## Misinformation Warning Labels: Twitter's Soft Moderation Effects on COVID-19 Vaccine Belief Echoes

Filipo Sharevski DePaul University Chicago, IL fsharevs@cdm.depaul.edu

Peter Jachim DePaul University Chicago, IL pjachim@depaul.edu Raniem Alsaadi DePaul University Chicago, IL ralsaad2@depaul.edu

Emma Pieroni DePaul University Chicago, IL epieroni@depaul.edu

#### **ABSTRACT**

Twitter, prompted by the rapid spread of alternative narratives, started actively warning users about the spread of COVID-19 misinformation. This form of soft moderation comes in two forms: as a warning cover before the Tweet is displayed to the user and as a warning tag below the Tweet. This study investigates how each of the soft moderation forms affects the perceived accuracy of COVID-19 vaccine misinformation on Twitter. The results suggest that the warning covers work, but not the tags, in reducing the perception of accuracy of COVID-19 vaccine misinformation on Twitter. "Belief echoes" do exist among Twitter users, unfettered by any warning labels, in relationship to the perceived safety and efficacy of the COVID-19 vaccine as well as the vaccination hesitancy for themselves and their children. The implications of these results are discussed in the context of usable security affordances for combating misinformation on social media.

#### **CCS CONCEPTS**

• Security and privacy  $\rightarrow$  Social aspects of security and privacy; Usability in security and privacy;

#### **KEYWORDS**

Soft moderation, Twitter, COVID-19, misinformation, warning labels, belief echoes

#### **ACM Reference Format:**

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

#### 1 INTRODUCTION

In 2016, when "fake news" gained enormous popularity, Facebook started adding tags that say "disputed" on stories that were debunked by fact-checkers [30]. About a year later, Facebook started adding fact-checks under potentially misleading stories [37]. The goal of these initiatives was presumably to minimize the probability that readers will believe the fake information. Twitter did not begin similar initiatives until 2020, when, in late March, the platform began issuing labels on Tweets deemed as spreading misinformation related to the COVID-19 pandemic [34]. According to Twitter, they are relying on their team and internal systems to monitor COVID-19 content for false or misleading information that is not corroborated by public health authorities or subject matter experts. The supposed aim of these labels is to reduce misleading or harmful information that could "incite people to action and cause widespread panic, social unrest or large-scale disorder" [34]. Originally, only Tweets that pertained to COVID-19 were flagged; however, following the November 2020 presidential election, Twitter broadened the types of misleading, false or disputed information to which it appended warning labels about the outcome of the election, claims of election fraud, or the safety of voting by mail [38].

However, there is no evidence that these labels are effective, and in fact, an early investigation suggests that they can have an effect of "backfiring" i.e. convince people to believe the misinformation even more than if the label were not there [7]. A study exploring this so-called "soft moderation" implemented by Twitter found that the platforms' content with warning labels generated more action than content without said labels [45]. The soft moderated Tweets were more likely to be distributed than a "valid" Tweet through discourse, not always because the soft moderated Tweets contained misinformation but because they contained a response to mock or disclaim an original author or a valid Tweet. Interestingly, a mere 1% of the Tweets gathered for the study were labeled with a COVID-19 warning and a number of these few Tweets were found to be mislabeled simply because they contained the words "oxygen" and "frequency". Another study, in this context, found that some users did not trust the soft moderation intervention and felt that Twitter itself was biased and mislabeling content [20]. Even without the warning labels, a varying degree of users' perceptions and beliefs regarding vaccines in general and about COVID-19 vaccines in particular plays a role in individuals' reaction to (mis)information Tweets. A study on vaccine misinformation spreading on social media platforms showed the effectiveness of "influential users" within an "echo chamber" of vaccine-fearing followers to be very high [13].

In this study, we set to examine the association between the beliefs regarding COVID-19 vaccines and the misinformation warning labels. To assess individuals' beliefs about the COVID-19 vaccines, we used the questionnaire about COVID-19 vaccines from [2] adopted to the current state of the vaccine development and distribution (we conducted the study in January-February 2021 period, after the questionnaire was published and the vaccine was approved). We sought to examine the effect of different types of warning labels, represented as warning covers over the content or warning tags under the content [34], as experiments have shown that the design of the warning label itself (how the warning is presented to users) can affect the effectiveness of the warning [25]. To examine the perceived accuracy of (mis)information Tweets about COVID-19 vaccine we utilized content from Twitter and a software that altered the warning tags to induce misperception. This software was initially developed to study polarized topics of discourse on social media on vaccines' safety and users' perception of accuracy contrary to their views and beliefs [36].

We found that the warning tags are ineffective in reducing the "belief echoes" on Twitter regarding the COVID-19 vaccine. The intended reduction effect, our results suggest, is achieved only with the misinformation warning covers preceding a misleading information (in our case, we use a Tweet referring to unverified adverse effects of the COVID-19 vaccination). We found that the less the users believed the COVID-19 vaccines are safe and efficacious, the more they perceived the misleading Tweets as accurate, even in presence of a warning tag below the Tweet's content. We also observed an echoing scepticism where the more the participants believed the COVID-19 vaccines are safe and efficacious, the less they perceived a Tweet being accurate, even in the case when they were presented with a verified COVID-19 vaccine information (A Tweet following Centers For Disease Control (CDC) guidelines in case there are any adverse effects after receiving the first COVID-19 vaccine dose). A similar echoing of beliefs and sceptic outlook of misleading and verified content, respectively, was found about the beliefs that herd immunity is a better option of immunization than mass COVID-19 vaccination. When it came to vaccine hesitancy, the warning labels did little to sway the participants on the benefits of the vaccination - the ones that were hesitant to receive the COVID-19 vaccine were convinced that it causes adverse effects leading to death, even if warned against such a claim. The anti-COVID-19 vaccine sentiment persisted, only in the case of the warning tags, when we asked whether children should get the vaccine too.

We additionally explored the political affiliation and the effect of soft moderation, following the previously observed divisions on misinformation content along the main party lines [13, 45]. We found that the Republicans and independent participants perceived the misleading Tweet as "somewhat accurate" while the Democrat participants perceived it as "not very accurate," regardless of the presence of a warning tag. We also found that almost one in four Republicans and and one in six independent participants didn't expect to have a efficacious COVID-19 vaccine, while that proportion for the Democrats was one in forty. Half the Republicans and a

third of the independent participants were hesitant to receive the COVID-19 vaccine, while only  $1/20^{th}$  of the Democrats won't proceed with personal immunization. Roughly 40% of the Republicans and independent participants were hesitant to vaccinate children for COVID-19, to which only 8.3% of the Democrats agree with.

We consider the misinformation warning labels as a form of usable security warnings akin to warnings about potentially harmful websites or favicons indicating unverified certificates [18]. Users, studies have shown, are reluctant to heed these warnings due to a lack of attention or motivation, incomprehension, or habituation [11, 31, 43]. The warning labels are written in plain language to draw attention to the user about the validity of the content. While it is early to assess the habituation effect of the warnings, the existence of the belief echoes posits an analysis of the unique blend of motivation and polarized habituation on Twitter [27]. We discuss the implications of our results in context of future designs of warning labels, as a form of a usable security interventions, aimed to curtail misinformation on social media.

#### 2 BACKGROUND

#### 2.1 Soft Moderation

The misinformation warning labels (or tags) provide a compromise between content removal and the commitment of the social media platforms to allow for free and constructive discourse. However, whether such warnings and corrections/fact checks are effective in achieving their aim remains unclear. When measured based on Tweet engagement (e.g., likes, retweets or quote retweets), it appears that such warnings may be somewhat effective: Twitter reported a 29 % decrease in quote Tweets that were labelled as misleading or disputed [38]. However, authors in [45] found that Tweets with warning labels received more engagement (likes, retweets, replies, and quote retweets) than other Tweets from the same users that did not have warning labels. Specifically, [45] found that between 2/3 to 4/5<sup>th</sup> of users receive more engagement on Tweets that contain content warning tags than Tweets that do not. Yet engagement is just one way to assess whether such warnings or corrections have an effect. This is especially the case since users' engagement varied such that some users reinforced false claims, mocked the false claims, or debunked the claims [45]. Thus, evaluating Tweet engagement alone is insufficient to evaluate whether content warning tags work to decrease the spread (or beliefs) of false information.

