214*** (.061) | .135* (.065) |
|   Immediate change | -1.009*** (.223) | -.893*** (.238) | -.703** (.260) |
|   Slope after | .087 (.054) | .068 (.058) | .037 (.065) |
|   Slope change | -.150+ (.080) | -.146+ (.085) | -.097 (.093) |
| **Collective** | | | |
|   Slope before | .309* (.137) | .239 (.157) | -.274 (.269) |
|   Immediate change | .270 (.558) | .648 (.636) | 1.677 (1.082) |
|   Slope after | -.049 (.145) | -.036 (.160) | -.218 (.270) |
|   Slope change | -.358+ (.199) | -.274 (.224) | .056 (.382) |
| **Individual vs. collective** | | | |
|   Slope before | .072 (.148) | .025 (.168) | -.409 (.277) |
|   Immediate change | 1.279* (.601) | 1.542* (.679) | 2.380* (1.112) |
|   Slope after | -.136 (.155) | -.104 (.170) | -.255 (.277) |
|   Slope change | -.208 (.215) | -.129 (.240) | .154 (.394) |
| **Outcome** | **Content diversity (z)** | | |
| Controls | Initial results (N = 424,969) | Limiting topics (N = 301,380) | Propensity score weighting (N = 301,380) |
| **Individual** | | | |
|   Slope before | .007*** (.002) | .008*** (.002) | .010*** (.002) |
|   Immediate change | -.030*** (.006) | -.039*** (.006) | -.044*** (.007) |
|   Slope after | -.003+ (.001) | .001 (.002) | .000 (.002) |
|   Slope change | -.010*** (.002) | -.007** (.002) | -.010*** (.002) |
| **Collective** | | | |
|   Slope before | .003 (.004) | .001 (.004) | .003 (.006) |
|   Immediate change | .040** (.015) | .045** (.017) | .004 (.022) |
|   Slope after | -.011** (.004) | -.015*** (.004) | -.002 (.005) |
|   Slope change | -.014** (.005) | -.015** (.006) | -.005 (.008) |
| **Individual vs. collective** | | | |
|   Slope before | -.005 (.004) | -.008+ (.004) | -.007 (.006) |
|   Immediate change | .070*** (.016) | .084*** (.018) | .048* (.023) |
|   Slope after | -.008* (.004) | -.016*** (.004) | -.002 (.005) |
|   Slope change | -.004 (.006) | -.008 (.006) | .005 (.008) |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests.*

**Extended Data Table 7. Control Analyses on Posters' and Taggers/Voters' Political Belonging**

| Outcome | Political diversity (%) | | Content diversity (z) | |
|---|---|---|---|---|
| Controls | Initial results (N = 424,969) | Control political belonging (N = 241,681) | Initial results (N = 424,969) | Control political belonging (N = 241,681) |
| **Individual** | | | | |
| Slope before | .237*** (.057) | .282*** (.080) | .007*** (.002) | .011*** (.002) |
| Immediate change | -1.009*** (.223) | -1.231*** (.310) | -.030*** (.006) | -.055*** (.008) |
| Slope after | .087 (.054) | .143+ (.077) | -.003+ (.001) | .001 (.002) |
| Slope change | -.150+ (.080) | -.139 (.112) | -.010*** (.002) | -.010*** (.003) |
| **Collective** | | | | |
| Slope before | .309* (.137) | .262 (.194) | .003 (.004) | .008 (.005) |
| Immediate change | .270 (.558) | .548 (.789) | .040** (.015) | .021 (.020) |
| Slope after | -.049 (.145) | .119 (.203) | -.011** (.004) | -.019*** (.005) |
| Slope change | -.358+ (.199) | -.143 (.281) | -.014** (.005) | -.027*** (.007) |
| **Individual vs. collective** | | | | |
| Slope before | .072 (.148) | -.020 (.210) | -.005 (.004) | -.003 (.005) |
| Immediate change | 1.279* (.601) | 1.780* (.848) | .070*** (.016) | .076** (.022) |
| Slope after | -.136 (.155) | -.024 (.217) | -.008* (.004) | -.020*** (.006) |
| Slope change | -.208 (.215) | -.004 (.302) | -.004 (.006) | -.017* (.008) |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests.*

**Extended Data Table 8. Control Analyses on Poster's Popularity**

| Outcome | Political diversity (%) | | Content diversity (z) | |
|---|---|---|---|---|
| Controls | Initial results (N = 424,969) | Removing popular posters (N = 240,653) | Initial results (N = 424,969) | Removing popular posters (N = 240,653) |
| **Individual** | | | | |
|   Slope before | .237*** (.057) | .400*** (.075) | .007*** (.002) | .009*** (.002) |
|   Immediate change | -1.009*** (.223) | -1.904*** (.285) | -.030*** (.006) | -.043*** (.008) |
|   Slope after | .087 (.054) | .187** (.070) | -.003+ (.001) | -.002 (.002) |
|   Slope change | -.150+ (.080) | -.213* (.104) | -.010*** (.002) | -.011*** (.003) |
| **Collective** | | | | |
|   Slope before | .309* (.137) | -.184 (.357) | .003 (.004) | .016 (.010) |
|   Immediate change | .270 (.558) | 1.708 (1.393) | .040** (.015) | .037 (.040) |
|   Slope after | -.049 (.145) | .243 (.357) | -.011** (.004) | -.007 (.010) |
|   Slope change | -.358+ (.199) | .427 (.506) | -.014** (.005) | -.023 (.015) |
| **Individual vs. collective** | | | | |
|   Slope before | .072 (.148) | -.584 (.365) | -.005 (.004) | .007 (.011) |
|   Immediate change | 1.279* (.601) | 3.612* (1.422) | .070*** (.016) | .081* (.041) |
|   Slope after | -.136 (.155) | .056 (.364) | -.008* (.004) | -.005 (.011) |
|   Slope change | -.208 (.215) | .640 (.516) | -.004 (.006) | -.012 (.015) |

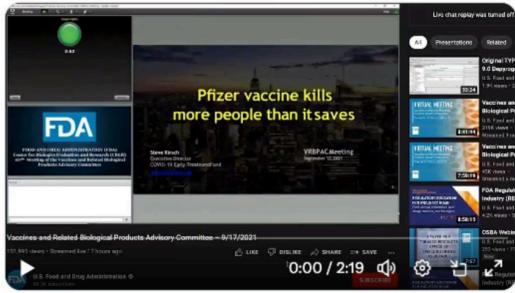*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests.*

**Extended Data Fig. 1. Examples of Individual and Collective Misinformation Tagging. a,** An instance of an individual misinformation tagging regarding COVID-19 vaccination through a reply. **b,** An instance of a collective misinformation tagging regarding COVID-19 vaccination, which is shown above other users' replies. Usernames and profiles have been anonymized. Following Twitter's privacy policy, we provide the manually rephrased tweets to ensure that user identities remain confidential.

## a  Individual Misinformation Tagging

**Poster A**
@poster_a

The FDA confirmed that the vaccine is killing more people than it's saving.



**Misinformation-tagger**
@misinformation-tagger

He made this claim during a public comment period to the FDA. His comment is inaccurate, and he is not affiliated with the FDA. [Link to the source]

## b  Collective Misinformation Tagging

**Poster B**
@poster_b

The ex-CEO of Pfizer states that the vaccine will destroy you.



