Cross-Partisan Interactions on Twitter

Yusuf Mücahit Çetinkaya^{*1,2}, Vahid Ghafouri^{*3,4,5}, Guillermo Suarez-Tangil⁴, Jose Such^{6,7}, Tuğrulcan Elmas¹

¹University of Edinburgh
 ²Middle East Technical University
 ³Oxford Internet Institute
 ⁴IMDEA Networks Institute
 ⁵Universidad Carlos III de Madrid
 ⁶King's College London
 ⁷Universitat Politecnica de Valencia

yusufc@ceng.metu.edu.tr, vahid.ghafouri@oii.ox.ac.uk, guillermo.suarez-tangil@imdea.org, jose.such@kcl.ac.uk, guillermo.suarez-tangil@imdea.org, guillermo

telmas@ed.ac.uk

Abstract

Many social media studies argue that social media creates echo chambers where some users only interact with peers of the same political orientation. However, recent studies suggest that a substantial amount of Cross-Partisan Interactions (CPIs) do exist - even within echo chambers, but they may be toxic. There is no consensus about how such interactions occur and when they lead to healthy or toxic dialogue. In this paper, we study a comprehensive Twitter dataset that consists of 3 million tweets from 2020 related to the U.S. context to understand the dynamics behind CPIs. We investigate factors that are more associated with such interactions, including how users engage in CPIs, which topics are more contentious, and what are the stances associated with healthy interactions. We find that CPIs are significantly influenced by the nature of the topics being discussed, with politically charged events acting as strong catalysts. The political discourse and pre-established political views sway how users participate in CPIs, but the direction in which users go is nuanced. While Democrats engage in cross-partisan interactions slightly more frequently, these interactions often involve more negative and nonconstructive stances compared to their intra-party interactions. In contrast, Republicans tend to maintain a more consistent tone across interactions. Although users are more likely to engage in CPIs with popular accounts in general, this is less common among Republicans who often engage in CPIs with accounts with a low number of followers for personal matters. Our study has implications beyond Twitter as identifying topics with low toxicity and high CPI can help highlight potential opportunities for reducing polarization while topics with high toxicity and low CPI may action targeted interventions when moderating harm.

Introduction

Social media is often criticized for creating *echo chambers*, where individuals are primarily exposed to and interact with content that aligns with their beliefs (Barberá 2015). Such spaces reinforce existing views, limit exposure to diverse

*These authors contributed equally.

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

perspectives, and contribute to polarization. However, recent studies argue that there is a significant amount of *Cross-Partisan Interactions*, CPIs, engagements among users with different political views, challenging the notion that social media only deepens divides (Rajadesingan et al. 2021).

CPIs hold significance on social media as they serve as potential bridges between divided ideological groups, facilitating exposure to diverse perspectives and fostering dialogue among disparate communities. Meanwhile, there is a debate on whether such interactions lead to healthy dialogue and mitigate polarization, or are toxic and exacerbate it. Such contrasting opinions suggest that the social media dynamics that affect the presence and quality of CPIs are complex. Understanding these dynamics is crucial to understanding CPIs and subsequently, how information and differing viewpoints propagate through networks, potentially influencing societal polarization and discourse quality. This broadened exposure can lead to a more informed populace and potentially temper the polarization that is rampant in online spaces. Furthermore, understanding CPIs is crucial for platform designers and policymakers who aim to create healthier online environments.

Our main contribution is investigating the dynamics of CPIs by identifying the interplay between actors, topics, and stances that are more associated with CPIs. To do so, we conduct the first large-scale topic-independent quantitative analysis of Cross-Partisan Interactions to the best of our knowledge. We also study in which cases the CPIs are *healthy*, i.e., non-toxic. For this, we formulate the following Research Questions (RQs):

RQ1: What are the characteristics of **users** who typically engage with (un)healthy cross-partisan interactions?

RQ2: Which **topics** are more prevalent in healthy or unhealthy cross-partisan interactions?

RQ3: What type of **content**, in terms of sentiment and stance, appears in healthy cross-partisan interactions?

We collect a comprehensive reply dataset and devise a methodology to determine the political alignment of users, and analyze the stance of toxicity of the content to answer these questions. While doing so, we find that popular users are more likely to foster CPIs as they attract replies from users with opposite views. In contrast, the repliers in CPIs tend to have lower number of followers, followings and tweet counts. Certain topics such as politics, economics, and golf has higher rates of CPI in contrast to public health, religion, and daily life. CPIs are correlated with heightened toxicity, which suggests that while these interactions can bridge ideological gaps, they also risk escalating tensions. Furthermore, the study confirms that more polarized users tend to use offensive language. We discuss the findings in the related section.

Related Work

Cross-partisan interactions (CPIs) can potentially create bridges in polarized political environments and disrupt echo chambers. In juxtaposition, they may also exacerbate existing divides and increase toxicity online. We now provide a brief survey on these research themes.

Echo Chambers: Echo chambers are environments where people are exposed only to information and opinions that reinforce their existing beliefs, leading to a lack of exposure to diverse perspectives and the potential for increased polarization. Past studies argue for the existence of echo chambers online by employing user data (Garimella et al. 2018; Elmas et al. 2021; Torres-Lugo, Yang, and Menczer 2022). However, some studies suggest that their prevalence and impact may be overstated (Dubois and Blank 2018; Guess, Nyhan, and Reifler 2018), A substantial portion of the audience accesses news from outlets with varying ideological perspectives from both online and offline platforms (Gentzkow and Shapiro 2011). Such cross-partisan media consumption indicates that individuals are not wholly isolated within echo chambers but rather encounter diverse viewpoints.

Cross-Partisan Interactions: The impact of cross-partisan interactions on polarization and dialogue health is a subject of ongoing debate. Past studies have not reached a consensus on whether such interactions bridge divides or exacerbate tensions, e.g.: Bail et al. (2018) and Rajadesingan et al. (2021) suggest cross-partisan interaction increases polarization, while works like Rajadesingan et al. (2021) and Çetinkaya et al. (2024) highlight the potential for CPIs to foster understanding and reduce polarization.

Others show that the impact of cross-partisan interactions on polarization is more complex. Balietti et al. (2021) find that exposure to cross-partisan content can decrease polarization if the exposed feel close to the author. Marchal (2022) report that negative cross-partisan interactions increase political polarization, but positive ones mitigate it on Reddit. Zade et al. (2024) reveal that the reply feature is often used to re-frame and shift the attention to different dimensions of the debate that were not salient in the original tweet, which may or may not exacerbate tensions.

These studies suggest that the dynamics of cross-partisan interactions are complex, and their impact varies on a caseby-case basis. Our study adopts a comprehensive approach to understanding which actors and topics are more associated with cross-partisan interactions and on which occasions such cross-partisan interactions are more toxic. To the best of our knowledge, such a case-by-case understanding of the likeliness and impact of CPIs is not clearly established.

Perhaps the most similar work to ours, by An et al. (2024), that uses COVID-related Twitter data suggests that Republicans were more likely to initiate cross-party interactions, but such interactions were driven by a small set of users. We extend this work by employing a general dataset that stems from the 1% sample of Twitter instead of limiting it to Covid so that we can analyze which actors and topics are more likely to create cross-partisan interactions extensively. We also propose an unbounded analysis using LLMs to understand the prevalent stances in the data.

Shugars and Beauchamp (2019) employs a dataset of tweets related to Trump and proposes a model that predicts user participation in prolonged conversations. Our problem differs from theirs as they predict if a given user will reply to a conversation thread while we try to predict the factors that are more associated with cross-partisan interactions.

Finally, Santoro and Broockman (2022) find that the effectiveness of cross-partisan conversations depends on the topic through a randomized controlled trial. They employed two topics of discussion: a non-political topic discussing the perfect day, and the other discussing their favorite political party. The former mitigates the polarization, while the other has no effect. We complement this work by testing their hypothesis in practice and employing a unique, comprehensive social media analysis that is not limited to only two topics. **Text-mining on Online Interactions:** To analyze user interactions at scale, social media studies employ text-mining methods such as detecting the sentiment, stance, toxicity, or

discourse act of the content, employing classification using a predefined set of classes, or assigning a score.