Examining whether correcting or fact checking false information affects readers' opinions and beliefs about the information, authors in [32] found that correcting mock news articles that included false claims for politicians often failed to reduce misperceptions among particular ideological groups. The corrections in studies often backfired, increasing misperceptions in the targeted group. The existence of so-called "backfire effect" suggests that providing corrections and fact checking information may have a countereffect in combating false information (although the effect in [32] is observed for mock news articles, not Facebook/Twitter posts). The backfiring effect was observed also in [45] for the Elections 2020, which found that 72% of the Tweets with warning labels were shared by Republicans while only 11% are shared by Democrats.

#### 2.2 Belief Echoes on Social Media

More recently, studies have begun examining whether adding labels to posts on social media (as opposed to in articles) can affect individuals' beliefs. For example, authors in [7] examined whether strategies that social media companies such as Twitter and Facebook use to oppose false stories or "fake news" would have the intended effect. This study also evaluated the efficacy of different types of labels: (i) a general warning, and (ii) two specific warnings pertaining to the article content. The authors found that a general warning had the intended effect of decreasing the perceived accuracy of the information but that adding "disputed" or "rated false" tags had a larger effect on minimizing perceived accuracy of the content, with the "rated false" tag most effective. Interestingly, and somewhat in contrast to the findings by Twitter regarding Tweet engagement [38], the authors in [7] found that the tags did not reduce participants' self-reported likelihood of sharing the headlines on social media. Authors in [6] evaluated whether social media corrections of presidential Tweets on support of executive policies affects individuals' attitudes. The idea was to test whether corrections are effective at rebutting false claims or whether they promote belief in the false claims among a particular demographic, a phenomenon dubbed as "belief echoes" [42]. The corrections had the intended effect on Democrats but the opposite effect on Republicans, showing evidence of the "belief echoes" on Twitter for the later category.

The "belief echoes" manifest on social media when the exposure to negative political information continues to shape attitudes even after the information has been effectively discredited. Belief echoes can result as a spontaneous affective response that is immediately integrated into a person's summary evaluation of social media content. The mere exposure to misinformation often generates a strong and automatic affective response, but the correction may not generate a response of an equal and opposite magnitude [19]. One reason for this is that warning labels as commonly phrased as negations or contain exclamation marks. To combat the affectively asymmetrical soft moderation, users need to engage in cognitively demanding and time consuming "strategic retrieval monitoring" [9] or recall of the warning label. That does not happen often, so the misinformation may continue to affect evaluations, thus creating automatic belief echoes. Even if a person recalls the correction, they may discard it because they are already negatively predisposed to it. For example, in the context of politics, if a person hears that a candidate was accused of fraud, they may reason that the accusation emerged because the candidate is generally untrustworthy or corrupt. If these secondary inferences linger after the initial information is discredited, they will continue to affect their evaluations [42].

### 2.3 COVID-19 Vaccine Echo Beliefs

Social media provides a vehicle for the spread of information regarding vaccines and vaccinations. Studies have found that most information on Twitter regarding vaccines is polarizing on the vaccine hesitancy and the beliefs about vaccines effects on child development [24]. The consumption of this information may affect individuals' perceptions, attitudes and beliefs about vaccinations [28]. For instance, vaccine-related Tweets by bots and trolls affect vaccine discourse on Twitter by promoting a relationship between vaccines and autism in children [3] or a relationship between

COVID-19 vaccines and significant adverse effects, including death, for adults [1]. Thus, misinformation regarding vaccines can have a significant effect on the acceptance of COVID-19 vaccines [8].

The current global pandemic provides ample opportunities for rampant misinformation regarding vaccines [10, 44]. This is particularly worrying because uptake of COVID-19 vaccines is critical for containing the spread of this disease and decreasing the morbidity and mortality imposed by the pandemic [26]. Ensuring that individuals perceive the COVID-19 vaccines as safe once they become available requires that they have the correct information regarding COVID-19 vaccines [26]. Currently, a significant minority of the worldwide population expresses skepticism about the safety, efficacy, and necessity of COVID-19 vaccines, which may make them more hesitant to take the COVID-19 vaccine. For instance, in Canada and the United States, 68.7% and 75.3%, respectively, reported being likely or very likely to accept the COVID-19 vaccine [26]. Given the spread of the COVID-19 pandemic and the spread of misinformation regarding COVID-19 vaccines on Twitter [44], it is imperative to explore the role of misinformation warning labels, as a form of usable security warnings, to curb misinformation pertaining to COVID-19 vaccines and vaccination more broadly.

#### 3 RESEARCH STUDY

#### 3.1 Belief Echoes: Preconditions

In this study, we set to examine the association between COVID-19 vaccine perceptions, beliefs, and hesitancy, the effect of the misinformation warning labels, and the perceived accuracy of (mis)information Tweets about COVID-19 vaccine content. First, we set to examine the *preconditions* for existence of belief echoes on Twitter regarding COVID-19 vaccines. In particular, we investigated whether exposure to (mis)information Tweets about the COVID-19 vaccine efficacy in the presence or absence of warning labels, both as tags and covers, affect individuals' perceptions of the Tweet's accuracy with the following set of hypotheses:

- H1: The presence of a warning tag under a Tweet containing
   misleading information about COVID-19 vaccines will not
   reduce the perceived accuracy of the Tweet's content relative
   to a no warning tag condition.
- H2: The presence of a warning cover before a Tweet containing *misleading* information about COVID-19 vaccines is shown to the user will not reduce the perceived accuracy of the Tweet's content relative to a no warning cover condition.
- H3: The presence of a warning tag under a Tweet containing verified information about COVID-19 vaccines will not reduce the perceived accuracy of the Tweet's content relative to a normal no warning tag condition.
- H4: The presence of a warning cover (malware inserted) before a Tweet containing *verified* information about COVID-19 vaccines is shown to the user will not reduce the perceived accuracy of the Tweet's content relative to a normal no warning cover condition.

To test the first hypothesis we utilized the Tweets containing *misleading information* shown in Figure 1a and Figure 1b. The Tweet in Figure 1a shows a warning tag underneath a Tweet, indicating

that the content is labeled as misinformation. The Tweet promulgates COVID-19 misinformation about a rare adverse effect that was linked to the SARS-CoV-2 virus, not the vaccine, at the time of writing [5]. To remove bias due to the "influencer" effect, the Tweet comes from a verified account named "TheVaccinator" (which we made up) and indicates a relatively high interaction engagement with 3k retweets, 13.5k quotations, and 12.8k likes, which is consistent with the expected engagement of Tweets containing COVID-19 vaccine information [45]. An alteration of the same Tweet is shown in Figure 1b without the accompanying warning tag. To test the second hypothesis we utilized the Tweets containing *misleading information* shown in Figure 1b and Figure 2 (which includes a warning cover instead of a warning tag).

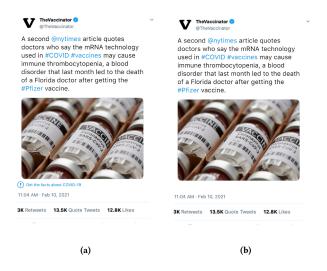


Figure 1: A Misleading Tweet: (a) With a Warning Tag; (b) Without a Warning Tag for Misleading Information.



Figure 2: A Warning Cover Preceding a Misleading Tweet

To test the third hypothesis we utilized the Tweets containing *verified information* shown in Figure 3a and Figure 3b. The Tweet

content indicates the verified information distributed by the CDC about proceeding with the second dose of the COVID-19 vaccine in case an individual has a serious reaction from the first dose, altered to include a warning tag in Figure 3b [15]. To control for bias, the Tweet comes from a verified account "TheVirusMonitor" instead of the CDC account and indicates a similar engagement as the misleading Tweet [45]. To test the fourth hypothesis we utilized the Tweets containing *verified information* shown in Figure 3a and Figure 4. We retained Figure 3a for the comparison of the conditions and altered the labeling in the Figure 3b to include a warning cover instead of a warning tag.