**Readers added context they thought people might want to know**

Michael Yeadon was not the CEO of Pfizer. He was a scientist and vice president who worked on allergy and respiratory research. He did not work with vaccines, and he left Pfizer in 2011.
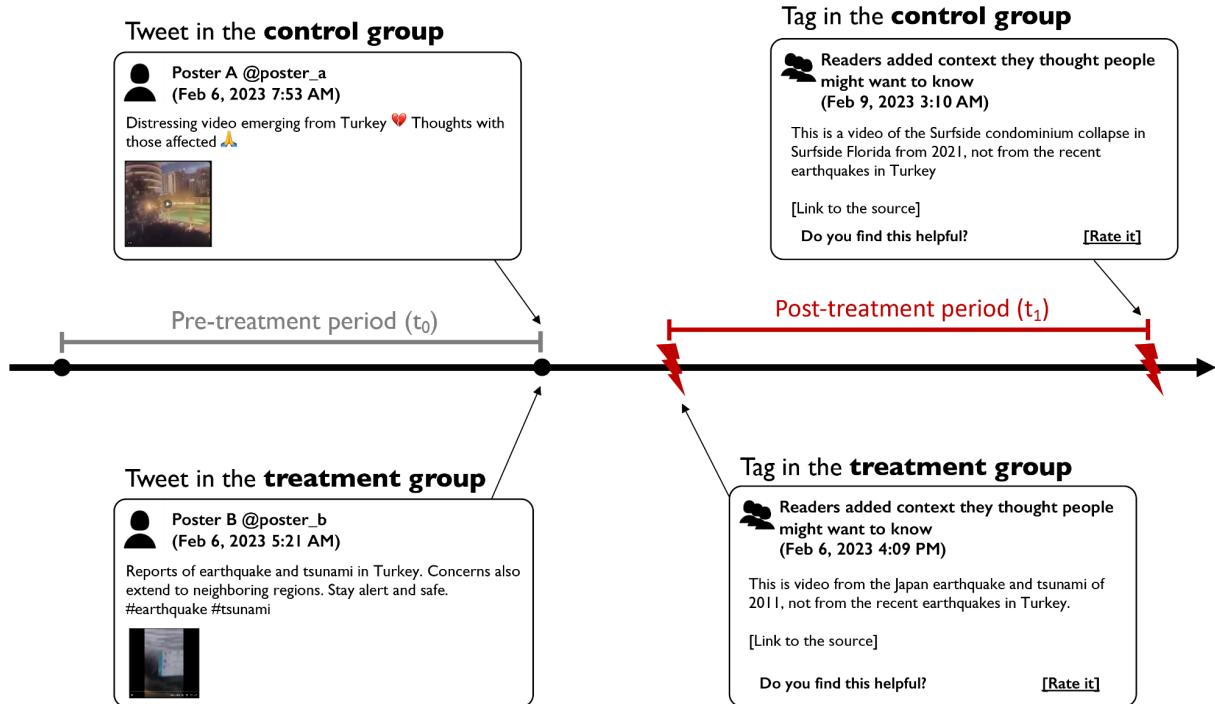
[Link to the source]

**Do you find this helpful?**          **[Rate it]**

**Another-replier**
@another-replier

**Extended Data Fig. 2. An example of a pair of matched tweets and tags in delayed feedback (DF) analysis.** Here, a pair of corrected tweets shows similar characteristics: they are about the same topic (i.e., the earthquake in Türkiye) and were posted at approximately the same time (tweet in the control group: February 6, 2023, 7:53 AM; tweet in the treatment group: February 6, 2023, 5:21 AM). Nevertheless, they were corrected at very different times (tag in the control group: February 9, 2023, 3:10 AM; tag in the treatment group: February 6, 2023, 4:09 PM), constituting a gap of 53 hours.

# Supplementary Method 1: Additional Analyses regarding the Effects of Collective Tagging on Content Diversity

Despite the steepness of the slope following collective tagging, our analysis indicates that content diversity does not drop below the pre-tagged period. Specifically, we estimate the significance of pairwise differences between weeks (i.e., comparing Week 5 to Week -5, Week 5 to Week -4, and so on, up to Week 5 versus Week 0). As shown in Supplementary Table 12, although content diversity approaches the level of content diversity observed in the pre-tagged period, it does not significantly decrease below that. For instance, when comparing the content diversity in Week 5 with that in Week -5 (see Fig. 2A), the difference in content diversity between these two points is not significant (p=.929). Similarly, the comparison of content diversity between Week 5 and Week -1 (right before tagging, as depicted in Fig. 2A) shows no significant difference (p=.398).

Nevertheless, some might be concerned that content diversity might eventually drop below the pre-tagged period after Week 5. However, after Week 5, we find that the slope becomes almost flat. Specifically, we find that the slope of change in content diversity between Week 5 and Week 8 becomes close to zero ($\beta$=.002, *Standard Error=.009, p=.799, N*=12,056). This means that the original posters maintain their position after reaching the initial level of content diversity at Week 5.

# Supplementary Method 2: Political and Content Proximity to Misinformation Taggers

We conduct additional analyses to examine how misinformation posters' political and content proximity to taggers changes after misinformation tagging occurs. Since we could identify only the Twitter accounts of individual misinformation taggers and not those of collective misinformation taggers, our analysis uses data from individual misinformation tagging.

First, we measure political stances of taggers in individual tagging by averaging the MediaBias/FactCheck scores (see Method: Political Diversity) of sources referenced in their historical tweets. This span includes two months before and two months after the posting of tagged tweets. Taggers who predominantly cite left-leaning media are considered "left," while those who cite right-leaning media are considered "right." Specifically, taggers with negative average political stance scores are categorized as left, while those with positive scores are categorized as right.

Second, we measure the political proximity between each poster's tweet and the misinformation tagger. Specifically, we assign a binary value to denote the political proximity: 1 (high proximity) is assigned if the poster cites a source with the same political stance as the tagger; 0 (low proximity) is assigned if not. Similarly, we assess the tweet's content proximity to the tagger. For this, we calculate the cosine similarity between the embedding of each poster's tweet and the tagger's embedding to indicate the content proximity between a specific tweet and the tagger.

As shown in Supplementary Fig. 5 and Supplementary Table 13, we estimate interrupted time series (ITS) models for political and content proximity to taggers. Before individual misinformation tagging, posters increase the political proximity ($\beta$=.146, $p$=.015) and content proximity ($\beta$=.003, $p$=.015) to taggers. Immediately after tagging, posters significantly decrease the political proximity ($\beta$=-.536, $p$=.022), although they slightly increase content diversity ($\beta$=.007, $p$=.206). After tagging, posters collapse in both political proximity ($\beta$=-.203, $p$=.015) and content proximity ($\beta$=-.005, $p$=.010) to taggers.

# Supplementary Method 3: Propensity Score Weighting

To control for the distribution of topics, we implement propensity score weighting (PSW) methods[73,74]. PSW balances the distribution of topics between individual and collective tagging by oversampling tweets corrected by collective tagging if they are topically more likely to be addressed by individual tagging (e.g., if a specific topic is predominantly corrected by individual tagging, the model increases the weight of collective tagging for that topic). Statistically, this technique removes the selection bias based on topics (i.e., different probability of receiving individual versus collective tagging across different topics). This approach helps to control not only the differences between political and non-political misinformation but also the more nuanced topical variations between tweets corrected by each type of tagging.