For instance, to analyze the online debate on the Italian 2016 referendum Lai et al. (2019) use a semi-supervised approach by manually annotating a corpus of Italian tweet-replies as in "FAVOR", "AGAINST", or neutral ("NONE") to the root tweet, and training a SVM classifier on the labels using *bag-of-words*, *bag-of-hashtags*, and *bag-of-mentions* as features. The labels they annotate by such prediction indicate that users use replies more frequently to express divergent opinions than other communication types. Matalon et al. (2021) employ VADER (Hutto and Gilbert 2014) to evaluate the sentiment of tweets related to Israel to obtain a sentiment score in the range: [-1, +1]. The results indicate that a root tweet's sentiment influences the likelihood of receiving a quote that inverts its opinion.

While past studies mainly employ lexicon based (Hutto and Gilbert 2014) or transformer-based (Karande et al. 2021) approaches, recent studies suggest that LLMs are capable of text mining tasks such as *stance detection* (Lan et al. 2024), *sentiment analysis* (Lee et al. 2024), *hate speech* (Kaya, Ozcelik, and Toraman 2024), and sociopolitical affiliation of texts (Ghafouri et al. 2023).

In our work, we employ two text-mining methods: unbounded content annotation and toxicity detection. We employ the former to analyze the stance and sentiment of the interactions qualitatively. We opt for this method to go beyond predefined sentiment and stance classes and understand nuances in online interactions, which is novel. For toxicity detection, we use the transformers-based Perspective API (Lees et al. 2022) which provides a toxicity score for a given text between 0 and 1. Prior research identifies a threshold to classify if a given social media post is toxic or not. For instance, Kumar et al. (2021) survey 17,280 participants to rate toxic content. They then identify a personalized threshold for each participant. They find that the median optimal threshold across participants 0.61. Similarly, Aleksandric et al. (2024) reports that Perspective API achieves the highest accuracy of 0.95 at a threshold of 0.6 for the "Toxicity" attribute on a Twitter dataset they annotated.

Data & Methodology

We define a *cross-partisan interaction*, CPI for short, as a direct interaction between two users of different political orientations. We use Twitter (X) as the platform to study. Twitter features four types of interactions between users: retweets, likes, quotes, and replies. We limited our analysis only to replies, as only replies provide evidence that people intend to engage in a direct reciprocal interaction that can lead to a dialogue (An et al. 2024; Zade et al. 2024).

To study CPIs, we first collect a dataset of interactions in the form of replies, replied tweets, and root tweets. We then employ political orientation detection to identify partisanship and discover cross-partisan interactions.

We limit our focus to the U.S. context and define partisanship as left-aligned (leaning towards liberals or Democrats) or right-aligned (leaning towards conservatives or Republicans). We limit the data period to 2020 as it captures the general discussions, such as the pandemic (often non-political) and the political discussions related to the 2020 U.S. presidential election. The dataset and related labels used in the analysis are shared on GitHub.¹

Replies & Roots: To have an unbiased sample of replies, we employed the 1% random sample of Twitter provided by the Internet Archive (Archive 2024). We recollected the replies and the tweets they replied to in March 2023 using the Twitter API. The dataset comprises 3,029,231 reply tweets in English, responding to 2,299,444 unique tweets. However, on Twitter, tweets can be part of a reply chain. To simplify the analysis, we discard the nested replies and limit the analysis where the replied tweet is not a reply of the original tweet (namely, *root*). This brings the dataset to 1,870,794 direct replies (61.8% of all replies), replying to 1,175,151 root tweets. There are 683,450 unique repliers and 211,584 root authors. The dataset is further filtered by only keeping the users with non-neutral political orientation, using the methodology that we will describe now.

Political Orientation: We employ the methodology of Barberá (2015) to measure users' political orientation. The method uses Bayesian inference on users following data to assign a political orientation score to them, which ranges between -5 and +5. Negative values signify leaning toward Democrats and positive values mean leaning toward Republicans. Our dataset with only direct replies contains 821,791 users. Among these, 57,037 users (7.0%) are not assigned

a score due to the absence of the following data and are excluded from the analysis. We see 493,954 (60.1%) users classified as left-aligned and 232,186 (28.3%) classified as right-aligned due to having an absolute score above 0.1. There are 38,614 users with a score between -0.1 and 0.1, that are considered neutral and discarded from the analysis.

The direct replies and root tweets from users with unassigned political orientations and from neutral users are discarded from the dataset and are not used in the rest of the analysis. The final dataset comprises 2,561,846 tweets from 726,140 users, including 1,594,271 replies from 596,749 users and 967,575 root tweets from 142,219 users.

CPI Data: We classify an interaction as a CPI if the replier and the root author are assigned a different political orientation. There are 661,661 replies classified as CPI (34%). Of these, 196,642 are from Republicans replying to Democrats, making up 33.2% of all Republican tweets, and 432,004 are from Democrats replying to Republicans, accounting for 34.3% of all Democrat tweets. While cross-partisan interactions originating from Democrats constitute 65% of all CPIs, this is due to the dataset containing a larger number of users labeled as Democrats.

Toxic Language Classification: To determine whether the interactions are healthy or toxic, we employ Google Jigsaw's Perspective API. It is a public tool based on a transformer model that assigns a score which indicates the toxicity of the language within the text (Lees et al. 2022). The API outputs scores ranging from 0 to 1, indicating the presence of various forms of hate speech, including "Toxicity," "Severe Toxicity," "Insult," "Threat," "Profanity," and "Identity". We streamline the process for our analysis by assigning each tweet the maximum score across these categories. We define a threshold over the score that determines whether a tweet contains offensive language. We set the threshold at 0.61 which is identified by Kumar et al. (2021) by employing a survey on 17,280 participants and Aleksandric et al. (2024) who found that Perspective API gets the highest accuracy on this threshold on their independent test set.

Sentiment & Stance Annotation: To better understand user interactions, we characterize tweets with annotations describing their sentiment and stance. We use these annotations to perform a comparative content analysis. Due to the large size of our dataset, we resort to automated characterization mechanisms. In particular, we leverage state-of-theart LLMs. We prompt the root tweets and replies to an LLM and task the model to describe them using three adjectives. For the replies, the task is to describe the *stance* against the root tweet with three adjectives while we offer both the reply and the root tweet to the LLM. Since root tweets are not usually directed to another tweet, we ask the LLM also to annotate its *sentiment*. This method is an alternative to constraining LLMs by predefined classes and helps us qualitatively analyze the sentiment and stances of the tweets.

Table 2 (in Appendix) provides sample prompts passed to LLM and the completion provided by LLM. Later, we discuss the reason and limitations of our choice of prompt engineering and possible future configurations.

The interactions we characterize in our dataset are of four

¹https://github.com/tweetpie/cross-partisan-interactions.

categories: Democrats replying to Democrats $(D\rightarrow D)$, Republicans replying to Republicans $(R\rightarrow R)$, Republicans replying to Democrats $(R\rightarrow D)$, and Democrats replying to Republicans $(D\rightarrow R)$. As, for a fair comparison, we intend to have a balanced amount of annotations from each interaction type, we randomly sample 100,000 (the approximate size of the smallest category of interaction) tweet-reply pairs resulting in a total of 400,000 instances.

We employ "*Mistral-7B-Instruct-v0.2*"². This opensource model allows for efficient processing of the large dataset while maintaining adequate language-understanding capabilities. We downloaded the model locally from the Hugging-Face. Each query completion took 5 seconds on an NVIDIA A100 80GB GPU. Due to the light size of the model, we parallelized the process into 10 folds and obtain the answers in 10 days.

We validate the annotations by manually inspecting a random sample of 100 AI-annotated root tweet & reply pairs and labeling them as correct or incorrect. Two annotators independently performed this task The LLM (Mistral-7B-Instruct-v0.2) achieves an accuracy of 97% and 88% (as evaluated by annotator 1 and annotator 2, respectively) for determining the sentiments of root tweets, and 88% and 85% for identifying the stances of replies toward root tweets. The calculated Kappa inter-annotator agreement scores are 0.26 for the sentiments of root tweets, and 0.44 for the stances of replies toward root tweets, which correspond to fair agreement and moderate agreement respectively (Landis and Koch 1977). The former score is lower due to high prevalence of data points annotated as true. Although not classified as incorrect, we observe that LLMs sometimes annotate the sentiment of the reply instead of its stance, e.g., "happy" instead of "happy for them" or "supportive." We do not correct these annotations and leave them as a limitation.