## 3.2 Belief Echoes: Safety and Herd Immunity

Assuming the preconditions of the belief echoes are met, we examined the relationship between COVID-19 vaccine beliefs on safety

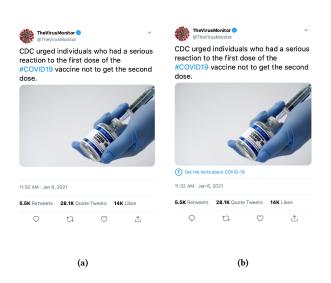


Figure 3: A Verified Tweet: (a) Without a Warning Tag; (b) With a Warning Tag for Misleading Information.



Figure 4: A Warning Cover Preceding a Verified Tweet

and herd immunity and the perceived accuracy of Tweets with COVID-19 vaccine information in presence/absence of warning labels. We used the same Tweets from Figures 1-4 together to test the following hypotheses:

- H5a: The belief that COVID-19 vaccines are not safe will not affect the perception of accuracy of a Tweet with misleading information about COVID-19 in any condition (with a warning tag/cover or without any warning)
- H5b: The belief that COVID-19 vaccines are not safe will not affect the perception of accuracy of a Tweet with verified information about COVID-19 in any condition (with a warning tag/cover or without any warning)
- H6a: The belief that there is no need for COVID-19 vaccine because herd immunity exists will not affect the perception of accuracy of a Tweet with misleading information about COVID-19 in any condition (with a warning tag/cover or without any warning)
- H6b: The belief that there is no need for COVID-19 vaccine because herd immunity exists will not affect the perception of accuracy of a Tweet with verified information about COVID-19 in any condition (with a warning tag/cover or without any warning)

## 3.3 Belief Echoes: Efficacy and Hesitancy

Next, we examined the relationship between COVID-19 vaccine efficacy/hesitancy and the perceived accuracy of Tweets with COVID-19 vaccine information in presence/absence of warning labels. We used the same Tweets as shown in Figures 1-4 together to test the following hypotheses:

- H7a: The perception of producing efficacious COVID-19 vaccine will not affect the perception of accuracy of a Tweet with misleading information about COVID-19 in any condition (with a warning tag/cover or without any warning)
- H7b: The perception of producing efficacious COVID-19 vaccine will not affect the perception of accuracy of a Tweet with verified information about COVID-19 in any condition (with a warning tag/cover or without any warning)
- H8a: The COVID-19 vaccine personal hesitancy will not affect the perception of accuracy of a Tweet with misleading information about COVID-19 in any condition (with a warning tag/cover or without any warning)
- H8b: The COVID-19 vaccine personal hesitancy will not affect the perception of accuracy of a Tweet with verified information about COVID-19 in any condition (with a warning tag/cover or without any warning)
- H9a: The COVID-19 vaccine hesitancy for children will not affect the perception of accuracy of a Tweet with misleading information about COVID-19 in any condition (with a warning tag/cover or without any warning)
- H9b: The COVID-19 vaccine hesitancy for children will not affect the perception of accuracy of a Tweet with verified information about COVID-19 in any condition (with a warning tag/cover or without any warning)

#### 3.4 Belief Echoes and Political Affiliation

To test the association between one's political affiliation and the perceived accuracy of the Tweets from Figures 1-4, following the evidence in [45] about the interplay between political affiliation and misinformation warning labels, we asked:

- RQ1: Is there a difference in the perceived accuracy of COVID-19 misleading/verified Tweets with warning labels (tags or covers) between Republican, Democrat, and Independent users?
- RQ2: Is there a difference between the beliefs and subjective attitudes of the Twitter users about the COVID-19 vaccine based on their political affiliation?

## 3.5 Sampling and Instrumentation

We first got approval from our Institutional Review Board (IRB) for an anonymous, non-full disclosure study. We set to sample a population of US residents using Amazon Mechanical Turk that is 18 years or above old, is a Twitter user, and has encountered at least one Tweet into their feed that relates to COVID-19 vaccines. There were both reputation and attention checks to prevent from bots and poor responses. The survey took between 5 and 10 minutes and the participants were compensates with the standard rate for participation. The study questionnaire, incorporating the instruments from [2, 7], is provided in the Appendix. We utilized a 2x3 between group experimental design where participants were randomized into one of six groups: (1) misleading Tweet with a warning tag; (2) misleading Tweet without a warning tag; (3) misleading Tweet with a warning tag; and (6) verified Tweet with a warning cover.

After participation, the participants were debriefed and offered the option to revoke their answers. We crafted the content of the Tweets to be of relevance to the participants such that they meaningfully engage with the Tweet's content (i.e., their responses are not arbitrary). We assumed participants understood the Twitter interface and metrics and were aware of the COVID-19 pandemic in general. However, we acknowledge that the level of interest regarding the COVID-19 vaccines could vary among the individual participants, affecting the extent to which their responses reflect their opinions.

#### 4 STUDY RESULTS

We conducted an online survey (N = 319) in January and February 2021. The power analysis conducted with  $G^*$  Power 3.1 [12] revealed that our sample was large enough to yield valid results for both Wilcoxon–Mann–Whitney U-test comparing two groups and Pearson's correlation (minimum of 44 per group). There were 180 (56.4%) males and 133 (41.7%) females, with 6 participants (1.8%) identifying as trans males, non-binary or preferring not to answer. The age brackets in the sample were distributed as follows: 13 (4.1%) [18 - 24], 108 (33.9%) [25 - 34], 98 (30.1%) [35 - 44], 46 (14.1%) [45 - 54], 37 (11.6%) [55 - 64], 17 (5.3%) [65 - 74], and 2 (.6%) [75 - 84]. In terms of education, 4 (1.3%) had less than high school, 24 (7.6%) had a high school degree or equivalent, 54 (17.0%) had some college but no degree, 43 (13.6%) had a 2-year degree, 139 (43.8%) had a 4-year college degree, and 53 (16.7%) had a graduate or professional

degree. Our sample, while balanced on the other demographics, was Democrat-leaning with 59 (18.5%) Republicans, 157 (49.2%) Democrats, and 102 (32%) Independent participants.

#### 4.1 Belief Echoes: Preconditions

The test of preconditions for formation of belief echoes on Twitter about the COVID-19 vaccine we first hypothesized that the presence of a warning tag under a Tweet containing *misleading* information about COVID-19 vaccines will not reduce the perceived accuracy of the Tweet's content relative to a no warning tag condition. Indeed, the Mann-Whitney U test comparing the perceived accuracy between the participants exposed to the Tweet from Figure 1a and Figure 1b, respectively was insignificant, as shown in Table 1. Confirming the H1 hypothesis, the results suggest that the warning tags didn't have the intended effect on reducing the perceived accuracy of a misleading COVID-19 vaccine information on Twitter. This proves the existence of preconditions for echoing one's belief irrespective of the warning tags.

However, this is not the case with the warning covers used to label misleading information. The Mann-Whitney U test comparing the perceived accuracy between the participants exposed to the Tweet from Figure 1b and Figure 2 was significant, rejecting the H2 hypothesis. The participants exposed to the warning cover (Figure 2) reported, on average, that the Tweet was "not very accurate" while the ones exposed only to the Tweet content (Figure 1b) that the Tweet was "somewhat accurate." The warning covers showed the intended effect of decreasing the perceived accuracy of misleading COVID-19 vaccine information on Twitter, dispelling one's echo beliefs for Tweets labeled with covers.

Table 1: Preconditions Tests: Hypotheses H1 to H4

	U-test	Significance	
H1	U = 1620.5	p = .217	
H2	U = 1841	$p = .004^*$	
Н3	<b>H3</b> $U = 1124$ $p = .063$		
<b>H4</b> $U = 1123.5$ $p = .087$			
Significance Level: $\alpha = 0.05$			

As we suspected, the warning labels didn't reduce the verified COVID-19 information, regardless of whether a warning cover or warning tag was presented as a soft moderation intervention. The Mann-Whitney U tests comparing the perceived accuracy between the participants exposed to the Tweet from Figure 3b and Figure 3a was insignificant as was the comparison between the exposures between Figure 3a and Figure 4 (Table 1). The retaining of H3 and H4 hypothesis suggest that participants critically discern the content of a seemingly valid Tweet, rather than the soft moderation labeling. Considering the previous reservations about the soft moderation implemented by Twitter when it comes to verified information [20], these results suggest that preconditions for echoing one's beliefs exist irrespective of any warning labeling also in the case where these beliefs are rooted in verifiable facts about COVID-19.