Specifically, let $C$ is an indicator of the type of misinformation tagging where 0 represents individual misinformation tagging and 1 represents collective misinformation tagging. $X_i$ is an indicator of whether a particular tweet belongs to a topic $i$. We fit logistic regression that predicts the probability of receiving collective tags, instead of individual tags, on $X_1, X_2, ..., X_9$.

$$log(C / (1-C)) = \beta_0 + \Sigma\beta_i X_i + \varepsilon \qquad (1)$$

Let $P$ denote the probability of receiving collective tags based on this model. Then, we estimate the inverse probability weight *IPW* as follows:

$$IPW = 1/P \text{ if } C = 1$$
$$IPW = 1/(1-P) \text{ if } C = 0 \qquad (2)$$

Finally, we apply inverse probability weights (IPW) to adjust our interrupted time series (ITS) models. By applying these weights, we control for any differences in the outcomes between individual and collective tags introduced by topical differences of the corrected misinformation.

# Supplementary Method 4: Political Stances of Taggers and Voters

We analyze the distribution of political stance among taggers and voters in each condition. Taggers refer to users who write individual tags (i.e., fact-checking replies), and voters refer to users who vote on the exposure of collective tags. Note that we have voters only in collective tagging because individual tags are exposed without votes. Here, we have four conditions: individual tags correcting left-leaning users, individual tags correcting right-leaning users, collective tags correcting left-leaning users, collective tags correcting right-leaning users. For each condition, the distribution of taggers' and voters' political stances are shown in Supplementary Table 3.

In every condition, we find that left-leaning taggers and voters are more prevalent. Compared to individual tagging, community notes are more balanced as they consider more right-leaning perspectives. Still, the majority of votes come from left-leaning voters. Based on this imbalance, we suspect that right-leaning voters are more likely to be "filtered" due to their low-quality contributions in the past (see Supplementary Method 5: Potential Backlash Among Voters in Community Notes). Therefore, right-leaning posters are likely to be corrected by someone of different political stances (i.e., left-leaning taggers and voters) in both individual and collective tagging.

We measure taggers' and voters' proxy political stances as follows. First, political stances of taggers in individual tagging is measured by averaging the MediaBias/FactCheck scores (see Method: Political Diversity) of sources referenced in their historical tweets. This span includes two months before and two months after the posting of tagged tweets. Taggers who predominantly cite left-leaning media are considered "left," while those who cite right-leaning media are considered "right." Specifically, taggers with negative average political stance scores are categorized as left, while those with positive scores are categorized as right.

Second, political stances of voters in collective tagging are measured by their activities in Community Notes. We were not able to access these anonymous voters' historical tweets to measure their political stances, but we were able to fully access their activities in Community Notes. Allen and colleagues (2024) suggest that voters' activities in Community Notes can be a proxy of their political stances, documenting a strong correlation between voters' political stances and their activities[14]. Therefore, we assume that voters are left-leaning if their prior activities in Community Notes show they "tag" or "vote" on right-leaning politicians' tweets as misinformation more frequently than those of left-leaning politicians. Conversely, we assume that voters are right-leaning if their prior records indicate they tag or vote on left-leaning politicians' tweets as misinformation more often than those of right-leaning politicians.

# Supplementary Method 5: Potential Backlash Among Voters in Community Notes

Some may concern that voters from one political party can engage in backlash in Community Notes, obstructing the correction of misinformation widespread within that party (e.g., misinformation about vaccination or the 2020 election). However, we find that backlash is an exception, rather than the norm. Community notes evaluate cross-perspective agreements[13], but they do so only after filtering out voters likely to backlash and oppose valid fact-checks. Specifically, the algorithm assesses whether users' votes are not aligned with the crowd's final decisions by more than 33% and removes their votes[13]. In other words, if a voter casted 10 votes in community notes, and more than 3 votes were found to be different from the crowd's decisions, their votes are excluded by the algorithm. Due to the high threshold to be eligible voters, users who misuse their votes to obstruct valid corrections are likely to be disregarded.

Indeed, almost all eligible voters are found not to *veto* vaccine- and election-related collective tags. We first have identified 3,558 vaccine-related tags containing the word "vaccine," and 2,597 election-related tags containing the word "election" in Community Notes. Among 4,535 voters who passed algorithmic filtering in Community Notes, 96.3% have never opposed vaccine-related tags, and 98.54% have not opposed election-related tags approved by the crowd. This suggests that community notes require cross-perspective agreements within a moderate (not extreme) political population who rarely disagree with valid fact-checks against blatant misinformation, such as COVID-19 and election misinformation.

# Supplementary Method 6: Inter-Rater Reliability for Political Stances Derived from MediaBias/FactCheck

We refer to the political stance scores from the MediaBias/FactCheck (MBFC) database (https://mediabiasfactcheck.com/) to assess each tweet's political stance. Specifically, for tweets with political scores ranging from -1 (extreme left) to 0, we label them as "left", and for tweets with political scores ranging from 0 to 1 (extreme right), we label them as "right".

To address the concerns regarding potential biases from the MediaBias/FactCheck scores, we estimate the inter-rater reliability of MBFC scores with an alternative database from Allsides.com, which labels 445 websites as "left," "leaning left," "leaning right," or "right" (we exclude websites labeled as "center" or "mixed")[33]. We map these labels to "left" (including "left" and "leaning left") and "right" (including "right" and "leaning right"). Then we calculate the inter-rater reliability scores between MBFC and Allsides.com on 257 websites (covered in both databases) using Cohen's Kappa and get a reliability score of .9161, which indicates the substantial agreements between the two databases.

# Supplementary Method 7: Limiting the Sample to Corrective Individual Tags

We find that PolitiFact links are mostly used to correct the original posters' arguments. To strictly identify PolitiFact links that correct the original posters, we submitted original posts, replies containing PolitiFact links, and the cited PolitiFact fact-checking articles to ChatGPT (gpt-4o-2024-05-13). We prompted the model to annotate whether the PolitiFact link was used to correct the original poster, rather than support them. Specifically, we utilized the following prompt template.

> Here is a tweet and its reply. The reply includes a link to a fact-checking article from PolitiFact. Does the reply use the article to correct the original tweet (rather than support it)?
> [OP's tweet] {original_text}
> [Reply] {reply_text}
> [Fact-checking article cited in the reply] {article}
> Determine whether the reply or the fact-checking article corrects the original tweet.
> Simply respond with 1 if it does, and 0 if it does not or if it is not clear.

Consequently, we identified 5,592 politifact links out of 6,760 links (82.72%) that we can confidently say that they are corrective. To assess the accuracy of this annotation method, the researcher performed the same task on a randomly selected sample of 50 PolitiFact links. The results revealed high accuracy between the human annotator and ChatGPT, with an F1 score of 85.2%, precision of 92.0%, and recall of 79.3%. We find that the precision is higher than the recall, which means that ChatGPT uses stricter and more conservative criteria to determine whether the PolitiFact link corrects the original poster.

