# Unpacking polarization: Antagonism and Alignment in Signed Networks of Online Interaction

Emma Fraxanet[1*], Max Pellert[3], Simon Schweighofer[4],
Vicenç Gómez[1], David Garcia[2,5]

[1*]Department of Information and Communication Technologies,
Pompeu Fabra University, Tànger 122-140, Barcelona, 08018, Spain.
[2]Complexity Science Hub Vienna, Josefstäfter Strasse 39, Vienna, 1080,
Austria.
[3]Chair for Data Science in the Economic and Social Sciences, University
of Mannheim, L15, 1–6, Mannheim, 68161, Baden-Württemberg,
Germany.
[4]Department of Media & Communication, Xi'an Jiaotong-Liverpool
University, Suzhou Industrial Park, Suzhou, 215123, P.R.China.
[5]Department of Politics and Public Administration, University of
Konstanz, Universitätstrasse 10, Konstanz, 78464, Baden-Württemberg,
Germany.

*Corresponding author(s). E-mail(s): emma.fraxanet@upf.edu;
Contributing authors: maxpe@gmx.com;
simon.schweighofer@outlook.com; vicen.gomez@upf.edu;
david.garcia@uni-konstanz.de;

## Abstract

Online polarization research currently focuses on studying single-issue opinion
distributions or computing distance metrics of interaction network structures.
Limited data availability often restricts studies to positive interaction data, which
can misrepresent the reality of a discussion. We introduce a novel framework that
aims at combining these three aspects, content and interactions, as well as their
nature (positive or negative), while challenging the prevailing notion of polariza-
tion as an umbrella term for all forms of online conflict or opposing opinions.

1

In our approach, built on the concepts of cleavage structures and structural balance of signed social networks, we factorize polarization into two distinct metrics: Antagonism and Alignment. Antagonism quantifies hostility in online discussions, based on the reactions of users to content. Alignment uses signed structural information encoded in long-term user-user relations on the platform to describe how well user interactions fit the global and/or traditional sides of discussion. Moreover, we re-define two additional metrics that describe Alignment in more detail: Cohesiveness, the tendency of users to align with their own group, and Divisiveness, which accounts for the separation between groups. These four metrics shed light on distinctive features of conflicts around specific topics and enhance understanding of ideological alignment and developement of partisan preferences toward societal issues. We can analyse the change of these metrics through time, localizing both relevant trends but also sudden changes that can be mapped to specific contexts or events. We apply our methods to two distinct platforms: Birdwatch, a US crowd-based fact-checking extension of Twitter, and DerStandard, an Austrian online news paper with discussion forums. In these two use cases, we find that our framework is capable of describing the global status of the groups of users (identification of cleavages) while also providing relevant findings on specific issues or in specific time frames. Furthermore, we show that our four metrics describe distinct phenomena, emphasizing their independent consideration for unpacking polarization complexities.

**Keywords:** Polarization, Signed Networks, Online Media, Alignment, Antagonism

# Introduction

Nowadays, it is difficult to watch a news broadcast, listen to a campaign speech, or read a political commentary without coming across the term *polarization*. It seems that, when political commentators need a catchy, one-word description of the current state of political affairs, they habitually default to *polarized*. But this inflationary usage of the concept of political polarization lumps together very different forms of political conflict. In a world where even apparently apolitical questions of lifestyle and taste have become associated with ideological positions [1], it may seem like every political conflict is being fought along the lines of left versus right, neatly splitting the political spectrum into two opposed factions. But neither in theory nor in practice is this the only way in which political antagonism can manifest in democratic societies.

The conflation of concepts when talking about polarization also explains the seemingly ambivalent role of political antagonism in democratic societies: On the one hand, polarization is usually conceptualized as detrimental to political stability and efficient governance. On the other hand, conflict and competition are recognized as essential parts of a functioning political system. This apparent contradiction is easily resolved by stipulating that political antagonism is not automatically detrimental to the stability of the system, as long as it is not exclusively located along the same dividing line, or *cleavage*. If political antagonism is located along multiple *cross-cutting cleavages* [2, 3], it can actually increase systemic cohesion by putting political actors into ever-changing configurations of alliances. In such a system, the opponents of yesterday may become the allies of tomorrow (and vice versa), which creates an incentive to maintain a minimum of civility [4]. In contrast, if conflicts are predominantly organized along a single cleavage, political actors will always find themselves alongside, and across from, the same group of people. It is easy to see why in such a system civility tends to be replaced by partisan hostility and political sectarianism [5].

The analysis of cleavage structure has been a central concern for political scientists (especially in Europe) since the seminal work of Lipset and Rokkan in 1967 [6]. Lipset and Rokkan theorized that party systems in Western democracies are the results of four basic societal conflicts: center vs. periphery, state vs. church, owner vs. worker, and land vs. industry, which are present to differing degrees in different societies. The four cleavages initially introduced by Lipset and Rokkan in the 1960s have since lost a large degree of their explanatory power [7]. New cleavages have been proposed by various authors, determined, for example, by conflicts around globalization [8], migration [9], or European integration [10]. However, it has been criticized that, similar to 'polarization', the term 'cleavage' has been overexpanded, and thus lost most of its meaning, serving now merely as a redescription of differences in political attitudes among the electorate [11, 12].