## 4.2 Belief Echoes: Safety and Herd Immunity

With the evidence of existing preconditions of belief echoes about the COVID-19, both for misleading and verified information, we set to explore how these belief echoes materialize. We asked the participants to what extent do they agree with the following statement: "I am not favorable to the COVID-19 vaccines because I believe they are unsafe". We found negative correlation with the perceived accuracy for the misleading Tweet with a warning tag (Figure 1a) and without a warning tag (Figure 1b) as shown in Table 2. The less participants were in favor of COVID-19 vaccines, the more accurate they perceived the misleading information regardless of the presence of a warning tag. We haven't found a significant correlation in the warning cover condition.

Table 2: Safety and Perceived Accuracy Tests: H5a/b

	r-test	Significance
H5a - with a warning tag	r =612	$p = .000^*$
H5a - without a warning tag	r =329	$p = .017^*$
H5b - with a warning tag	r =344	$p = .011^*$
H5b - with a warning cover	r =473	$p = .000^*$
Significance Level: $\alpha = 0.05$		

The test of the same relationship for the case of the verified Tweet also revealed a negative correlation with the Tweet's perceived accuracy in both the warning tag (Figure 3b) and the warning cover condition (Figure 4). The more participants were in favor of COVID-19 vaccines the less accurate they perceived the verified information regardless if there was a warning tag or a warning cover. On a first thought, this might be a surprising result, but a careful consideration indicates a presence of a belief echo and resistance to *soft moderated* verified information from [42], but not the standalone verified Twitter content.

Next, we asked the participants to what extent do they agree with the following statement: "There is no need to vaccinate for COVID-19 because I believe a natural herd immunity exists." We found negative correlation with the perceived accuracy for the misleading Tweet with a warning tag (Figure 1a) and without a warning tag (Figure 1b) as shown in Table 3. The more participants were in favor of the COVID-19 herd immunity, the more accurate they perceived the misleading information regardless if there was or there wasn't a warning tag. The test of the verified Tweet also revealed a negative correlation with the Tweet's perceived accuracy only in the original, no warning labels condition (Figure 3a). The more the participants were in favor of COVID-19 herd immunity the more accurate they perceived the verified information when no "soft moderation" intervention was applied. This finding adds to the evidence on the existence of echo beliefs in the essential need of COVID-19 vaccines as a preferred way of immunization.

## 4.3 Belief Echoes: Efficacy and Hesitancy

To test the subjective attitudes towards COVID-19 immunization, we asked the participants "will it be possible to produce efficacious COVID-19 vaccines". We found a significant result for the misleading Tweet with a warning tag (Figure 1a) and with a warning cover (Figure 2) as shown in Table 4. In both cases, the participants who didn't believe in efficacious COVID-19 vaccines perceived the misleading Tweet "somewhat acurate" while the participants that did

Table 3: Herd Immunity and Perceived Accuracy Tests: H6a/b

	r-test	Significance
H6a - with a warning tag	r =529	$p = .000^*$
H6a - without a warning tag	r =387	$p = .005^*$
H6b - without a warning tag	r =445	$p = .001^*$
Significance Level: $\alpha = 0.05$		

believe perceived it as "not very accurate." We found a significant result for the verified Tweet in its original form, without any "soft moderation" (Figure 3a). In this case, the participants that didn't believe in efficacious COVID-19 vaccines perceived the original, non-moderated Tweet "not very accurate." while the ones that did believe in the efficacy perceived it as "somewhat accurate."

Table 4: Efficacy and Perceived Accuracy Tests: H7a/b

	r-test	Significance	
H7a - with a warning tag	$\chi^2(2) = 8.566$	p = .003*	
H7a - with a warning cover	$\chi^2(2) = 9.237$	$p = .002^*$	
<b>H6b</b> - without a warning tag $\chi^2(2) = 3.969$ $p = .005^*$			
Significance Level: $\alpha = 0.05$			

To test the personal hesitancy to COVID-19 immunization, we asked the participants "Will you get vaccinated, if possible?" We found a significant result for the misleading Tweet with a warning tag (Figure 1a) and with a warning cover (Figure 2) as shown in Figure 5. In both cases, the participants that were hesitant to receive the COVID-19 vaccine perceived the misleading Tweet "somewhat accurate" while the participants that want to receive the vaccine as "not very accurate." The participants that were unsure, perceived the misleading Tweet in both cases as "not at all accurate.". We found a significant result for the verified Tweet in its original form, without any "soft moderation" (Figure 3a). In the case, the participants that didn't want to receive the COVID-19 vaccine perceived the original, non-moderated Tweet "somewhat accurate." while the participants that wanted to receive the vaccine perceived it as "not very accurate." Similarly, the participants that were unsure perceived the Tweet as "not at all accurate."

Table 5: Personal Hesitancy and Perceived Accuracy Tests: H8a/b

	r-test	Significance	
H8a - with a	$\chi^2(2) = 9.381$	n - 000*	
warning tag	$\chi$ (2) = 9.381	$p = .009^*$	
H8a - with a	$\chi^2(2) = 7.163$	n — 020*	
warning cover	$\chi$ (2) = 7.163	$p = .028^*$	
H8b - without a	.2(0) 12.512 - 001*		
warning tag	$\chi^2(2) = 13.513$	$p = .001^*$	
Significance Level: $\alpha = 0.05$			

To test the hesitancy to COVID-19 immunization for children, we asked the participants "Should children be vaccinated for COVID-19 too?" We found a significant result for the misleading Tweet with only the warning tag (Figure 1a) as shown in Table 6. The participants that were hesitant to administer the COVID-19 vaccine to children perceived the misleading Tweet with a warning tag "somewhat accurate" while the participants that agreed with administering the COVID-19 vaccine to children "not very accurate." We found a significant result for the verified Tweet in its original form, without any "soft moderation" (Figure 3a). In the case, the participants that were hesitant to administer the COVID-19 to children perceived the an original, non-moderated Tweet as "somewhat accurate." while the participants that agreed to administer the vaccine to children perceived it as "not very accurate."

Table 6: Hesitancy for Children and Perceived Accuracy Tests: H9a/b

	r-test	Significance	
H9a - with a warning tag	$\chi^2(2) = 10.663$	$p = .001^*$	
H9b - without a warning tag $\chi^2(2) = 8.001$ $p = .005^*$			
Significance Level: $\alpha = 0.05$			

#### 4.4 Belief Echoes and Political Affiliation

Following the association between one's political affiliation and the warnings of a misleading Twitter content [45], we analyzed the perceived accuracy among the participants based on their political affiliation (Republican, Democrat, independent). We found a significant difference in perception between the political affiliations of the participants for the misleading Tweet with the warning tag and without the warning tag as shown in Table 7. In both cases, the Republicans and independent participants perceived the Tweet as "somewhat accurate" while the Democrats as "not very accurate."

Table 7: Political Affiliation and Perceived Accuracy Tests: RO1

	r-test	Significance	
Misleading Tweet with a warning tag	$\chi^2(2) = 7.063$	$p = .029^*$	
Misleading Tweet without a warning tag	$\chi^2(2) = 9.127$	$p = .005^*$	
Significance Level: $\alpha = 0.05$			

We also analyzed the association between political affiliation, the beliefs of safety and herd immunity, and the subjective attitudes on the vaccine efficacy and hesitancy. While there were no significant correlations about the safety and herd immunity beliefs, we found significant differences on the question of producing efficacious vaccines, personal hesitancy, as shown in Table 8. Almost one in four Republicans and and one in six Independents don't expect to have an efficacious COVID-19 vaccine (Table 9), while that proportion

for the Democrats is one in forty. Half the Republicans and a third of the Independents are hesitant to receive the COVID-19 vaccine (Table 10), while only a tenth of the Democrats won't proceed with personal immunization. Roughly 40% of the Republicans and Independents are hesitant to vaccinate children for COVID-19, to which only 8.3% of the Democrats agree with (Table 11).