For instance, ChatGPT responds that the following reply and PolitiFact link corrects the original poster (i.e., 1).

> Here is a tweet and its reply. The reply includes a link to a fact-checking article from PolitiFact. Does the reply use the article to correct the original tweet (rather than support it)?
>
> [OP's tweet] We've seen almost double the number of children pass away from the vaccine compared to those lost to COVID. [Link]
> [Reply] @user1 @user2 That's not true. https://t.co/47fIls9Xfc
> [Fact-checking article cited in the reply]
> Claim Date: 2021-11-05
> Speaker: Viral image (Posters on Facebook and other social media; Party: None)

Claim: "Children are 50 times more likely to be killed by the Covid vaccines than by the virus itself."
Rating: Pants on Fire

Determine whether the reply or the fact-checking article corrects the original tweet. Simply respond with 1 if it does, and 0 if it does not or if it is not clear.

On the other hand, ChatGPT responds that the following PolitiFact link does not correct the original poster (i.e., 0). Indeed, the PolitiFact article is not utilized to correct the original poster but rather as a news article to present issues relevant to the topic.

Here is a tweet and its reply. The reply includes a link to a fact-checking article from PolitiFact. Does the reply use the article to correct the original tweet (rather than support it)?

[OP's tweet] Why schools in India are failing children on climate change

[Link]
[Reply] @user1 In the USA.. CLIMATE... over and over.. [Link]
[Fact-checking article cited in the reply]
Claim Date: 2017-01-17
Speaker: Chad Mayes (Party: Republican)
Claim: California has "the highest poverty rate in the nation" when considering the U.S. Census Bureau's Supplemental Poverty Measure.
Rating: True

Determine whether the reply or the fact-checking article corrects the original tweet. Simply respond with 1 if it does, and 0 if it does not or if it is not clear.

Subsequently, we limited the sample to the aforementioned 5,592 links identified by ChatGPT from the individual tagging data, which do not meaningfully alter the results (see Supplementary Tables 14 and 15)

49

# Supplementary Tables and Figures

**Supplementary Table 1. Illustrative Examples of Content Diversity Scores for a Sample User.**

| | Least Diverse Contents | | Most Diverse Contents | |
|---|---|---|---|---|
| Rank | Tweet | Content Diversity | Tweet | Content Diversity |
| 1 | Breaking News: Recently released Australian Government reports indicate a significant increase in excess deaths, up to 5162%, compared to the year 2020, potentially linked to COVID vaccination. @user @user, can you provide more insights? [Link to article removed] | .175936 | A call to release Pelosi's tax returns, citing precedent set by Trump's tax return release. #Transparency [Link to article removed] | .505400 |
| 2 | A government study reveals that 92% of COVID deaths were among those who were 'fully vaccinated.' Surprising information not widely covered in the media. [Link to article removed] | .196608 | Miley Cyrus donates to an LGBTQ organization following the banning of her song 'Rainbowland' in an elementary school. [Link to article removed] | .473959 |
| 3 | Japanese experts express confusion over high 'COVID deaths' despite a high vaccination rate, highlighting the need for more information. [Link to article removed] | .204821 | Zambia receives a 'debt-for-nature' proposal from WWF for $13 billion restructuring, raising questions about the motivations behind such deals. [Link to article removed] | .461800 |
| 4 | A German doctor receives a two-year jail sentence for illegally issuing thousands of mask exemptions, sparking a debate on the fairness of such punishment. [Link to article removed] | .205890 | An appreciation of superstars who use their influence for positive impact on humanity. [Link to TikTok video removed] | .456110 |
| 5 | Reports suggest an increase in excess mortality in Australia, but authorities remain quiet on the matter. #HealthCrisis [Link to article removed] | .209318 | Labor approves $9.5 million for 'facts of the voice' without acknowledging it as funding for a de-facto 'yes' campaign. @user, it's crucial for MPs to engage with constituents on this issue. [Link to article removed] | .454970 |

*Notes*: The table illustrates the examples of how content diversity scores are measured for a particular user. The user had primarily demonstrated interests centered around COVID-19 and associated misinformation. However, as the user explored a more diverse range of topical interests—including tax, LGBTQ+, international issues, and labor—the content diversity score increased. Following Twitter's privacy policy, we provide the rephrased tweets using ChatGPT (GPT-3.5) to ensure that user identities remain confidential (Prompt: Rewrite these tweets in a way that preserves the original meaning while respecting user privacy.).

**Supplementary Table 2. Diversity of Topics Corrected by Individual and Collective Tagging**

|  | Shannon H | Exp(H) | 1/Simpson D |
|---|---|---|---|
| Individual Tagging | 2.336 | 10.344 | 6.697 |
| Collective Tagging | 2.724 | 15.247 | 12.579 |

*Notes:* Results show that collective tagging corrects more diverse topics across various metrics compared to individual tagging. To examine whether each tag type selectively and repeatedly corrects certain political topics (low diversity) or corrects diverse topics (high diversity), we use Shannon H (Shannon Diversity Index), Exp(H) (Exponent of Shannon Index), and 1/Simpson D (Inverse Simpson Diversity Index). Shannon H is measured by $-\sum(p_i log(p_i))$, where $p_i$ is the proportion of tags belonging to a particular topic, with higher values indicating higher diversity. Simpson D is measured by $\sum(p_i^2)$. Higher values of 1/Simpson D indicate higher diversity.

**Supplementary Table 3. Political Stances of Taggers and Voters**

| | Individual tagging | |
|---|---|---|
| | **Tags correcting left-leaning users (N=2,470)** | **Tags correcting right-leaning users (N=2,815)** |
| **Taggers' political stance** | Left: 81.70%<br>Right: 18.30% | Left: 92.93%<br>Right: 7.07% |
| | **Collective tagging** | |
| | **Tags correcting left-leaning users (N=462)** | **Tags correcting right-leaning users (N=365)** |
| **Voters' political stance** | Left: 53.52%<br>Right: 46.48% | Left: 70.37%<br>Right: 29.63% |

*Notes*: Taggers are users who write individual tags (e.g., fact-checking replies), while voters are users who vote on the exposure of collective tags. For the methods used to measure the political stances of taggers and voters, refer to Supplementary Method 4: Political Stances of Taggers and Voters.

.

**Supplementary Table 4. Pairwise Correlation Coefficients Among Variables**

| Variables | Mean | SD | Pairwise Correlation Coefficients | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | | 1. Political Diversity (%) | 2. Content Diversity | 3. # of tweets per day | 4. % of retweets per day |
| 1. Political Diversity (%) | 16.570 | 37.181 | | | | |
| 2. Content Diversity | .357 | .109 | .020*** | | | |
| 3. # of tweets per day | 18.116 | 30.196 | -.107*** | -.052*** | | |
| 4. % of retweets per day | 32.650 | .407 | -.046*** | .012*** | -.200*** | |
| 5. Low-Credibility Sources (%) | 27.15 | 44.472 | -.227*** | -.030*** | .172*** | -.041*** |

*Notes: ***p<.001. SD: Standard Deviation*

.