RQ1: User Analysis

This section investigates the characteristics of users who engage in CPIs, addressing RQ1. We first focus on profile attributes to identify the individuals most commonly associated with CPIs. We then adopt a user-specific approach, analyzing each user. We use regression to understand the correlation across features.

Profile & Tweet Attributes vs. CPI Rates

Certain profile and tweet attributes may be more likely to foster CPIs. Thus, we investigate the prevalence of certain user-based attributes in the promotion of CPIs. In particular, we look at the following profile and tweet attributes: (1) number of followers, (2) followings, (2) total tweet count, and (3) tweets' replies, (4) retweets, (5) quote, and (6) like count. We also incorporate the same attributes (1-3) to understand if they have a larger influence than the root user.

Fundamentals of the Analysis: We use Logistic Regression to compute the coefficients for these attributes to identify their effect, using 20% of the dataset as the test set to measure the models' performances. We train three models: one

using the entire dataset, one focusing solely on replies from Democrats, and another focusing on replies from Republicans. The dependent variable Y is whether a tweet is part of a cross-partisan conversation (Y = 1) or not (Y = 0). The logistic regression model can be expressed as follows:

$$logit(P(Y=1)) = log\left(\frac{P(Y=1)}{1 - P(Y=1)}\right) = \beta_0 + \sum_{i=1}^n \beta_i X_i,$$
(1)

where:

- β_0 is the intercept of the model.
- β_i are the coefficients associated with the *n* predictor variables X_i .
- X_i represent the various author-related metrics included in the model.

The logistic regression model for the entire dataset is estimated using Maximum Likelihood Estimation (MLE) on a dataset of 1,902,704 observations. Figure 1 displays the distribution of attributes used in the model. The model converges after 6 iterations with a Pseudo R-squared of 0.05413. This indicates that the model explains about 5.4% of the variance with the given predictors in the outcome variable, which is typical for models in complex social or behavioral datasets where many factors are at play. The model's ROC AUC score is 0.656, which suggests a moderate discriminative ability. The model that uses Democrats subset achieves a significantly higher Pseudo R-squared value of 0.1598, indicating that the predictors explain about 15.98% of the variance, compared to only 3.48% in the Republican model. Similarly, the ROC AUC score for Democrats is 0.6357, reflecting moderate discriminative ability, while the Republican model's ROC AUC score is only 0.538. Table 1 presents the coefficients for each predictor variable for three models. All coefficients are statistically significant (p < 0.0001).

Variable	All	Dem.	Rep.
const	-0.656	-0.482	-0.672
replier followers_count	-0.025	<u>-0.158</u>	0.020
replier following_count	-0.058	-0.077	-0.061
replier _{tweet-count}	-0.065	-0.082	-0.028
root_userfollowers_count	0.284	0.817	-0.356
root_userfollowing_count	-0.020	<u>0.032</u>	-0.119
root_user_tweet_count	-0.074	-0.092	-0.019
root _{reply-count}	0.525	<u>1.978</u>	<u>-0.304</u>
rootretweet_count	0.276	0.719	0.220
root _{quote-count}	-0.199	-0.378	-0.164
root _{like_count}	-0.292	<u>-1.450</u>	<u>0.788</u>

Table 1: Logistic Regression results for CPIs. Coefficients are underlined if the sign is different in two camps.

Key Predictors: We observe that accounts fostering CPIs by replying have fewer count of followers, following, and tweets as these have negative coefficients. This means that new or small accounts are more likely to engage in CPIs, possibly to attract attention and grow their accounts. Alternatively, people may be using their fake accounts to engage

²https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.2

in CPIs. The differences among the two type of repliers are rather small as in Figure 1. On the other hand, root tweets from users with high follower count are more associated with CPIs from Democrats but less with CPIs from Republicans. Root users' tweet counts have negative association while their following counts have contrasting but negligible associations with CPIs. We observe that the coefficients for the root's attributes are higher than the repliers' attributes. This suggests that the CPIs depend mostly on root users.

Interestingly, the attributes of root tweets have contrasting effects. The reply count of the root tweet has the biggest coefficient by far for the entire dataset. This suggests a richget-richer effect, i.e., popular tweets with many replies also attract many CPIs. This findings is even more pronounced for Democrats, but does not generalize to Republicans, which has a negative coefficient for the same attribute. The tweet's retweet count has a similar association and generalizes to Republicans although the coefficient is low. Quote count is less associated with CPIs for both camps, possibly because users tend to quote tweets without replying to them, thereby reducing CPIs in the form of replies.

Surprisingly, like count has a very high negative coefficient for Democrats and a very high positive coefficient for Republicans. This indicates that Republicans are more inclined to engage in CPIs with highly liked tweets from the opposing camp, whereas Democrats may tend to create discussions on tweets they (and the general public) dislike. Alternatively, the likes may amplify the content in the users' echo chamber, while retweets may go beyond and create cross-partisan interactions. Note that since the data is from 2020, the disparity between these attributes may have changed due to changes in Twitter's algorithms.

The differences between the camps underline how crosspartisan engagement is influenced not only by root tweet attributes but also by the alignment of the repliers, reflecting distinct behavioral tendencies and engagement strategies.

We would like to clarify that the attributes of root tweets with negative coefficients (e.g., quote and like count for the entire dataset) fostering CPIs are still higher as seen in Figure 1, since they are correlated with other attributes, e.g., a tweet with many likes and quotes also receive many quotes and replies. However, these attributes contribute negatively to the model, i.e., likes and quotes may be lower for the tweets with comparable retweet and reply counts in CPIs.

User-Specific CPI Rates

While user profile attributes partly explain the type of users that are more associated with CPIs, it does not provide the whole picture. For instance, popular users may differ in their likelihood of initiating CPIs. For a more in-depth analysis, we focus on identifying the key predictors among users.

Fundamentals of the Analysis: We employ Lasso regression to model these users. This method is particularly suited to our dataset, which is characterized by high dimensionality and sparsity, i.e., there are many users but a limited number of interactions among them, which is typical of social media (Tibshirani 1996; Emmert-Streib and Dehmer 2019).

We prepare a sparse matrix where rows represent interactions, columns represent users, and each entry is a binary in-



Figure 1: The distribution of the attribute values among the CPIs and non-CPIs.

dicator denoting whether the interaction is associated with a specific user. The output variable is binary, with a value of 1 indicating the conversation is CPI and -1 indicating non-CPI. We apply a threshold to eliminate noise from users with only a few interactions, considering only those users who acted as a replier or a root more than 10 times. The filtered dataset has 1,702,019 interactions created by 72,923 users. 24,729 of these users were exclusively repliers, 16,407 users were exclusively roots, and 31,787 users interacted both ways.

The Lasso regression model is defined as follows:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{N} \sum_{i=1}^{N} (y_i - \mathbf{x}_i^T \beta)^2 + \lambda \sum_{j=1}^{p} |\beta_j| \right\}, \quad (2)$$

where y_i is the response variable, \mathbf{x}_i is the vector of predictors for the *i*-th observation, β is the vector of coefficients, and λ is the regularization parameter that controls the degree of shrinkage applied to the coefficients (Friedman, Hastie, and Tibshirani 2010). Users with non-zero coefficients are deemed significant contributors to the CPIs. This feature selection process is intrinsic to the Lasso technique and aids in simplifying the model by removing irrelevant predictors.

Key Predictors: Figure 2 shows the political orientation scores of the authors with their coefficient signatures. There are 26,581 users with negative, 26,649 users with positive, and 19,693 users with zero coefficients.



Figure 2: The distribution of political orientation scores for Lasso coefficients of authors on their CPI involvement.