Table 8: Political Affiliation, Beliefs, and Subjective Attitudes Tests: RQ2

	r-test	Significance	
Producing efficacious	$\chi(1) = 22.059$	$p = .001^*$	
vaccines Personal Hesitancy	$\chi(2)(2) = 55.486$	$p = .000^*$	
Hesitancy for Children $\chi(1) = 45.665$ $p = .000^*$			
Significance Level: $\alpha = 0.05$			

Table 9: Pearson Chi-Square Test - Political Affiliation vs Production of an Efficacious Vaccine

	Republicans	Democrats	Independent
Agree	46 (78.0%)	153 (97.5%)	86 (89.6%)
Disagree	13 (22.0%)	4 (2.5%)	16 (15.7%)

Table 10: Pearson Chi-Square Test - Political Affiliation vs Personal Vaccination Hesitancy

	Republicans	Democrats	Independent
Certain	28 (47.5%)	136 (86.6%)	63 (61.8%)
Hesitant	27 (45.8%)	8 (5.1%)	35 (34.3%)
Undecided	4 (6.8%)	13 (8.3%)	4 (3.9%)

Table 11: Pearson Chi-Square Test - Political Affiliation vs Children Vaccination Hesitancy

	Republicans	Democrats	Independent
Certain	34 (57.6%)	144 (91.7%)	61 (59.8%)
Hesitant	25 (42.4%)	13 (8.3%)	41 (40.2%)

#### 5 DISCUSSION

Consistent with the previous evidence on receptivity to misinformation and resistance to warnings [7, 32], we found that the more likely participants were to believe that COVID-19 vaccines are unsafe, the more receptive they were to misleading COVID-19 vaccines information from Twitter, resisting the soft moderation intervention, proving the existence of belief echoes.

#### 5.1 Strength and Type of Warning Label

That the participants perceived the misleading Tweets as accurate in the presence of a warning tag but not in the presence of a cover condition suggests that the warning tags are not effective or insufficient to sway participants' perceived accuracy. These findings are consistent with previous research showing that the design of the warning label affects individuals' perceptions of the content [7, 22, 29, 35]. The design of warning labels and how they are presented to users can impact the warning labels' effectiveness, with more explicit labels being more effective [29]. For instance, it was found that individuals routinely ignored contextual warnings, which were akin to the warning tags in our study in that they did not obscure the misleading content nor require individuals to click through to see the information [22].

However, interstitial warnings, akin to cover warnings in our study in that they required individuals to click through to continue, were effective in countering disinformation. Authors in [22] posit that this may be because interstitial designs are more noticeable for users and thereby more effective at countering misinformation. It may also be that these designs require users' engagement and thus necessitate a cognitive awareness of the tag's content. Similarly, authors in [7] found that the perceived accuracy decreased with the increasing strength of the warning tags, such that tags which said "rated false" were significantly more effective than "disputed" tags at reducing beliefs in the misleading information. These findings suggest that more explicit, unambiguous warnings are more effective at countering misleading information. In our case, the warning covers - not the tags, which are more verbose and exact, were a more potent way of urging users to critically discern COVID-19 vaccine content on Twitter. However, users may habituate to the cover warnings in the long-term if they perceive that the moderator, Twitter, is biased in labeling content from users with particular political leanings [4].

# 5.2 Preconceived Notions-Explaining Results of Verified Tweets

Dispelling belief echoes on Twitter might be a more complex task, dependent on the content or type of misinformation and the subjective involvement of the participants. The fact that participants who were more likely to believe that COVID-19 vaccines were unsafe, were less likely to perceive the verified Tweet as inaccurate in both the warning tag and cover conditions may be due to the fact that the verified Tweet, though accurate, still reflected information that was negative about COVID-19 vaccines (i.e., invoking the idea that they may lead to serious side-effects and should be avoided by some participants in some situations). In other words, it would make sense that participants favouring vaccines would be more likely to disbelieve the Tweet expressing a concern about COVID-19 vaccines when accompanied by a warning label. A similar conclusion could be drawn also in the case of belief in herd immunity versus mass immunization with COVID-19 vaccines.

#### 5.3 COVID-19 Vaccine Beliefs

In terms of efficacy, participants who thought that the COVID-19 vaccines were ineffective were more likely to rate the misleading Tweets as accurate, regardless of the soft moderation applied. This suggests that belief echoes persist despite the warnings, and perhaps, a presence of a warning tag or cover may actually increase people's likelihood of finding a misleading Tweet accurate if it conforms to existing beliefs. This finding is consistent with the backfire effect previously observed for polarizing content on social media [33]. For example, evidence suggests that corrections on misleading Tweets strengthened misperceptions (or perceptions of accuracy) among those most strongly committed to the belief [32]. The corrections that contradict users' preconceived notions were found to lead individuals to double down on their beliefs.

In terms of hesitancy (both personal and for children vaccination), we found a similar effect, such that those who were hesitant about vaccines were more likely to perceive misleading Tweets with tags and covers as accurate while those who wanted the vaccine perceived it as inaccurate. Again, the fact that only Tweets with tags and covers were viewed as accurate suggests evidence for a backfire effect such that the mere presence of the tags/covers may increase individuals' beliefs in the Tweets' accuracy if the content reinforces the participants' anti-vax stance. A similar conclusion follows from the fact that, for the original (verified) Tweet, the pro-vax participants who wanted to get the vaccine viewed it as not very accurate but the anti-vax participants viewed it as accurate.

#### 5.4 Political Affiliations

Along the lines of the findings in [6, 32, 45], we found further evidence of the association between user's political affiliations and the receptivity to misleading content. The Republican and Independent participants perceived the Tweet as "somewhat accurate" while the Democrat participants perceived it as "not very accurate" in both the misleading Tweet with and without a warning tag. That the difference between the expectation of an efficacious COVID-19 vaccine is twentyfold between the Republicans and Democrats is a bit surprising, but consistent with the breakdown of trust in scientists to deliver an efficacious COVID-19 vaccine along the party line [17]. The hesitancy we found in our study is consistent with the previous reported breakdown for the COVID-19 vaccine hesitant Republicans and Democrats, both personally and in regards to children's vaccination [23]. Interestingly, Independents showed a high hesitancy on par with the Republicans in both cases.

Authors in [6] found that while corrections had the intended effect among Democrats, soft moderation techniques backfired among Republicans. Specifically, the authors found that while corrections of misleading claims decreased Democrats' perceptions of claim accuracy, they actually strengthened Republicans' perceptions of accuracy. As in [32], these findings suggest that corrections of misleading information on social media may not only be ineffective among some individuals but may actually reinforce individuals' preconceived notions. While our study did not assess participants' beliefs before and after receiving corrections, as all participants were only assessed once, the findings that political affiliation affects individuals' perceptions of accuracy and the impact that warning labels have on those perceptions are consistent with the backfiring effect among individuals with certain political ideologies.

#### 5.5 Usable Security and Privacy Implications

Reluctance to heed security or privacy warnings is not a new phenomenon and has been well researched in the past [11, 18, 31]. While efforts have been invested in increasing the clarity of the messages and design of affordances to attract attention and motivate users, habituation is a complex problem transcending security designs. Habituation describes a diminished response with repetitions of the same stimulus, decreasing the intended effect of security and privacy warnings among users. Authors in [43], in this context, have uncovered the phenomenon of "generalization" where habituation to one stimulus carries over to other novel stimuli that are similar in appearance. We didn't explore the diminished response with repetitions of the same warning label to a Tweet, being that a tag or a cover, but that certainly warrants close research attention. The findings of our study suggest that heeding a misleading information warning only happened when the information is obscured by a plain text warning of the risks, not when the warning follows the Tweet with tag.

The warning tag, consisting of an exclamation mark symbol urges users to "Get the facts about COVID-19," in Twitter's blue font, communicates a seemingly ambiguous message without explicitly addressing that the Tweet's content aims to mislead users about the COVID-19 vaccine. Perhaps a line of research could explore a variation of tags that are more direct, for example "This is COVID-19 misinformation", written in bold red font and conventional warning favicons. Alternatively an impartial message like "No judgment, but this might be COVID-19 misinformation" could also show users' receptivity to not-so-overt moderation. The warning cover, along these lines, communicated a message where Twitter appeared not taking sides by saying: This Tweet violated the Twitter Rules about spreading misleading and potentially harmful information related to COVID-19. Twitter has determined that may be in the public interest for the Tweet to Remain Accessible. It also provided a link for the participant to Learn More, which largely subsumes the warning tag by leading users to a repository of verified COVID-19 information. Alternative wording like This Tweet was rated 'false', but we keep it in the public interest, based on the previous evidence [7], could yield a stronger reduction of misperception.