**Supplementary Table 5. Control Analyses for the Proportion of Retweets**

| Outcome | Political diversity (%) | | Content diversity (z) | |
|---|---|---|---|---|
| Robustness checks | Initial results (N = 424,969) | Control the proportion of retweets (N = 424,969) | Initial results (N = 424,969) | Control the proportion of retweets (N = 424,969) |
| **Individual** | | | | |
| Slope before | .237*** (.057) | .220*** (.057) | .007*** (.002) | .007*** (.002) |
| Immediate change | -1.009*** (.223) | -.995*** (.223) | -.030*** (.006) | -.030*** (.006) |
| Slope after | .087 (.054) | .096+ (.054) | -.003+ (.001) | -.003+ (.001) |
| Slope change | -.150+ (.080) | -.124 (.080) | -.010*** (.002) | -.010*** (.002) |
| **Collective** | | | | |
| Slope before | .309* (.137) | .319* (.137) | .003 (.004) | .003 (.004) |
| Immediate change | .270 (.558) | .283 (.558) | .040** (.015) | .040** (.015) |
| Slope after | -.049 (.145) | -.028 (.145) | -.011** (.004) | -.011** (.004) |
| Slope change | -.358+ (.199) | -.347+ (.199) | -.014** (.005) | -.014** (.005) |
| **Individual vs. collective** | | | | |
| Slope before | .072 (.148) | .100 (.148) | -.005 (.004) | -.005 (.004) |
| Immediate change | 1.279* (.601) | 1.278* (.601) | .070*** (.016) | .070*** (.016) |
| Slope after | -.136 (.155) | -.124 (.155) | -.008* (.004) | -.008* (.004) |
| Slope change | -.208 (.215) | -.224 (.215) | -.004 (.006) | -.004 (.006) |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests.*

**Supplementary Table 6. Estimates of Interrupted Time Series (ITS) Models.**

| Type of misinformation tagging | Individual | Collective | Difference (Collective - Individual) |
|---|---|---|---|
| Slope before posting the tweet | $\beta_1$ | $\beta_1 + \beta_4$ | $\beta_4$ |
| Immediate intercept change after misinformation tagging | $\beta_2$ | $\beta_2 + \beta_5$ | $\beta_5$ |
| Slope after misinformation tagging | $\beta_1 + \beta_3$ | $\beta_1 + \beta_3 + \beta_4 + \beta_6$ | $\beta_4 + \beta_6$ |
| Slope change (after - before) | $\beta_3$ | $\beta_3 + \beta_6$ | $\beta_6$ |

**Supplementary Table 7. Estimates of Delayed Feedback (DF) Models.**

| Type of misinformation tagging | Individual | Collective | Difference (Collective - Individual) |
|---|---|---|---|
| Difference in Pre-Post Change (Treatment - Control) | $\beta_1$ | $\beta_1 + \beta_2$ | $\beta_2$ |

**Supplementary Table 8. Robustness Checks for Interrupted Time Series (ITS) Models**

| Outcome | Political diversity (%) | | | | |
|---|---|---|---|---|---|
| Robustness checks | Initial results (*N*=424,969) | Removing bots (N=401,188) | Controlling credibility (N=424,969) | Without negative sentiments (N=268,626) | Removing non-responders (*N*=167,766) |
| **Individual** | | | | | |
| Slope before | .237*** (.057) | .221*** (.058) | .229*** (.055) | .224** (.074) | .514*** (.114) |
| Immediate change | -1.009*** (.223) | -.882*** (.226) | -.972*** (.216) | -1.470*** (.286) | -1.481** (.436) |
| Slope after | .087 (.054) | .073 (.055) | .035 (.053) | .237** (.070) | .145 (.106) |
| Slope change | -.150+ (.080) | -.148+ (.081) | -.194* (.077) | .013 (.102) | -.368* (.157) |
| **Collective** | | | | | |
| Slope before | .309* (.137) | .339* (.150) | .225+ (.132) | .336+ (.178) | .311* (.147) |
| Immediate change | .270 (.558) | .223 (.620) | .389 (.539) | .375 (.722) | .270 (.599) |
| Slope after | -.049 (.145) | -.117 (.162) | -.023 (.140) | -.103 (.189) | -.049 (.156) |
| Slope change | -.358+ (.199) | -.456* (.221) | -.248 (.193) | -.439+ (.259) | -.361+ (.214) |
| **Individual vs. collective** | | | | | |
| Slope before | .072 (.148) | .118 (.161) | -.004 (.143) | .112 (.192) | -.203 (.186) |
| Immediate change | 1.279* (.601) | 1.105+ (.660) | 1.361* (.581) | 1.845* (.777) | 1.752* (.741) |
| Slope after | -.136 (.155) | -.190 (.171) | -.058 (.150) | -.340+ (.202) | -.195 (.189) |
| Slope change | -.208 (.215) | -.308 (.235) | -.054 (.207) | -.452 (.279) | .008 (.266) |
| Outcome | Content diversity (z) | | | | |
| Robustness checks | Initial results (*N*=424,969) | Removing bots (N=401,188) | Controlling credibility (N=424,969) | Without negative sentiments (N=268,626) | Removing non-responders (*N*=167,766) |
| **Individual** | | | | | |
| Slope before | .007*** (.002) | .006*** (.002) | .007*** (.002) | .007** (.002) | .004 (.003) |
| Immediate change | -.030*** (.006) | -.026*** (.006) | -.030*** (.006) | -.027** (.008) | -.026* (.011) |
| Slope after | -.003+ (.001) | -.003* (.001) | -.003* (.001) | -.005* (.002) | -.005+ (.003) |
| Slope change | -.010*** (.002) | -.009*** (.002) | -.010*** (.002) | -.011*** (.003) | -.009* (.004) |
| **Collective** | | | | | |
| Slope before | .003 (.004) | .001 (.004) | .003 (.004) | .012* (.005) | .003 (.004) |
| Immediate change | .040** (.015) | .052** (.016) | .041** (.015) | .009 (.020) | .040* (.016) |
| Slope after | -.011** (.004) | -.013** (.004) | -.011** (.004) | -.006 (.005) | -.011** (.004) |
| Slope change | -.014** (.005) | -.015* (.006) | -.013* (.005) | -.018* (.007) | -.014* (.006) |
| **Individual vs. collective** | | | | | |
| Slope before | -.005 (.004) | -.005 (.004) | -.005 (.004) | .006 (.005) | -.001 (.005) |
| Immediate change | .070*** (.016) | .077*** (.017) | .071*** (.016) | .035 (.022) | .067** (.019) |
| Slope after | -.008* (.004) | -.010* (.005) | -.008* (.004) | -.001 (.006) | -.006 (.005) |
| Slope change | -.004 (.006) | -.005 (.006) | -.003 (.006) | -.007 (.008) | -.005 (.007) |