The coefficients of the users show a Pearson correlation of 0.78 with their CPI versus non-CPI ratios, indicating a strong positive relationship. This suggests that as the coefficient values increase, the likelihood of a user contributing to CPIs also rises proportionally. Our analysis reveals that Democrats tend not to be involved in CPIs, as their bins have more negative coefficient portions compared to Republicans, who are more likely to engage in CPIs. Moreover, the political orientation scores show an inverse relationship with the probability of creating CPIs; users close to the political center have a higher percentage of CPI involvement, which can be interpreted as being in more moderate positions. In contrast, those with more extreme scores on the political spectrum are less likely to be involved in CPIs. We further document this in Figure 3, which shows the distribution of CPIs with respect to political orientation scores.



Figure 3: The distribution of CPIs by authors' political orientation scores, grouped into 0.5 interval bins.

CPI Toxicity: We categorize the replies into bins based on the political orientation of the replier. For each bin, we calculate the average toxicity score of the replies. Additionally, we determine the percentage of toxic tweets in each bin for interpretability, using a toxicity threshold of 0.61. We removed the bins where the replier has a political orientation score greater than 3.5 as there were only 205 replies in total. We repeat this analysis by averaging the reply toxicity per user, so that each data point in a bin represents a user rather than a tweet. Since the results do not differ significantly, we visualize only the first experiment. Furthermore, we tested the sensitivity of our findings by varying the toxicity threshold from 0.4 to 0.9. This adjustment yielded no substantial change in the behavior of the toxicity distribution. Figure 4 illustrates the mean reply toxicity scores (bars) and percentage of toxic tweets (bar labels) across different political orientation score bins. The overall trend shows that toxicity levels for CPI replies (red bars) are consistently higher than for non-CPI replies (blue bars) as the political orientation score moves away from the center. As moving toward the center, toxicity scores gradually decrease for both CPI and non-CPI interactions, reaching their lowest around the center, where CPI toxicity drops to 3.1% and non-CPI toxicity to 2.7%. Interestingly, toxicity decreases at the edges of the spectrum. Overall, the toxicity scores have a symmetric behavior with respect to the point zero for both parties.



Figure 4: Mean toxicity levels of the repliers' political orientation score bins for both CPI and non-CPIs. The numbers on bars show the percentage of the content that is counted as toxic (*score* > 0.6) in each bin.

RQ2: Topic Analysis

To answer RQ2, we investigate the characteristics of the topics that are likely to foster CPIs. We employ BERTopic (Grootendorst 2022) pipeline for topic modeling to identify the topics. We first generate dense vector representations of the texts using BERT, a pre-trained transformer model. We then reduce the dimensionality of these embeddings using UMAP to make the clustering more efficient. Finally, we apply HDBSCAN, a clustering algorithm, to group similar embeddings, forming distinct topics. We process "root tweets" to generate embeddings for further analysis.

The topic modeling identifies 177 distinct topics within the dataset. We further categorize these topics to ease comprehension and visualization. We utilize the top 10 keywords provided by BERTopic, along with the 10 closest tweets to the centroid of each topic, to create headline summaries and human-readable titles for each topic using ChatGPT-40. We manually group the topics into 11 overarching *cat*egories using these summaries to simplify the analysis. Political topics comprise approximately 35% of these 177 topics, encompassing discussions on governance, civil rights, election campaigns, and policy debates. A notable subset of these, around 10-15%, specifically pertains to electionrelated discourse, reflecting the temporal and sociopolitical context of the data. Non-political topics, by contrast, encapsulate a broad spectrum of interests, including cultural reflections, entertainment, and everyday experiences, such as musical themes, food discussions, and personal milestones.

Topics vs. CPI Rates

We compute the percentage of replies to root tweets in a specific category relative to all categories. We also compute the percentage of CPIs in each category. Figure 5 presents the results. The replies to tweets with the topics under *Hobbies* & *Daily Life*, and *Politics & Governance* categories emerge as the most substantial, making up nearly half of the data for both camps. In contrast, *Law, Order & Justice, Sport Events Environmental Issues*, and *Religious & Praying* occupy smaller portions of the dataset as they are rather niche. The popularity of categories remains consistent across both groups, except for a notable difference: *Hobbies & Daily Life* is the most popular among Democrats, whereas *Politics & Governance* takes the lead among Republicans.



Figure 5: The distribution of topic categories in the dataset with the box-and-whisker plot of their CPI percentages on different topics for the entire dataset and both camps.

The topics under *Politics & Governance* and *Politic Figures* generate the highest levels of CPI, suggesting that they are particularly effective at breaking through echo chambers. Within these politically charged categories, there is considerable variance in CPI percentages, reflecting the context-dependent nature of how these topics are discussed and debated. For instance, topics such as "Impeachment Proceedings and Senate Actions" and "Election Strategies" exhibit a high degree of CPI (53, 60%), highlighting their contentious nature and the widespread debate across ideological divides.

The topics of public interest such as "Political Approval", "Sporting Events", and "Economic Stimulus" exhibit significant cross-partisan engagement. These topics span various domains, including governance, sports, and economic issues, attracting users from diverse ideological backgrounds due to their broad relevance. Discussions related to *Sport Events* like "Sporting Events and National Interests" show notable CPI with percentages above 55%. In contrast, more routine discussions in *Public Health*, such as "Health Guidelines and Mask Usage" and "Hygiene Practices", display lower CPI, with percentages between 27% and 36%.

A contrast between camps is observed in the "Hobbies & Daily Life" and "Politics & Governance." While both groups exhibit a significant proportion of discussions in these categories, Republicans demonstrate a higher CPI percentage in the "Hobbies & Daily Life." In contrast, Democrats show a slightly higher CPI percentage in the "Politics & Governance." The category "Sport Events" stands out due to its consistently high CPI percentages for both groups.

Conversely, topics such as "Religious & Praying" and "Environmental Issues" are less represented in the dataset, exhibiting minimal proportions for both groups. However, Democrats show a slightly higher CPI percentage in "Religious & Praying," suggesting a marginally greater propensity for cross-partisan discourse on this topic. Meanwhile, discussions on "International Affairs" appear consistently across both groups but differ in CPI percentages. Democrats display a higher median CPI percentage and a broader range of cross-partisan engagement in this category.

Topics under the categories *Hobbies & Daily Life* and *Religious & Praying* exhibit considerably lower CPI percentages, indicating that they are more likely to be confined within echo chambers. For instance, topics such as "*Musical Themes*", "*Reading Habits*", "*Popular Food Discussions*", "*Community Support and Prayers*", "*Religious Cultural Commentary*", and "*Religious Beliefs*" show minimal cross-partisan engagement, with CPI percentages between 23% and 30%, which suggests that users participating in these discussions are more insular, predominantly interacting with like-minded individuals. However, there are exceptions, particularly within the sports category, where discussions around Golf and Football topics show higher CPI percentages, ranging from 51% to 55%.

The findings suggest that while politically and socially significant topics tend to disrupt echo chamber dynamics, more personal or neutral topics may reinforce them, leading to isolated discussions within homogeneous groups. While echo chambers are indeed a concern, particularly around less contentious topics, the data reveals that not all discussions are confined to such insular spaces. Political and controversial subjects seem to break through these barriers, introducing users to diverse perspectives, a beneficial trend for democratic processes, and informed civic engagement.

Topics vs. Stance & Toxicity

The variance in CPI and toxicity percentages across different topics, as illustrated in Figure 6, offers valuable insights into the dynamics of online discourse. Generally, CPI percentages reflect a topic's political or social sensitivity, with those rooted in societal values and governance tending to exhibit higher CPI rates, sparking greater cross-partisan discussions. A Pearson correlation of 0.76 between topic CPI and toxic content percentages emphasizes the strong positive relationship between these two dimensions.

Topics in the chart's top-left quadrant, such as "Social Labels and Cultural Stereotypes," "Gender Discussions," and "Racial Discrimination," represent highly toxic, low-CPI subjects. These are problematic as they often fuel internal group animosity and deepen polarization. In contrast, topics in the bottom-right quadrant, such as "School Operations and Student Concerns," "Financial Aid and Social Welfare," "NFL Draft Strategies," and "Health Insurance Policies," are characterized by high-CPI and lower toxicity, which offer opportunities for bridging societal divides, acting as contact points where cross-partisan dialogue can flourish.