In the context of generalization, Twitter labels alternative narratives that are: (1) statements or assertions that have been confirmed to be false or misleading by subject-matter experts, such as public health authorities (misleading information); (2) statements or assertions in which the accuracy, truthfulness, or credibility of the claim is contested or unknown (disputed claims); and (3) information (which could be true or false) that is unconfirmed at the time it is shared (unverified claims)[34]. With a similarity in labeling alternative narratives, from a user perspective, habituation to a tag or a cover for misleading COVID-19 vaccination could potentially carry over to other warnings about other polarizing events, such as elections. An interesting line of research could investigate the generalization effect not just between different labels, but in various combinations of formatting and wording. A user might be well aware and agree that some claims about elections are disputed, but they can nonetheless retain their beliefs about the COVID-19 vaccine safety and efficacy. Further so, another line of research could

be an investigation of the generalization phenomenon between social media platforms, for example between warning labels on Twitter and Facebook.

This discussion brings an important aspects of usable security affordances that depart from the conventional warning about systemlevel exploits towards content-level warnings. System-level exploits hardly relate to any potent beliefs (outside perhaps of the stereotypical foreign nation-state interference) or better said, users might not have strong polarizing stances on phishing or malware, usually perceiving it as a "bad thing" [14]. Content-level exploits, on the other side, are far more complex and potent in polarizing users, given that they are subjectively involved with the content [39]. Users might ignore a red screen proceeding to a suspicious website, but they usually trust Chrome or Firefox that they have honest intentions in warning them about potential risks. Evidence already indicates that users are not trusting of the soft moderation intervention, feeling that Twitter itself was biased and mislabeling content [20]. Remaining impartial while trying to dispel belief echoes might be a harder problem depending on the content - while there are safe and unsafe websites, there is, and will continue to be, a wealth of polarizing content on Twitter that will require content-relevant warning labeling.

## 5.6 Ethical Implications

While we set out to investigate the effect of soft moderation on Twitter and debriefed the participants at the end of the study, the results could have several ethical implications, nonetheless. We exposed the participants to a misleading and manipulated soft moderation of twitter content on the topic of the COVID-19 vaccine that could potentially affect participants' stance on vaccination and the pandemic. The exposure might not sway participants on the hesitancy or their perceptions of safety and efficacy, but could make the participants reconsider their approach of obtaining the vaccine for themselves or their children. The exposure could also affect the participants' stance of social media soft moderation in general and nudge people to move to less regulated platforms [46]. A recent example of such a migration from Twitter to Parler, Rumble and Newsmax was witnessed after Twitter actively labeled and removed false information on the platform during the 2020 U.S. elections [21].

That the participants were able to critically discern the content of the verified Tweet despite our alternation to include warning labels is reassuring and suggests that misinformation has the potential to be contained, if not eradicated, from social media platforms. However, the potential of crafting software that could silently attach or remove warning covers before they are presented to Twitter users could have unintended consequences. In the past, such an effort was tested in manipulating a Twitter textual content (not any additional affordances in the user interface) to induce misperceptions about the relationship between vaccines and autism [36]. With the evidence of nation-states censoring Twitter regarding narratives countering their interest in the past, it is possible that such a nation-state could use a similar approach and implement a "post-soft moderation" logic within a state-approved and disseminated social media application [41]. This may be far from the realm of possibility, even if the capabilities exist, but for such a sensitive

topic as COVID-19 vaccination, meddling with the warning labels could give an edge to a vaccine competitor in the global race for development and procurement of COVID-19 vaccines. We condemn such ideas and use of our research results. Evidence for such a nefarious misinformation Twitter campaign that promotes homegrown Russian vaccine and undercuts rivals has already surfaced [16].

Perhaps outside the scope of this study, the ethical questions remains whether Twitter, or any social media platform, acting as a private entity, could set a precedent of an ultimate arbiter of what constitutes misinformation and what does not. Twitter most likely applies an automated means of warning labeling in conjunction with manual moderation, as evidence with the strange labeling of Tweets that contained the words "oxygen" and "frequency" for COVID-19 related Tweets [45]. Even with an attempt at honest moderation, cross-checked with the health authorities like CDC, a potential problem might arise in case a previously held belief, or a fact about COVID-19 is later disputed. For example, at the beginning of the pandemic, authorities claimed masks were not effective in protecting the virus from spreading, a claim that later was not reversed, resulting in masks becoming essential to any human-to-human interaction [47]. So if the warning labels were applied to moderate any Tweet that contains the words "mask" and "stop" or "spread" at the early periods of the pandemic, they must be retracted. Such a thing could cast doubt on studies like ours, even if we as researchers, and Twitter as moderators, acted in good faith. Certainly, this could damage the reputation of users as well as Twitter and further exacerbate the impression of not-so-honest impartiality in labeling content as misleading, especially against users identifying themselves as conservatives [4].

#### 5.7 Future Research

We acknowledge that there is further research to be done into investigating the full ramifications of soft moderation by social media platforms, especially beyond the topics of the COVID-19 pandemic or presidential elections. A promising line of research is the combination of soft and hard moderation, given that Twitter has exercised the right to ban or suspend accounts indefinitely that have been labeled for misinformation in the past, like in case of Donald Trump. One could probe the warning labeling algorithm and reverse engineer it to find if there is a relationship between the number/type of warning labels an account could receive, for example, before it gets permanently banned (if an automated ban exists, given the wide latitude and the imperative of Twitter to remain not overly controlling of the public discourse). When social media platforms, as private entities, are predetermined by users to hold biased positions [20], any action taken to apply soft moderation techniques may be undermined in the process, instead working to legitimize the beliefs of skeptics. In this direction, a content analysis of the Tweets being labeled could reveal the topics, images, words, or the network of accounts behind such impressions.

In the United States, where the right to speech is protected, more research may be done to see how alternative narratives, belonging to a same type of content or a topic (e.g. COVID-19 vaccines cause adverse effects leading to death) are soft moderated between platforms, for example comparing Twitter, Facebook, or Parler. Soft moderated content is usually closely related to content used for

trolling so this relationship could be also explored, such as understanding if warning labeled Tweets provoke emotional response and of what kind. Similarly, the warning labeling can be associated to identify the evolution of political information operations on Twitter, that have been waged on the topic of COVID-19 already [40].

## 5.8 Scope Limitations

The current study has important limitations. First, it is possible that a different topic or even different information regarding the effect of the COVID-19 vaccines would have different outcomes. We used a couple of Tweets that were relevant to the state of the pandemic and mass immunization during the period of January-February 2021, which could be perceived with a different level of accuracy after a certain period of time. We used only two Tweets and a study that explores the effect on multiple misleading or verified Tweets on COVID-19 vaccine could uncover different efficiency or strength of soft moderation. Overall, the findings in the present study may be specific to the effect that warning tags have on COVID-19 (mis)information and cannot be generalized to other topics. Second, participants who are more regular social media users in general may be desensitized to the information presented in the Tweets, which may have affected their perception of the issue irrespective of the warnings.

Third, our experiment was limited to Twitter as a social media platform of choice. Since the content and images we present are borrowed and adapted to the study objectives from Twitter, we are limited to evaluating the effects of the warning tags and warning cover on perceptions of accuracy on Twitter only and may not be generalized regarding other social media platforms. We were limited to the formatting and wording of the warning labels chosen by Twitter at the time of the study. If Twitter chooses to place the tag, say on top of the Tweet instead of the bottom, the results could be different. Similarly, if the wording of the warning cover changes, the results might not hold for such new conditions. Fourth, we did not examine the effects over a period of time. Thus, we are unable to examine the Tweet's effects following the study. We also acknowledge another limitation imposed of the timeline of the study and the speed of COVID-19 vaccine development. By the time participants completed the study, much more might be known about the COVID-19 vaccines to sway public opinion. For instance, if many participants have gotten the vaccine without major side effects by the time participants complete the study, this might affect their responses. Fifth, although we tried to sample a representative set of participants for our study using Amazon Mechanical Turk, the outcomes might have been different if we used another platform, or other type of sampling. Also, a much larger sample size could have provided a more nuanced view of the soft moderation, but we had limited funding for this study.