*Notes: Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests.*

**Supplementary Table 9. Robustness Checks for Delayed Feedback (DF) Models**

| Outcome | Political diversity (%) | | | | |
|---|---|---|---|---|---|
| Robustness checks | Initial results (*N*=8,901) | Removing bots (*N*=8,442) | Controlling credibility (*N*=8,901) | Without negative sentiments (*N*=5,615) | Removing non-responders (*N*=3,275) |
| Individual | -5.886** (1.911) | -5.348** (1.942) | -4.958** (1.854) | -8.075** (2.388) | -1.747 (6.571) |
| Collective | 1.219 (3.059) | -.899 (3.364) | 2.246 (2.966) | 3.178 (3.874) | 2.338 (3.568) |
| Individual vs. Collective | 7.105* (3.589) | 4.449 (3.865) | 7.204* (3.480) | 11.253* (4.536) | 13.086+ (7.465) |
| Outcome | Content diversity (z) | | | | |
| Robustness checks | Initial results (*N*=8,901) | Removing bots (*N*=8,442) | Controlling credibility (*N*=8,901) | Without negative sentiments (*N*=5,615) | Removing non-responders (*N*=3,275) |
| Individual | .018 (.041) | .026 (.042) | .022 (.041) | .043 (.057) | .014 (.142) |
| Collective | .274*** (.066) | .292*** (.073) | .279*** (.066) | .294** (.092) | .298*** (.077) |
| Individual vs. Collective | .256** (.077) | .266** (.084) | .256** (.077) | .251* (.108) | .283+ (.161) |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Observations refer to the number of tweets that received individual or collective tags. Each cell presents the difference in pre-post change (treatment group - control group) in each outcome. Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests.*

**Supplementary Table 10. Interrupted Times Series (ITS) Models after Removing Tweets Mentioning the Experience of Receiving Tags or Interacting with Taggers**

| Outcome | Political diversity (%) | | Content diversity (z) | |
|---|---|---|---|---|
| Robustness checks | Initial results (N = 424,969) | Removing tweets mentioning tags or replying to taggers (N = 422,991) | Initial results (N = 424,969) | Removing tweets mentioning tags or replying to taggers (N = 422,991) |
| Individual | | | | |
| Slope before | .237*** (.057) | .228*** (.057) | .007*** (.002) | .007*** (.002) |
| Immediate change | -1.009*** (.223) | -1.048*** (.225) | -.030*** (.006) | -.030*** (.006) |
| Slope after | .087 (.054) | .107* (.055) | -.003+ (.001) | -.003+ (.001) |
| Slope change | -.150+ (.080) | -.121 (.080) | -.010*** (.002) | -.010*** (.002) |
| Collective | | | | |
| Slope before | .309* (.137) | .307* (.137) | .003 (.004) | .003 (.004) |
| Immediate change | .270 (.558) | .262 (.557) | .040** (.015) | .041** (.015) |
| Slope after | -.049 (.145) | -.046 (.145) | -.011** (.004) | -.011** (.004) |
| Slope change | -.358+ (.199) | -.352+ (.199) | -.014** (.005) | -.014** (.005) |
| Individual vs. collective | | | | |
| Slope before | .072 (.148) | .078 (.148) | -.005 (.004) | -.004 (.004) |
| Immediate change | 1.279* (.601) | 1.310* (.601) | .070*** (.016) | .070*** (.016) |
| Slope after | -.136 (.155) | -.153 (.155) | -.008* (.004) | -.008* (.004) |
| Slope change | -.208 (.215) | -.232 (.215) | -.004 (.006) | -.004 (.006) |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P-values) are derived from two-sided tests.*

**Supplementary Table 11. Delayed Feedback (DF) Models after Removing Tweets Mentioning the Experience of Receiving Tags or Interacting with Taggers**

| Outcome | Political diversity (%) | | Content diversity (z) | |
|---|---|---|---|---|
| Robustness checks | Initial results (N = 8,901) | Removing tweets mentioning tags or replying to taggers (N = 8,787) | Initial results (N = 8,901) | Removing tweets mentioning tags or replying to taggers (N = 8,787) |
| Individual | -5.886** (1.911) | -5.329** (1.944) | .018 (.041) | .016 (.042) |
| Collective | 1.219 (3.059) | 1.070 (3.053) | .274*** (.066) | .276*** (.066) |
| Individual vs. Collective | 7.105* (3.589) | 6.400+ (3.601) | .256** (.077) | .260** (.077) |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Each cell presents the difference in pre-post change (treatment group - control group) in each outcome. Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests.*

**Supplementary Table 12. Pairwise Differences in Content Diversity Before and After Collective Tagging**

| A pair of weeks being compared | Difference in content diversity (Coefficient (Standard Error)) | P-value |
|---|---|---|
| (Week 5) - (Week -5) | -.001 (.016) | .929 |
| (Week 5) - (Week -4) | -.004 (.014) | .776 |
| (Week 5) - (Week -3) | -.007 (.013) | .615 |
| (Week 5) - (Week -2) | -.009 (.013) | .482 |
| (Week 5) - (Week -1) | -.012 (.014) | .398 |
| (Week 5) - (Week 0) | -.015 (.016) | .357 |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Statistical significance levels (P values) are derived from two-sided tests.*

**Supplementary Table 13. Interrupted Time Series (ITS) Model Results for Political and Content Proximity to Taggers**

| Outcome | Political Proximity (%) | Content Proximity |
|---|---|---|
| Slope before | .146* (.060) | .003* (.001) |
| Immediate change | -.536* (.234) | .007 (.005) |
| Slope after | -.057 (.057) | -.002 (.001) |
| Slope change (After - Before) | -.203* (.083) | -.005* (.002) |
| Observations | 317,442 | |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. We multiply political proximity by 100 to interpret the estimates as absolute percentage point changes. We normalize content proximity to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests.*

**Supplementary Table 14. Interrupted Times Series (ITS) Models after Strictly Limiting the Sample to Corrective Individual Tags**