Figure 7 presents the top five stance percentages used in CPIs across major topics, offering insights into the emotional and attitudinal tone of the replies. The distribution of stances varies significantly depending on the topic, highlighting the diversity of cross-partisan communications. For example, in non-political domains like Hobbies & Daily Life and Sport Events, positive and playful stances such as "eager," "happy," and "exciting" are prominent, suggesting that these areas foster more open and lighthearted dialogue. Conversely, topics related to Politics & Governance and Public *Health* are dominated by more serious and critical stances, including "concerned," "confident," and "cautious," reflecting these discussions' contentious and urgent nature. "Religious & Praying" demonstrates a higher prevalence of empathetic and introspective stances, indicating a more emotionally driven discourse. The figure also reveals that certain topics, such as Economic Issues and International Affairs, exhibit a blend of concerned and curious stances, pointing to an underlying tension between critique and a willingness to understand the opposing view. The variation in stance distribution across topics suggests that the emotional tenor of CPIs is context-dependent, with some domains providing opportunities for constructive engagement while others remain more adversarial or cautious. Understanding these stance patterns is critical for fostering healthier cross-partisan dialogue in polarized environments.

RQ3: Content Analysis

Lastly, we investigate the content characteristics that are more associated with CPIs as part of our efforts to address RQ3. For this, we leverage the LLM annotations (root sentiments and reply stances) generated.

Stance Contrast, PI vs. CPI

Our initial phase of content analysis is to investigate how users of different parties differentiate in terms of stance when interacting with in-group (PI) and out-group (CPI) users. This can potentially shed light on the productivity of CPIs across the two parties.

We compare the frequencies of AI-annotated stances in partisan (PI) vs. cross-partisan (CPI) replies (e.g., $D \rightarrow D$ vs. $D \rightarrow R$, $R \rightarrow R$ vs. $R \rightarrow D$), normalized by their mean. Each annotation's contrast is shown in a pyramid bar chart. Let $f_i^{G_{right}}$ be the frequency of annotation *i* in the right-side group; the x-axis is computed as in Equation 3.

$$X_{i} = \frac{f_{i}^{G_{right}} - f_{i}^{G_{left}}}{(f_{i}^{G_{right}} + f_{i}^{G_{left}}) \times 0.5}$$
(3)

A positive (negative) value indicates a higher frequency in the right (left) group.

Figure 8a presents the word frequency differences between Democrat-to-Democrat $(D\rightarrow D)$ and Democrat-to-Republican $(D\rightarrow R)$ interactions. Notably, $D\rightarrow D$ interactions are characterized by words associated with empathy, positive sentiment, and agreement, such as "*happy*", "*reciprocal*", and "*empathetic*". Conversely, $D\rightarrow R$ interactions exhibit words indicative of conflict, like "*accusatory*", "*dismissive*", "*hostile*", and "*critical*".

Unlike Democrats, we do not observe a strong contrast between $R \rightarrow R$ and $R \rightarrow D$ interactions (Figure 8b). Except for "agreement" vs. "unsupportive", Republican replies show more consistent tone across partisan and cross-partisan contexts. Some positive stances, like "loving" and "congratu*latory*", appear more often in $R \rightarrow D$ replies, often in personal contexts such as births or marriages (Table 3). This aligns with our earlier finding that Republicans tend to reply to less-followed users (Table 1), suggesting a lower distinction between in- and out-group engagement, particularly in personal matters. Furthermore, the stance contrast bars in Figure 8a are visibly larger than in Figure 8b, indicating Democrats vary their tone more between PIs and CPIs. While Democrats are more likely to engage in CPIs (Barberá 2015), their cross-partisan replies are also more likely to be negative or confrontational compared to their intraparty ones-a pattern less evident among Republicans.

Root Sentiment vs. Reply Stance & Toxicity

Our next content-based analysis focuses on identifying the underlying sentiments that are more likely to foster certain stances in replies. We employ the top 10 most frequent adjectives from the root tweets' sentiments and stances of the replies to create contingency matrices for both PI and CPI. We use the Chi-Squared test $\chi^2 = \sum \frac{(\text{Observed}-\text{Expected})^2}{\text{Expected}}$ to assess associations between sentiment \leftrightarrow stance pairs. Figure 9 shows a triangular heatmap: top-left triangles represent partisan, bottom-right cross-partisan interactions. Colors range from blue (-1) to red (+1), with white indicating no association. Each cell shows $\frac{\text{Observed}-\text{Expected}}{\text{Expected}}$ and the corresponding p-value, computed using the squared numerator from the test.



Figure 6: The distribution of the topics with axes representing the CPI percentage and the toxicity percentage.



Figure 7: Top 5 reply stance annotation percentages under each topic for CPIs.

The heatmap reveals that certain root sentiments are more likely to encourage similar sentiments/stances in the replies. For instance, we observe that "*eager*" sentiment in the root tweets often lead to "*eager*" and "*motivating*" stances in replies, with values highly above 0, indicating a strong positive correlation. Similarly, "*Positive*" sentiments generally

elicit more "motivating" responses, while "critical" sentiments lead to stances like "dismissive", "critical", and "sarcastic", all with values greater than 0. In contrast, "skeptical" and "critical" sentiments discourages "eager" and "motivating" stances in replies, with values strongly below 0.

We also observe that these patterns are relatively consistent across both partisan and cross-partisan interactions in the heatmap. To further validate this, we extend the scope of the contingency matrix to the 50 most frequent sentiments and stances (a contingency matrix with 2500×2 cells) and detect a Pearson correlation of 95% between the values calculated for PIs and the values calculated for the CPIs. This may suggest that root sentiment is not a very strong factor for fostering CPIs compared to users and topics.

In addition to stances, we also investigate the toxicity rates provoked in the replies by each of the root sentiments. Figure 10 illustrates the top 10 root sentiments that attract the highest and lowest toxicity score replies in cross-partisan interactions. The mean toxicity score quantifies the level of hostility or aggressive language in replies for the given root sentiments that provoke these responses. The data shows that negative sentiments such as "juvenile," "shameful," "hateful," and "disgusted" are associated with the highest mean toxicity scores, ranging between 0.18 and 0.25. These sentiments, often judgmental or derogatory, may incite defensive or antagonistic replies, particularly in cross-partisan settings where opposing ideological stances intensify interpersonal conflict. In contrast, more positive or neutral sentiments like "loving," "gentle," "happy," and "friendly" correspond to significantly lower toxicity scores, ranging from



Figure 8: Stance-wise differences of partisan **vs.** crosspartisan replies across parties. The bar labels indicate the overall frequency of the annotation. For samples of tweets for each annotation (stance), see Table 3 in Appendix.

0.04 to 0.08. These sentiments indicate harmonious or empathetic discourse, fostering more civil and respectful replies. The stark contrast between the sentiments that provoke the most and least toxic responses suggests that the nature of root sentiments plays a pivotal role in shaping the tone of online discussions. In CPIs, where political and ideological tensions are inherent, sentiment-driven language influences the trajectory of dialogue, with more aggressive language fueling hostility, and gentler or more neutral language encouraging more constructive exchanges. Therefore, strategically using language emphasizing empathy or neutrality may mitigate the toxicity in polarized conversations.

Discussion

We discuss the takeaways of our analysis and its limitations.

Takeaways

Our work provides valuable new insights into online CPIs from a multi-fold perspective, i.e.: user, topic, and content. In particular, we find common characteristics of users engaging in CPIs, analyzing the discourse in terms of topics, and content; where the content breaks down to the analysis of *sentiments* of the roots, *stances* of the replies, and *toxicity*.



Figure 9: Chi-test statistics $\chi = \frac{Observed-Expected}{Expected}$ for co-occurrences of sentiments in root tweets and stances in replies (Top-Left: PI, Bottom-Right: CPI). Starred cells indicate p-values below 0.05.



Figure 10: Top 10 root sentiments that attract the highest and lowest toxicity score replies.

When we put these three dimensions in context, we see how social media exchanges influence societal polarization and discourse quality from a complementary perspective when compared to prior work. In particular, we derive the following novel insights and discuss the implications they have.