#### 6 CONCLUSION

In the present study, we sought to determine whether two forms of soft moderation on Twitter affect the perceived accuracy of Tweets pertaining to COVID-19 vaccines. We were also interested in examining whether perceived Tweet accuracy varies based on individual's' beliefs regrading COVID-19 vaccine safety, efficacy, willingness to receive vaccinations, and their political affiliations.

Overall, our results suggest that warning covers are more effective than warning tags in dispelling individuals' beliefs about misleading Tweets. Individuals' pre-existing beliefs regarding COVID-19 vaccine safety, efficacy, and hesitancy affect individuals' perceptions of Tweet accuracy such that their perceptions of the accuracy align with their biases. Furthermore, our results also show that individuals' political affiliations also affect their perceptions of accuracy for misleading Tweets such that Republicans and Independents, who are more likely to express skepticism regarding vaccines, are more likely to perceive misleading Tweets as accurate, irrespective of any moderation effort.

In all cases, individuals perceive the Tweets in ways that are most favourable to or consistent with their pre-existing beliefs. This may lead to a backfire effect, as evidenced by the fact that individuals who were skeptical of vaccines were more likely to rate misleading Tweets with tags and covers, but not misleading Tweets without tags, as accurate. Taken together, our findings provide additional evidence for the existence of belief echoes pertaining to COVID-19 vaccines that are largely resistant to soft moderation in the form of warning tags but not warning covers. We believe that the insight gained from this research regarding how individuals' pre-existing belief biases impact perceptions of Tweet accuracy in the context of soft moderation can be used to develop more effective moderation techniques that do minimize unintended consequences. We hope that our results could inform the usable security community towards future steps in eradicating misinformation on Twitter and social media in general.

#### REFERENCES

- [1] Jon-Patrick Allem and Emilio Ferrara. 2018. Could social bots pose a threat to public health? *American journal of public health* 108, 8 (2018), 1005.
- [2] Luigi Roberto Biasio, Guglielmo Bonaccorsi, Chiara Lorini, and Sergio Pecorelli. 2020. Assessing COVID-19 vaccine literacy: A preliminary online survey. *Human Vaccines & Immunotherapeutics* 0, 0 (2020), 1–9. https://doi.org/10.1080/216455 15.2020.1829315
- [3] David A Broniatowski, Amelia M Jamison, SiHua Qi, Lulwah AlKulaib, Tao Chen, Adrian Benton, Sandra C Quinn, and Mark Dredze. 2018. Weaponized health communication: Twitter bots and Russian trolls amplify the vaccine debate. American journal of public health 108, 10 (2018), 1378–1384.
- [4] Jenna Burrell, Zoe Kahn, Anne Jonas, and Daniel Griffin. 2019. When Users Control the Algorithms: Values Expressed in Practices on Twitter. Proc. ACM Hum.-Comput. Interact. 3, CSCW, Article 138 (Nov. 2019), 20 pages. https://doi.org/10.1145/3359240
- [5] Bill Chappell. 2021. Instagram Bars Robert F. Kennedy Jr. For Spreading Vaccine Misinformation. (2021). https://www.npr.org/sections/coronavirus-live-update s/2021/02/11/966902737/instagram-bars-robert-f-kennedy-jr-for-spreading-vaccine-misinformation.
- [6] Dino P Christenson, Sarah E Kreps, and Douglas L Kriner. 2020. Contemporary Presidency: Going Public in an Era of Social Media: Tweets, Corrections, and Public Opinion. Presidential Studies Quarterly (2020).
- [7] Katherine Clayton, Spencer Blair, Jonathan A Busam, Samuel Forstner, John Glance, Guy Green, Anna Kawata, Akhila Kovvuri, Jonathan Martin, Evan Morgan, et al. 2019. Real solutions for fake news? Measuring the effectiveness of general warnings and fact-check tags in reducing belief in false stories on social media. *Political Behavior* (2019), 1–23.
- [8] Warren Cornwall. 2020. Officials gird for a war on vaccine misinformation. Science 369, 6499 (2020), 14–15. https://doi.org/10.1126/science.369.6499.14 arXiv:https://science.sciencemag.org/content/369/6499/14.full.pdf
- [9] Ullrich KH Ecker, Stephan Lewandowsky, and David TW Tang. 2010. Explicit warnings reduce but do not eliminate the continued influence of misinformation. Memory & cognition 38, 8 (2010), 1087–1100.
- [10] European Parliament. [n. d.]. Vaccine Hesitancy and Drop in Vaccination Rates in Europe - Thursday, 19 April 2018. https://www.europarl.europa.eu/doceo/doc ument/TA-8-2018-0188\_EN.html.
- [11] Michael Fagan and Mohammad Maifi Hasan Khan. 2016. Why Do They Do What They Do?: A Study of What Motivates Users to (Not) Follow Computer Security Advice. In Twelfth Symposium on Usable Privacy and Security (SOUPS)

- $2016). \ USENIX\ Association, Denver, CO, 59-75.\ https://www.usenix.org/conference/soups2016/technical-sessions/presentation/fagan$
- [12] Franz Faul, Edgar Erdfelder, Axel Buchner, and Albert-Georg Lang. 2009. Statistical power analyses using G\*Power 3.1: Tests for correlation and regression analyses. Behavior Research Methods 41, 4 (2009), 1149–1160. https://doi.org/10.3758/BRM.41.4.1149
- [13] Jieyu D. Featherstone, George A. Barnett, Jeanette B. Ruiz, Yurong Zhuang, and Benjamin J. Millam. 2020. Exploring childhood anti-vaccine and pro-vaccine communities on twitter – a perspective from influential users. *Online Social* Networks and Media 20 (2020), 100105. https://doi.org/10.1016/j.osnem.2020.100
- [14] Adrienne Porter Felt, Alex Ainslie, Robert W Reeder, Sunny Consolvo, Somas Thyagaraja, Alan Bettes, Helen Harris, and Jeff Grimes. 2015. Improving SSL warnings: Comprehension and adherence. In Proceedings of the 33rd annual ACM conference on human factors in computing systems. 2893–2902.
- [15] Centers for Disease Control and Prevention. 2021. COVID-19 Vaccines and Allergic Reactions. https://www.cdc.gov/coronavirus/2019-ncov/vaccines/safet y/allergic-reaction.html.
- [16] Sheera Frenkel, Maria Abi-Habib, and Julian E. Barnes. 2021. Russian Campaign Promotes Homegrown Vaccine and Undercuts Rivals. https://www.nytimes.co m/2021/02/05/technology/russia-covid-vaccine-disinformation.html.
- [17] Cary Funk and Alec Tyson. 2020. Intent to Get a COVID-19 Vaccine Rises to 60% as Confidence in Research and Development Process Increases. https: //www.pewresearch.org/science/2020/12/03/intent-to-get-a-covid-19-vaccinerises-to-60-as-confidence-in-research-and-development-process-increases/.
- [18] Simson Garfinkel and Heather Richter Lipford. 2014. Usable security: History, themes, and challenges. Synthesis Lectures on Information Security, Privacy, and Trust 5. 2 (2014), 1–124.
- [19] Bertram Gawronski, Roland Deutsch, Sawsan Mbirkou, Beate Seibt, and Fritz Strack. 2008. When "Just Say No" is not enough: Affirmation versus negation training and the reduction of automatic stereotype activation. Journal of Experimental Social Psychology 44, 2 (2008), 370–377. https://doi.org/10.1016/j.jesp.2006.12.004
- [20] Christine Geeng, Tiona Francisco, Jevin West, and Franziska Roesner. 2020. Social Media COVID-19 Misinformation Interventions Viewed Positively, But Have Limited Impact. arXiv:cs.CY/2012.11055
- [21] Mike Isaac and Kellen Browning. 2020. Fact-Checked on Facebook and Twitter, Conservatives Switch Their Apps. https://www.nytimes.com/2020/11/11/techn ology/parler-rumble-newsmax.html.
- [22] Ben Kaiser, Jerry Wei, Elena Lucherini, Kevin Lee, J. Nathan Matias, and Jonathan Mayer. 2020. Adapting Security Warnings to Counter Online Disinformation. https://arxiv.org/abs/2008.10772.
- [23] Kendall Karson. 2020. Americans willing to receive COVID-19 vaccine but divided on timing: Poll. https://abcnews.go.com/Politics/americans-receive-covid-19-v accine-divided-timing-poll/story?id=74703426.
- [24] Jessica Keim-Malpass, Emma M. Mitchell, Emily Sun, and Christine Kennedy. 2017. Using Twitter to Understand Public Perceptions Regarding the #HPV Vaccine: Opportunities for Public Health Nurses to Engage in Social Marketing. Public Health Nursing 34, 4 (2017), 316–323. https://doi.org/10.1111/phn.12318 arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1111/phn.12318
- [25] Antino Kim, Patricia L. Moravec, and Alan R. Dennis. 2019. Combating Fake News on Social Media with Source Ratings: The Effects of User and Expert Reputation Ratings. *Journal of Management Information Systems* 36, 3 (2019), 931–968. https://doi.org/10.1080/07421222.2019.1628921
- [26] Jeffrey V Lazarus, Scott C Ratzan, Adam Palayew, Lawrence O Gostin, Heidi J Larson, Kenneth Rabin, Spencer Kimball, and Ayman El-Mohandes. 2020. A global survey of potential acceptance of a COVID-19 vaccine. *Nature medicine* (2020), 1–4.
- [27] Joan Massachs, Corrado Monti, Gianmarco De Francisci Morales, and Francesco Bonchi. 2020. Roots of Trumpism: Homophily and Social Feedback in Donald Trump Support on Reddit. In 12th ACM Conference on Web Science. 49–58.
- [28] Philip M Massey, Amy Leader, Elad Yom-Tov, Alexandra Budenz, Kara Fisher, and Ann C Klassen. 2016. Applying Multiple Data Collection Tools to Quantify Human Papillomavirus Vaccine Communication on Twitter. J Med Internet Res 18, 12 (05 Dec 2016), e318. https://doi.org/10.2196/jmir.6670
- [29] Patricia L. Moravec, Antino Kim, and Alan R. Dennis. 2020. Appealing to Sense and Sensibility: System 1 and System 2 Interventions for Fake News on Social Media. *Information Systems Research* 31, 3 (2020), 987–1006. https://doi.org/10.1 287/isre 2020 0927
- [30] Adam Mosseri. 2016. Addressing Hoaxes and Fake News. Facebook (2016). https://about.fb.com/news/2016/12/news-feed-fyi-addressing-hoaxes-and-fake-news/.
- [31] James Nicholson, Lynne Coventry, and Pam Briggs. 2017. Can we fight social engineering attacks by social means? Assessing social salience as a means to improve phish detection. In *Thirteenth Symposium on Usable Privacy and Security* (SOUPS 2017). USENIX Association, Santa Clara, CA, 285–298. https://www.us enix.org/conference/soups2017/technical-sessions/presentation/nicholson
- [32] Brendan Nyhan and Jason Reifler. 2010. When corrections fail: The persistence of political misperceptions. *Political Behavior* 32, 2 (2010), 303–330.