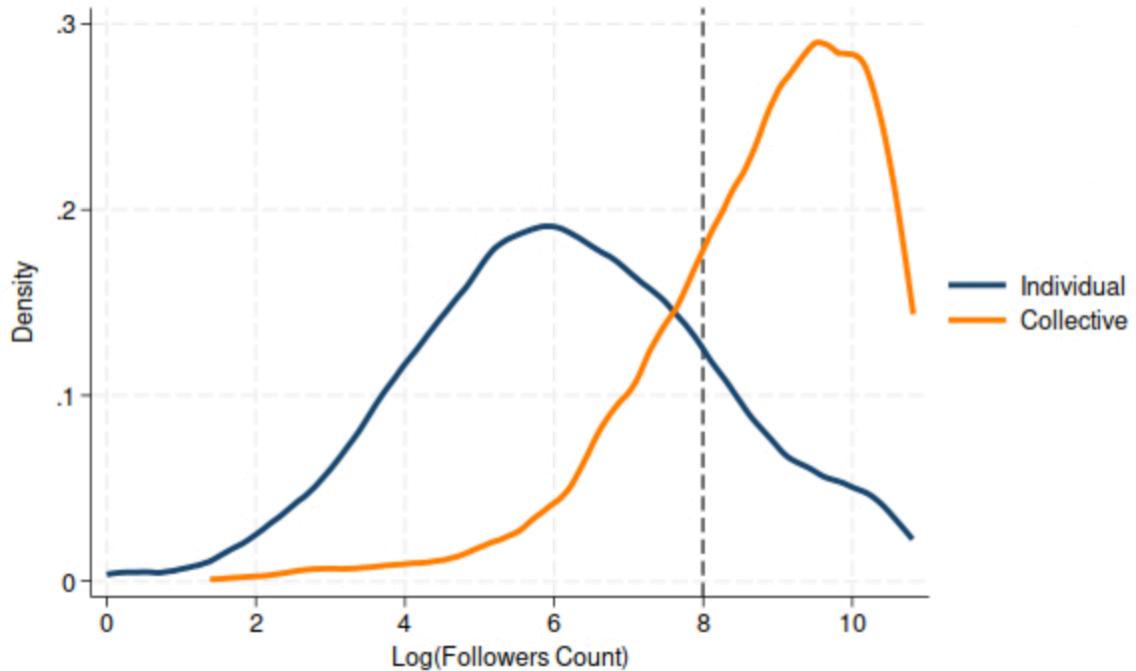| Outcome | Political diversity (%) | | Content diversity (z) | |
|---|---|---|---|---|
| Robustness checks | Initial results (N = 424,969) | Limiting the sample (N = 369,070) | Initial results (N = 424,969) | Limiting the sample (N = 369,070) |
| **Individual** | | | | |
| Slope before | .237*** (.057) | .284*** (.062) | .007*** (.002) | .009*** (.002) |
| Immediate change | -1.009*** (.223) | -1.086*** (.241) | -.030*** (.006) | -.039*** (.006) |
| Slope after | .087 (.054) | .070 (.059) | -.003+ (.001) | -.003* (.002) |
| Slope change | -.150+ (.080) | -.214* (.086) | -.010*** (.002) | -.012*** (.002) |
| **Collective** | | | | |
| Slope before | .309* (.137) | .309* (.136) | .003 (.004) | .003 (.004) |
| Immediate change | .270 (.558) | .270 (.554) | .040** (.015) | .040** (.015) |
| Slope after | -.049 (.145) | -.049 (.144) | -.011** (.004) | -.011** (.004) |
| Slope change | -.358+ (.199) | -.358+ (.198) | -.014** (.005) | -.014** (.005) |
| **Individual vs. collective** | | | | |
| Slope before | .072 (.148) | .025 (.149) | -.005 (.004) | -.006 (.004) |
| Immediate change | 1.279* (.601) | 1.356* (.604) | .070*** (.016) | .079*** (.016) |
| Slope after | -.136 (.155) | -.120 (.156) | -.008* (.004) | -.008+ (.004) |
| Slope change | -.208 (.215) | -.144 (.216) | -.004 (.006) | -.001 (.006) |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P-values) are derived from two-sided tests.*

**Supplementary Table 15. Delayed Feedback (DF) Models after Strictly Limiting the Sample to Corrective Individual Tags**
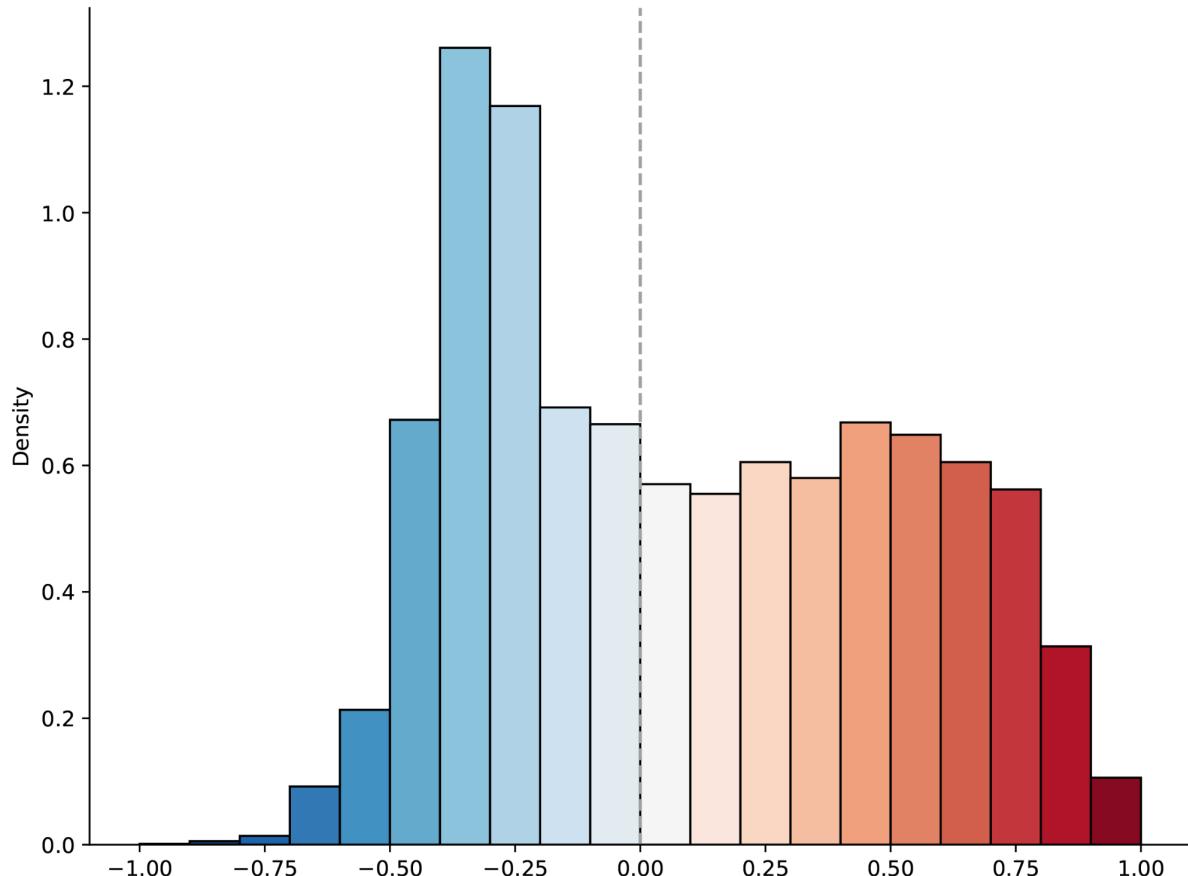
| Outcome | Political diversity (%) | | Content diversity (z) | |
|---|---|---|---|---|
| Robustness checks | Initial results (N = 8,901) | Limiting the sample (N = 8,145) | Initial results (N = 8,901) | Limiting the sample (N = 8,145) |
| Individual | -5.886** (1.911) | -7.135*** (2.006) | .018 (.041) | .026 (0.044) |
| Collective | 1.219 (3.059) | 1.336 (2.994) | .274*** (.066) | .272*** (0.066) |
| Individual vs. Collective | 7.105* (3.589) | 8.471* (3.584) | .256** (.077) | .246** (0.079) |