The Who: Those who participate in CPI show *distinct patterns based on their political orientation*. Users with extreme political views, especially on the far-right, tend to use more toxic language. Interestingly, the center-right replies with less *toxicity* than the center-left, yet the *extreme toxicity* category is still higher in Republicans. Although Democrats are proven to be more likely to engage in CPI than Republicans (Barberá 2015), their participation often includes more negative and non-constructive stances, such as being "*critical*", "*accusatory*", or "*hostile*". In contrast, Republicans show less differentiation in stance between interactions within their own group (PI) and cross-partisan replies (CPI), with instances of even positive CPI stances ("*loving*" and "*congratulatory*"), suggesting a more consistent approach to discourse regardless of political affiliation.

Moreover, we detect an association between the CPI rate of a user and their popularity: *Popular users receive more attention, leading to higher CPIs, while smaller accounts mainly engage with the opposite group.* One potential explanation for this could be that since the less followed accounts are viewed less, they may feel more secure to engage in a controversial debate with the opposite party. Alternatively, they may be fake or new & small accounts craving for increasing their impressions and followers.

The What: CPIs are significantly influenced by the nature of the topics being discussed, with politically charged events acting as strong catalysts. However, not all topics foster constructive dialogue. Topics such as gender, race, and discrimination tend to foster less CPI but higher toxicity, indicating more adversarial and polarized exchanges. Conversely, discussions around financial aid and social welfare demonstrate a healthier balance, with more CPI and lower toxicity, suggesting these areas offer opportunities for more constructive and cooperative cross-party engagement. Identifying topics with low toxicity and high CPI can help highlight potential opportunities for reducing polarization, while high-toxicity, low-CPI topics represent areas of concern. Our results show the distribution of CPI and toxicity across different discussion themes, highlighting potential opportunities for reducing polarization and identifying areas of concern.

The How: To foster a healthy online dialogue, how you talk matters. Users who talk in a friendly way, receive less toxic or critical replies. For instance, in Figure 9, "eager" and "positive" sentiments in the root tweets associate with "motivating" and "eager" stances in the replies and correlate negatively with "dismissive" stance. On the other hand, negative sentiments in root tweets are more likely to provoke toxic (see Figure 10) or negative stance in the replies (e.g. "dismissive" or "critical"). Interestingly, this phenomenon is independent of whether the interaction is a CPI or a PI as we observed a 95% correlation for the values in Figure 9. Understanding this dynamic is paramount to reducing the widespread polarization in online spaces and disrupting the idea that toxicity is inherent to certain online communities.

Limitations

Continuous CPIs: The CPIs between users close to the political center, and the CPIs between extremes may be inherently different. We initially experimented with a continuous CPI value by multiplying the political orientation score of the poster and the replier. However, we did not notice a drastic change in the results of our initial experiments. We plan to address a deeper analysis as part of our future work.

Reply chain and media: To simplify the analysis, we limit it to root tweets and their direct replies, and discard the replies of replies. Our content analysis is limited to the textual content present in the tweets. The analysis of the links, images, and videos in the tweets falls beyond our scope.

LLM's annotation: We used a heuristic approach with Mistral-7B-Instruct-v0.2 for stance/sentiment annotation. While effective overall, some replies were misclassified due

to overlap between stance and sentiment. Without setting a limit for the number of adjectives, the LLM generated sentences rather than words as in Table 2. Moreover, in some cases, the LLM did not differentiate the concept of *stance* and *sentiment* when asked to annotate the *stances* of the replies (e.g. the word "loving" in Figure 8 is more of a *sentiment* than a *stance*). However, we find such confusions to be rare, and we observe that they do not hinder the objective of our analysis. This is because understanding any salient semantic quality in the content of replies is insightful for us; whether it is *stance* or *sentiment*. Experiments with other LLMs and configurations are left to the future work.

Platform: Our study focuses on X (formerly Twitter), which has its own unique dynamics. These dynamics may influence the interaction patterns and content characteristics observed in our analysis. Expanding this research to other social media platforms, such as Reddit, YouTube, Instagram, and others, would provide a broader understanding of crossplatform user behavior and interactions. Such platforms offer different interaction paradigms, media-sharing mechanisms, and community structures that could yield additional insights and validate the generalizability of our findings.

Context: The data collection period was chosen to be 2020 as this year has seen massive debates both on political (riots, elections) and non-political debates (the pandemic). However, the presence of a Republican government in the U.S. may have influenced the results, which could differ under a Democrat government. Additionally, the results may have differed during a period with less politically charged events, such as in 2021. Extending the study over a longer period could help capture evolving patterns. We acknowledge these limitations and leave extending the study to a longer period for future work due to current restrictions with collecting Twitter data. Finally, as the dataset was limited to the U.S. context, the findings may not be a global context.

Acknowledgments

Guillermo Suarez-Tangil has been appointed as 2019 Ramon y Cajal fellow (RYC-2020-029401-I) funded by MICIU/AEI/10.13039/501100011033 and ESF Investing in your Future.

References

Aleksandric, A.; Roy, S. S.; Pankaj, H.; Wilson, G. M.; and Nilizadeh, S. 2024. Users' Behavioral and Emotional Response to Toxicity in Twitter Conversations. In *ICWSM*, volume 18, 29–42.

An, Z.; Breuhaus, J.; Niu, J.; Sariyuce, A. E.; and Joseph, K. 2024. Curated and Asymmetric Exposure: A Case Study of Partisan Talk during COVID on Twitter. In *ICWSM*.

Archive, T. I. 2024. The twitter stream grab.

Bail, C. A.; Argyle, L. P.; Brown, T. W.; Bumpus, J. P.; Chen, H.; Hunzaker, M. F.; Lee, J.; Mann, M.; Merhout, F.; and Volfovsky, A. 2018. Exposure to opposing views on social media can increase political polarization. *National Academy* of Sciences, 115(37).

Balietti, S.; Getoor, L.; Goldstein, D. G.; and Watts, D. J. 2021. Reducing opinion polarization: Effects of exposure to similar people with differing political views. *National Academy of Sciences*, 118(52).

Barberá, P. 2015. Birds of the same feather tweet together: Bayesian ideal point estimation using Twitter data. *Political analysis*, 23(1).

Çetinkaya, Y. M.; Lee, Y.; Külah, E.; Toroslu, İ. H.; Cowan, M. A.; and Davulcu, H. 2024. Toward a Programmable Humanizing Artificial Intelligence Through Scalable Stance-Directed Architecture . *IEEE Internet Computing*, 28(05): 20–27.

Dubois, E.; and Blank, G. 2018. The echo chamber is overstated: the moderating effect of political interest and diverse media. *Information, communication & society*, 21(5).

Elmas, T.; Hardi, K.; Overdorf, R.; and Aberer, K. 2021. Can Celebrities Burst Your Bubble? In *MISINFO*.

Emmert-Streib, F.; and Dehmer, M. 2019. High-dimensional LASSO-based computational regression models: regularization, shrinkage, and selection. *Machine Learning and Knowledge Extraction*, 1(1).

FORCE11. 2020. The FAIR Data principles. https://force11. org/info/the-fair-data-principles/.

Friedman, J.; Hastie, T.; and Tibshirani, R. 2010. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1).

Garimella, K.; De Francisci Morales, G.; Gionis, A.; and Mathioudakis, M. 2018. Political discourse on social media: Echo chambers, gatekeepers, and the price of bipartisanship. In *WWW*.

Gebru, T.; Morgenstern, J.; Vecchione, B.; Vaughan, J. W.; Wallach, H.; Iii, H. D.; and Crawford, K. 2021. Datasheets for datasets. *Communications of the ACM*, 64(12): 86–92.

Gentzkow, M.; and Shapiro, J. M. 2011. Ideological segregation online and offline. *Journal of Economics*.

Ghafouri, V.; Agarwal, V.; Zhang, Y.; Sastry, N.; Such, J.; and Suarez-Tangil, G. 2023. AI in the Gray: Exploring Moderation Policies in Dialogic Large Language Models vs. Human Answers in Controversial Topics. In *CIKM*.

Grootendorst, M. 2022. BERTopic: Neural topic modeling with a class-based TF-IDF procedure. *arXiv*:2203.05794.

Guess, A.; Nyhan, B.; and Reifler, J. 2018. Selective exposure to misinformation: Evidence from the consumption of fake news during the 2016 US presidential campaign. *European Research Council*, 9(3).