- [33] Gordon Pennycook, Adam Bear, Evan T Collins, and David G Rand. 2020. The implied truth effect: Attaching warnings to a subset of fake news headlines increases perceived accuracy of headlines without warnings. Management Science (2020).
- [34] Yoel Roth and Nick Pickles. 2020. Updating our approach to misleading information. Twitter (2020). https://blog.twitter.com/en\_us/topics/product/2020/updating-our-approach-to-misleading-information.html.
- [35] Haeseung Seo, Aiping Xiong, and Dongwon Lee. 2019. Trust It or Not: Effects of Machine-Learning Warnings in Helping Individuals Mitigate Misinformation. In Proceedings of the 10th ACM Conference on Web Science (WebSci '19). Association for Computing Machinery, New York, NY, USA, 265–274. https://doi.org/10.114 5/3292522.3326012
- [36] Filipo Sharevski, Peter Jachim, and Kevin Florek. 2020. To Tweet or Not to Tweet: Covertly Manipulating a Twitter Debate on Vaccines Using Malware-Induced Misperceptions. In Proceedings of the 15th International Conference on Availability, Reliability and Security (ARES '20). Association for Computing Machinery, New York, NY, USA, Article 75, 12 pages. https://doi.org/10.1145/3407023.3407025
- [37] Jeff Smith. 2017. Designing Against Misinformation. Medium (2017). https://medium.com/facebook-design/designing-against-misinformation-e5846b3aa1e2.
- [38] Todd Spangler. 2020. Twitter Flags 200 Trump Posts as False or Disputed Since Election Day - Variety. Variety (2020). https://variety.com/2020/digital/news/tw itter-trump-200-disputed-misleading-claims-election-1234841137/.
- [39] Leo G Stewart, Ahmer Arif, and Kate Starbird. 2018. Examining trolls and polarization with a retweet network. In Proc. ACM WSDM, workshop on misinformation and misbehavior mining on the web.
- [40] Benjamin Strick. 2020. Discovering A Pro-Chinese Government Information Operation On Twitter and Facebook: Analysis Of The #MilesGuo Bot Network. https://www.bellingcat.com/news/2020/05/05/uncovering-a-pro-chinese-government-information-operation-on-twitter-and-facebook-analysis-of-the-milesguo-bot-network/.
- [41] Kurt Thomas, Chris Grier, and Vern Paxson. 2012. Adapting social spam infrastructure for political censorship. In 5th {USENIX} Workshop on Large-Scale Exploits and Emergent Threats ({LEET} 12).
- [42] Emily Thorson. 2016. Belief echoes: The persistent effects of corrected misinformation. Political Communication 33, 3 (2016), 460–480.
- [43] Anthony Vance, David Eargle, Jeffrey L. Jenkins, C. Brock Kirwan, and Bonnie Brinton Anderson. 2019. The Fog of Warnings: How Non-essential Notifications Blur with Security Warnings. In Fifteenth Symposium on Usable Privacy and Security (SOUPS 2019). USENIX Association, Santa Clara, CA. https://www.usenix.org/conference/soups2019/presentation/vance
- [44] Robin C. Vanderpool, Anna Gaysynsky, and Wen-Ying Sylvia Chou. 2020. Using a Global Pandemic as a Teachable Moment to Promote Vaccine Literacy and Build Resilience to Misinformation. American Journal of Public Health 110, S3 (2020), S284–S285. https://doi.org/10.2105/AJPH.2020.305906
   [45] Savvas Zannettou. 2021. "I Won the Election!":An Empirical Analysis of Soft
- [45] Savvas Zannettou. 2021. "I Won the Election!": An Empirical Analysis of Soft Moderation Interventions on Twitter. arXiv 2101.07183v1 (18 January 2021). https://arxiv.org/pdf/2101.07183.pdf.
- [46] Savvas Zannettou, Tristan Caulfield, Emiliano De Cristofaro, Nicolas Kourtelris, Ilias Leontiadis, Michael Sirivianos, Gianluca Stringhini, and Jeremy Blackburn. 2017. The Web Centipede: Understanding How Web Communities Influence Each Other through the Lens of Mainstream and Alternative News Sources. In Proceedings of the 2017 Internet Measurement Conference (IMC '17). Association for Computing Machinery, New York, NY, USA, 405–417. https://doi.org/10.114 5/3131365.3131390
- [47] Joyce C Zhang and Anil Adisesh. 2020. Controversies in respiratory protective equipment selection and use during COVID-19. Journal of hospital medicine 15, 5 (2020).

#### APPENDIX

The study questionnaire included the following questions:

#### • Perceived Accuracy of a Tweet:

- 1. To the best of your knowledge, how accurate is the claim described in the Tweet?
- 4-point Likert scale (1-not at all accurate, 2-not very accurate, 3-somewhat accurate, 4-very accurate).

#### • Beliefs:

- 2. How much do you agree with the following statement:"I am not favorable to vaccines because they are unsafe"?
- 3. How much do you agree with the following statement: "here is no need to vaccinate because a natural immunity exists"? 4-point Likert scale (1 Totally, 2 A Little, 3 Partially, 4 Not at All).

## • Subjective Attitudes:

4. Will it be possible to produce safe and efficacious COVID-19 vaccines?

Yes/No.

5. Will you get vaccinated, if possible?

Yes/No/I Don't Know.

6. Should children be vaccinated for COVID-19 too? Yes/No.

## • Demographics:

Age, gender identity, education, political leanings.