*Notes: \*\*\*p<.001 \*\*p<.01 \*p<.05 +p<.1. Each cell presents the difference in pre-post change (treatment group - control group) in each outcome. Political diversity has been multiplied by 100 so that the estimates are interpretable as absolute percentage point changes. Content diversity has been normalized to z-scores (the number of standard deviations from the mean). Standard errors are in parentheses. All regressions control for user fixed effects and the number of tweets per day. Statistical significance levels (P values) are derived from two-sided tests.*
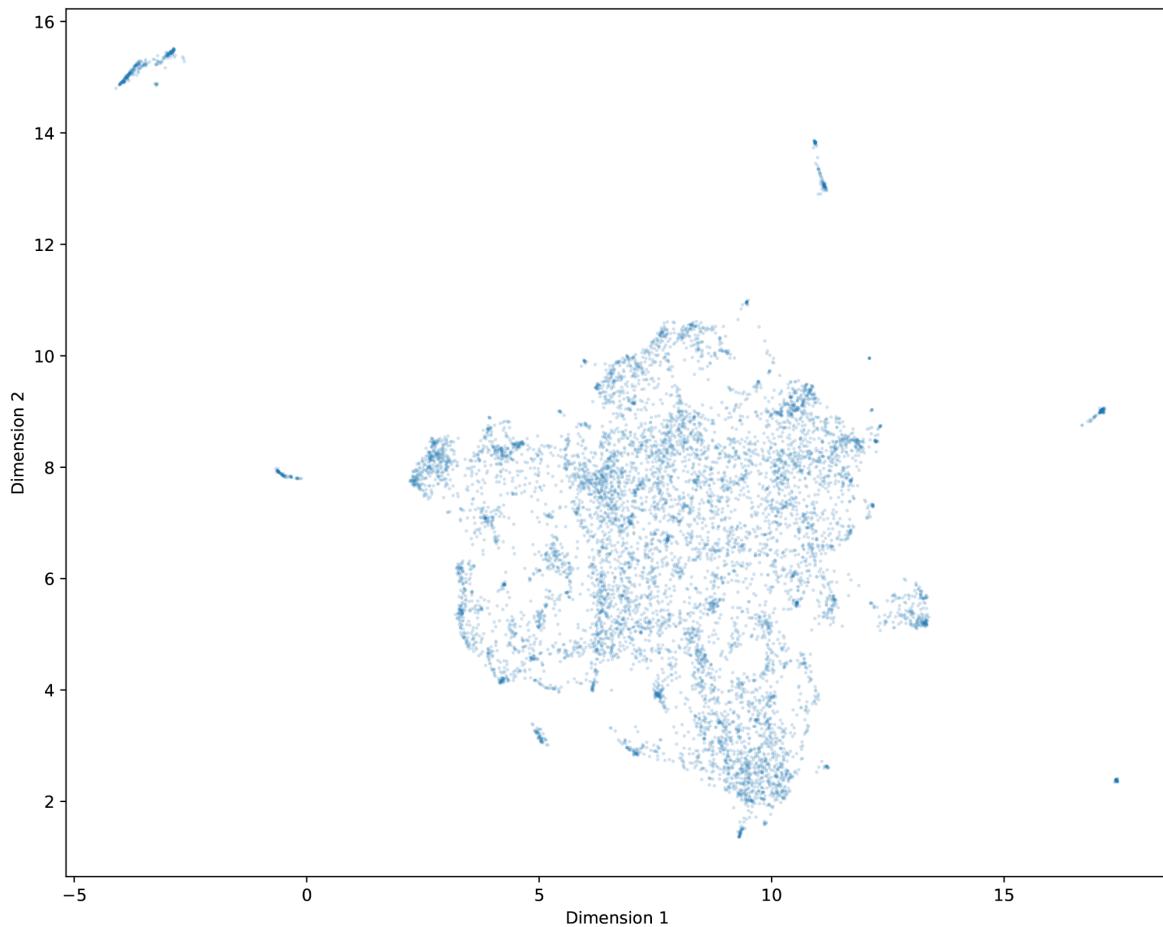
**Supplementary Fig. 1. Distribution of the Log-Transformed Number of Followers of Tagged Users.** The blue line indicates the distribution of the log-transformed number of followers among posters fact-checked by individual tags. The orange line indicates the distribution among posters fact-checked by collective tags. The gray dashed line indicates the average number of followers among posters fact-checked by individual tags (i.e., 2,967).
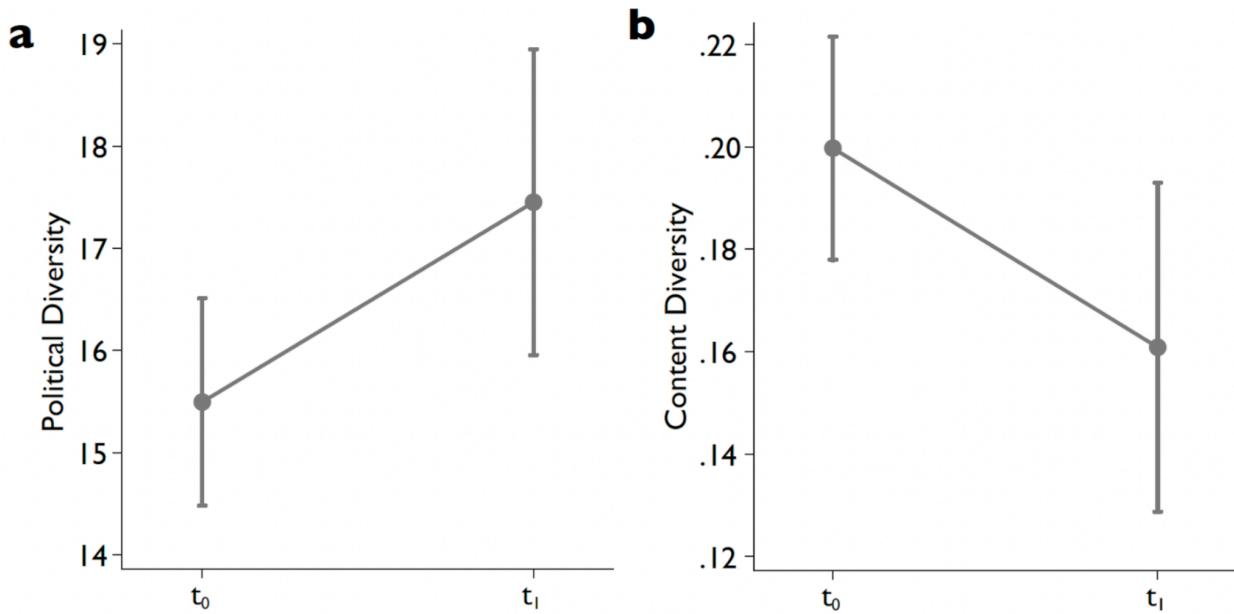
**Supplementary Fig. 2. Distribution of Tagged Posters' Political Stances.** The political stance scores range from -1, indicating a left-leaning stance, to 1, indicating a right-leaning stance. The average political stance is 0.166 (SD = 0.372), indicating a tendency for a larger proportion of users to lean to the right.

**Supplementary Fig. 3. Tweet Embedding Space.** This figure presents a two-dimensional UMAP (Uniform Manifold Approximation and Projection) visualization, illustrating the embeddings of a randomly selected subset of 10,000 tweets posted before and after misinformation tagging. We show the largest cluster identified by HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise).

**Supplementary Fig. 4. Average baseline changes of political and content diversity obtained from DF analysis.**

**Supplementary Fig. 5. Political and Content Proximity to Taggers Changes with the Intervention of Individual Misinformation Tagging.** Results from Interrupted Time Series (ITS) analysis. The *x*-axis denotes the timeline of tweets posted before and after tagging, with negative values indicating the number of weeks before posting tagged tweets and positive values indicating the number of weeks after. The *y*-axis denotes political and content proximity to taggers, with dots capturing the average diversity score of the corresponding week, and error bars indicating 95% confidence intervals.