Hutto, C.; and Gilbert, E. 2014. Vader: A parsimonious rulebased model for sentiment analysis of social media text. In *ICWSM*.

Karande, H.; Walambe, R.; Benjamin, V.; Kotecha, K.; and Raghu, T. 2021. Stance detection with BERT embeddings for credibility analysis of information on social media. *PeerJ Computer Science*, 7.

Kaya, A.; Ozcelik, O.; and Toraman, C. 2024. ARC-NLP at ClimateActivism 2024: Stance and Hate Speech Detection by Generative and Encoder Models Optimized with Tweet-Specific Elements. In *CASE 2024*, 111–117.

Kumar, D.; Kelley, P. G.; Consolvo, S.; Mason, J.; Bursztein, E.; Durumeric, Z.; Thomas, K.; and Bailey, M. 2021. Designing toxic content classification for a diversity of perspectives. In *SOUPS '21*.

Lai, M.; Tambuscio, M.; Patti, V.; Ruffo, G.; and Rosso, P. 2019. Stance polarity in political debates: A diachronic perspective of network homophily and conversations on Twitter. *Data & Knowledge Engineering*, 124.

Lan, X.; Gao, C.; Jin, D.; and Li, Y. 2024. Stance Detection with Collaborative Role-Infused LLM-Based Agents. *ICWSM*, 18(1).

Landis, J. R.; and Koch, G. G. 1977. The Measurement of Observer Agreement for Categorical Data. *Biometrics*, 33(1): 159–174.

Lee, Y.; Çetinkaya, Y. M.; Külah, E.; Toroslu, İ. H.; and Davulcu, H. 2024. Masking the Bias: From Echo Chambers to Large Scale Aspect-Based Sentiment Analysis. In *ASONAM*. Springer.

Lees, A.; Tran, V. Q.; Tay, Y.; Sorensen, J.; Gupta, J.; Metzler, D.; and Vasserman, L. 2022. A new generation of perspective api: Efficient multilingual character-level transformers. In *ACM SIGKDD*.

Marchal, N. 2022. "Be nice or leave me alone": An intergroup perspective on affective polarization in online political discussions. *Communication Research*, 49(3).

Matalon, Y.; Magdaci, O.; Almozlino, A.; and Yamin, D. 2021. Using sentiment analysis to predict opinion inversion in Tweets of political communication. *Scientific Reports*.

Rajadesingan, A.; Duran, C.; Resnick, P.; and Budak, C. 2021. 'Walking Into a Fire Hoping You Don't Catch': Strategies and Designs to Facilitate Cross-Partisan Online Discussions. *ACMHCI*.

Santoro, E.; and Broockman, D. E. 2022. The promise and pitfalls of cross-partisan conversations for reducing affective polarization: Evidence from randomized experiments. *Science advances*, 8(25).

Shugars, S.; and Beauchamp, N. 2019. Why keep arguing? Predicting engagement in political conversations online. *Sage Open*.

Tibshirani, R. 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1).

Torres-Lugo, C.; Yang, K.-C.; and Menczer, F. 2022. The manufacture of partisan echo chambers by follow train abuse on Twitter. In *ICWSM*.

Zade, H.; Williams, S.; Tran, T. T.; Smith, C.; Venkatagiri, S.; Hsieh, G.; and Starbird, K. 2024. To reply or to quote: Comparing conversational framing strategies on Twitter. *Computing and Sustainable Societies*.

Ethics Checklist

- 1. For most authors...
- (a) Would answering this research question advance science without violating social contracts, such as violating privacy norms, perpetuating unfair profiling, exacerbating the socio-economic divide, or implying disrespect to societies or cultures? Yes
- (b) Do your main claims in the abstract and introduction accurately reflect the paper's contributions and scope? Yes
- (c) Do you clarify how the proposed methodological approach is appropriate for the claims made? Yes
- (d) Do you clarify what are possible artifacts in the data used, given population-specific distributions? Yes
- (e) Did you describe the limitations of your work? Yes, check Discussion-Limitations section.
- (f) Did you discuss any potential negative societal impacts of your work? NA
- (g) Did you discuss any potential misuse of your work? NA
- (h) Did you describe steps taken to prevent or mitigate potential negative outcomes of the research, such as data and model documentation, data anonymization, responsible release, access control, and the reproducibility of findings? Yes, the data we used was opensourced and our analysis does not reveal the identity of the users. Also we share our labels on GitHub for reproducibility.
- (i) Have you read the ethics review guidelines and ensured that your paper conforms to them? Yes
- 2. Additionally, if your study involves hypotheses testing...
- (a) Did you clearly state the assumptions underlying all theoretical results? Yes
- (b) Have you provided justifications for all theoretical results? Yes
- (c) Did you discuss competing hypotheses or theories that might challenge or complement your theoretical results? Yes
- (d) Have you considered alternative mechanisms or explanations that might account for the same outcomes observed in your study? Yes
- (e) Did you address potential biases or limitations in your theoretical framework? Yes
- (f) Have you related your theoretical results to the existing literature in social science? Yes, in Discussion section, we make references to previous work and compare our findings.
- (g) Did you discuss the implications of your theoretical results for policy, practice, or further research in the social science domain? Yes
- 3. Additionally, if you are including theoretical proofs...
 - (a) Did you state the full set of assumptions of all theoretical results? NA
- (b) Did you include complete proofs of all theoretical results? NA

- 4. Additionally, if you ran machine learning experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? Yes, the code is fully shared online.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? Yes
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? NA. Our machine learning experiment was done for generating text where the concept of accuracy and error-bar does not apply.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? Yes, see Data & Methodology section.
 - (e) Do you justify how the proposed evaluation is sufficient and appropriate to the claims made? Yes
 - (f) Do you discuss what is "the cost" of misclassification and fault (in)tolerance? Yes
- 5. Additionally, if you are using existing assets (e.g., code, data, models) or curating/releasing new assets, without compromising anonymity...
 - (a) If your work uses existing assets, did you cite the creators? Yes
 - (b) Did you mention the license of the assets? Yes
 - (c) Did you include any new assets in the supplemental material or as a URL? Yes
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? NA
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? Yes
- (f) If you are curating or releasing new datasets, did you discuss how you intend to make your datasets FAIR (see FORCE11 (2020))? NA
- (g) If you are curating or releasing new datasets, did you create a Datasheet for the Dataset (see Gebru et al. (2021))? NA
- 6. Additionally, if you used crowdsourcing or conducted research with human subjects, without compromising anonymity...
- (a) Did you include the full text of instructions given to participants and screenshots? NA
- (b) Did you describe any potential participant risks, with mentions of Institutional Review Board (IRB) approvals? NA
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? NA
- (d) Did you discuss how data is stored, shared, and deidentified? NA

Ethical Impact

In this work, we collect and analyze the posts and replies among users and their profile attributes. We collect this data using Twitter's Academic API before May 2023. We mitigate the potential ethical risks of analyzing user data by limiting our focus to public accounts and analyzing aggregate data. That is, we do not examine and report individual profiles. To further protect user privacy, we anonymize the user data before sharing our dataset. Our regression models do not use any personally identifiable information. Our study aims to bridge divides among the public by understanding the dynamics of non-toxic cross-partisan interactions. It should not be misused to create or amplify content that would provoke toxic content within echo chambers or manipulate the public discourse in any other way.

Appendix

Root Tweet	Video footage of WHITE people fliping over a cop car in SLC, Utah that later caught on fire.
	Our protest was peaceful. Black organizers wanted to remain peaceful. White people are respon-
	sible.[LINK]
Reply Tweet	@USER Like half of the people flipping the car are black, what the fuck are you talking about?
Prompt	You are analyzing a tweet on Twitter.: tweet: [ROOT_TWEET] Without any explanation, provide
-	three adjectives representing different sentimental qualities of the tweet. For every adjective you
	intend to provide, use their simplest form in English: \n adjective 1: \n adjective 2: \n adjective 3:
Response	agitated: The tweet expresses a sense of agitation towards the actions of the white people in the
Mistral-7B	situation. peaceful: The tweet emphasizes that the protest was peaceful, suggesting a calming or
	soothing quality. responsible: The tweet assigns responsibility to the white people, implying a sense
	of seriousness or gravity.
Prompt	You are analyzing an interaction between two people on Twitter. User A has posted a tweet and User
	B has replied to it: User A: [ROOT_TWEET] User B: [REPLY_TWEET]. Without any explanation,
	describe the stance of user B's reply toward user A's tweet with three adjectives. Make sure that the
	adjectives you provide are with regards to the stance that is directed toward user A and/or his/her
	opinion, not toward an external concept or person. For every adjective you intend to provide, use
	their simplest form in English: \n adjective 1: \n adjective 2: \n adjective 3:
Response	1. Skeptical: User B expresses doubt about the accuracy of User A's statement that only white
Mistral-7B	people were flipping the cop car. 2. Confrontational: User B challenges User A's assertion directly
	and forcefully. 3. Curious: User B seems interested in obtaining more information and clarification
	about the event described in User A's tweet.

Table 2: Sample prompt and response for LLM-aided annotation of interactions

Stance	Root Tweet	Reply Tweet	Interaction
patriotic	I'm Proud To Be An American!	@USER Love the song. Happy tears of pride.	$R \to R$
Prayerful	The Lord will make you great	@USER The Lord will make me great. Amen.	$R \rightarrow R$
Conspiratorial	Giuliani Rips Fauci, Says US Paid for 'Damn Virus That's Killing Us' [LINK]	@USER The Deep State at work with their cronies #CrooksandCommunists	$R \rightarrow R$
certain	You Vote: Do you think Antifa is a do- mestic terrorist organization? — Just The News [LINK]	@USER Yep without a doubt	$R \to R$
unsupportive	Your support means everything. I will never stop fighting for our movement. [LINK]	@USER You don't have my support.	$R \rightarrow D$
imaginative	If life were a 90s Fox drama, Harry and Meghan would abandon the royalty to start a fresh ad agency in LA.	@USER Set it at Christmas and you have a Hallmark movie in the making	$R \rightarrow D$
proactive	The Pope says tax evaders have stolen from the Government and weakened Italy's health scheme and are murderers [LINK]	@USER We need to stop cash in hand as well, it's rife, especially in London #Lock- downNow	$R \rightarrow D$
congratulatory	Oh my gosh you guys I have a son!!! This is crazy!! [IMAGE]	@USER Congratulations! Enjoy!	$R \rightarrow D$
loving	My baby is back home!! [IMAGE]	@USER That's why you put that dress onaaahhh [IMAGE]	$R \rightarrow D$
apologetic	This version of Hamburglar has definitely burned a church in Norway [LINK]	@USER "Sorry for partying"-the hambur- glar probably	$R \rightarrow D$
reciprocal	Shots from the PEACEFUL #BlackLives- Matter protest in Austin today. City made me proud! [LINK]	@USER I was there too. Incredibly hope- ful to see so many people show up for change. #BLMprotest	$D \rightarrow D$
similar	Admit it. You wish it was Sunday already. #KillingEve	@USER Oh my god yes I'm obsessed too	$D \rightarrow D$
relatable	I'm at work. I'm always at work. I'm los- ing my mind.	@USER Oh good, I was starting to think I was the only one	$D \rightarrow D$
fascinated	#WhenTheLockdownEnds I'm going to party like a Maya Ruler in a weird lobster costume! [LINK]	@USER Wow! Is that from Bonampak? What a party!	$D \rightarrow D$
empathetic	RIP Grandpa We lost him to COVID-19 last night	@USER I'm so sorry	$D \rightarrow D$
accusatory	Why bother to tweet about this, of all things? [LINK]	@USER He's happy to have people dying to improve his ratings.	$D \to R$
hostile	This is critical evidence confirming what we already knew to be true — China lied. China continues to lie. China must be held responsible [LINK]	@USER shut up, traitor	$D \rightarrow R$
dismissive	Trump, reading, says COVID-19, then adds, "You know what that is? Right. Be- come a very famous term. C-O-V-I-D. COVID."	@USER He's a child.	$D \rightarrow R$
skeptical	New: Biden says he will choose his running mate next week - CNNPolitics [LINK]	@USER He said this immediately after be- coming the nominee.	$D \to R$
critical	President Trump says a new, faster coro- navirus test is going to be used soon. He says: I hope the new test works out. Hope- fully it will check out or test out. It will be a very simple test. It won't be unpleasant at all.	@USER Trump doing the right thing after exhausting all other options.	$D \rightarrow R$

Table 3: Sample tweets for AI-generated stances for replies.

Sentiment	Root Tweet	Party
significant	Meanwhile, in other news, this important summit meeting to develop strategies to fight	Democrat
_	the pandemic.	
shameful	There is no greater embarrassment in the House of Representatives than Jim Jordan.	Democrat
moral	Anyone who thinks they need to go forward with an execution in this moment shouldn't	Democrat
	ever have the authority to carry out executions.	
alert	We are making progress. Some good recovery numbers are coming in. GHS should be	Democrat
	announcing soon. Doesn't mean we can let down our guard and live without care. Let's	
	keep observing the preventive etiquette.	
perspective	To people complaining about the wrong statues getting removed: if the right statues had	Democrat
	been removed earlier or NEVER placed, you could have avoided all this.	
responsible	If we're going to rise to this moment with the attention and action that it fully deserves, we	Democrat
	better start listening. We better own up to our own responsibility that led to this moment.	
	And we better start inviting change.	_
proactive	NEW: New York State will require all hospitals to have on hand a 90-day supply of PPE	Democrat
	at quantities sufficient to meet the rate of use during the worst of this crisis.	
balanced	The level of anger directed at the media from these protestors was alarming. As always, I	Democrat
	will tell a fair and unbiased story today.	
overwhelmed	I cannot work. The Christmas break has broken me	Democrat
humorous	My wife is a teacher. There's herding cats, and then there's "webinar with 20 1st graders"	Democrat
intimate	My mom just passed away in her sleep. Does not appear covid related. I'm numb.	Democrat
gentle	Good morning sweet girl\nHave a great day.	Democrat
unempathetic	She can stand in the unemployment line	Democrat
unconventional	There are no rules for breakfast I just had garlic cheese bread along with my berry	Democrat
	protein shake \nWhatever	_
tired	I need to be in bed. #QuarantineLife	Democrat
surreal	Even the bedbugs are wearing masks.	Democrat
isolated	Am I the only person who's never watched Friends?	Democrat
misleading	The President is rambling, reading a script full of badly formed theories. And lies.	Republican
interesting	This is the most bizzare part of a fascinating thread about how Labour responded to the	Republican
	1992 defeat.[link]	
arrogant	Trump: "Were doing a job the likes of which nobody's ever done."	Republican
radical	Al is more dangerous than nuclear weapons.	Republican
diplomatic	US Secretary of State Mike Pompeo in his opening remarks in #Doha says that the Taliban	Republican
	have committed to cut their ties with terrorist groups. #Afghanistan #PeaceDeal [link]	
formal	President Trump announces intent to appoint Dr. Sebastian Gorka, of Virginia, to be a	Republican
10	Member of the National Security Education Board for a term of four years.	D
self-centered	"We understand there are still many sick and dyingBut what about the rest of us?" [link]	Republican
alarmist	Our national media is ROOTING for hydroxycloroquine to not work as a treatment for	Republican
	#Covid_19. Think about that.	D
amused	Just got a PR pitch that includes the word "covidpreneur"	Republican
accusatory	The Chinese Communist Party is responsible for every single death, every job lost, every	Republican
	retirement nest egg lost, from this coronavirus. And Xi Jinping must be made to pay the	
	price.	D 11
negative	A rebound in 2021 will not happen if we keep doing this [link]	Republican
disgusted	What kind of man uses a funeral to push a political agenda? [link]	Republican
outraged	Wake the f*** up America! [link]	Republican
sarcastic	A face mask is part of the 'scamdemic,' they say. But they'll be happy to sell you one	Republican
C 1		D 11
prayertul	Please Lord let it be so. [link]	Republican
prejudiced	Black people love acting like kings, queens and African warriors, but never like normal	Republican
	Americans. [link]	D 11
conspiratorial	The world your mind has been MK Ultra Programmed into believing you live in, bears	Republican
	no resemblance to the World you are actually living in.	

Table 4: Sample tweets for AI-generated root sentiments.