In this study, we identify and analyze two distinct factors of political polarization: First, the degree of *antagonism* in a community, a metric reflecting the prevalence of negativity in the interactions that are triggered by a controversial issue. And second, the degree of *alignment* of a community around an issue, reflecting how much the issue 'fits', and thereby reinforces, the main dividing lines in a community. Political polarization can then be defined as the product of antagonism and alignment, both of which have to be present for a political system to fission into radically opposed

3

factions. Our quantification of antagonism and alignment is based on the identification of cohesive groups in networks of signed interactions, as well as the cleavages separating them.

Alignment and antagonism can be measured on signed networks where each node represents an individual and their interactions are represented by positive (+) or negative (-) edges. In social media, positive edges are captured by liking, praising, forming friendships, or establishing trust, while negative edges are captured by disliking, toxic behavior, hostility, or distrust. By considering explicitly negative interactions within social media, we gain a deeper understanding of community structures and relations than by only analyzing positive interactions. For example, relying only on positive interaction data creates biases that lead to an overestimation of online fragmentation and distorted pictures of the polarization of a community [13] [14]. This is particularly important when assessing the degree of political polarization in social media use, which might have been overstated due to missing information on negative interactions [15].

Balance theory [16] [17] postulates that positive interactions happen with a higher likelihood between individuals belonging to the same political faction, whereas negative interactions happen predominantly between opposed factions. Balance can also be defined by the absence of cycles containing an odd number of negative edges [17]. In practice, real-world signed networks are not completely balanced and different definitions of partial balance have been introduced, e.g., signed triangle count[18], walk-based partial balance measures [19] or frustration-based measures [20]. The latter provides a network partitioning algorithm according to a maximization of balance. Building on those partitions, we designed the *Signed Alignment Index* (SAI): a metric that captures the tendency of a network to lack *frustrated edges*, i.e. positive interactions across groups and negative interactions within groups. A high SAI thus corresponds to few frustrated edges, indicating a clear-cut division of the network into politically opposed groups. The SAI can be applied to subsets of interactions in an online network to track changes in alignment over time and to compare how alignment manifests across issues in society. This way, we can discover cleavages based on high-resolution and contextualized data as a supplementary approach to theorizing specific cleavages *ab initio*.

Furthermore, we can analyze the two independent mechanisms that contribute to alignment, namely *cohesiveness* and *divisiveness* [21]. These mechanisms account for the proportion of positive edges within groups (cohesiveness) versus the proportion of negative edges between groups (divisiveness). In this work, however, we also renormalize the original description of these metrics against a null model in order to make them independent of the proportion of negative edges. This allows us to compare the results for different sub-sets of data with different network features. At present, out-group disaffection is the most relevant variable in the steep increase of political sectarianism, especially in the US [5]. Hence, a proper consideration of negative interactions and relations is crucial to the analysis of polarization within online systems. Moreover, the separation between the mechanisms forming polarization - not only alignment and antagonism but also cohesiveness and divisiveness- is relevant for the mapping of these findings into specific discussions.

An advantage of our approach is that it is applicable to different political systems (multi-party as well as two-party), given the freedom of choice in the number of groups that the network is divided into. Our methodology does not rely on assumptions such as the dimensionality of the ideological space, nor the pre-definition of specific cleavages. Note that we only require information in the structure of positive and negative interactions, and not necessarily the content. Therefore, this framework is agnostic in terms of political system, language, or issue dimensions, as long as positive and negative interaction information is available. We apply these metrics to two unique datasets that contain positive and negative interactions between users: Birdwatch, the American Twitter system to annotate information quality; and DerStandard, an Austrian online newspaper with discussions on news pieces.

Birdwatch is a crowd-based fact-checking platform designed to combat misinformation on Twitter. The platform is operating since January 2021. Birdwatch users (also referred to as birdwatchers) can add notes to tweets assessing the trustworthiness of the tweet content and provide additional information such as claimed sources and (counter-)arguments. These notes, which other birdwatchers can see in their Twitter feed attached to the tweets in question and/or the Birdwatch site, can be positively or negatively rated. Previous work analyzing this platform has shown high political alignment and polarization among involved users [22]. Researchers found that partisanship of both original tweet authors and the birdwatchers are relevant features for the prediction of notes being perceived as misleading or helpful. Users are more likely to police tweets and notes from counter-partisans. Moreover, there also seems to be a partisan cheerleading effect in the ratings [23]. By using Bayesian Ideal Point Estimation to infer the ideology of users whose tweets appear in Birdwatch [24], we are also able to compare our results with such previous literature.

DerStandard is an Austrian newspaper that has a long tradition (dating back to the 1990's) of offering users discussion forums on their webpage. This online community is highly engaged and the platform has seen growth in the number of users and their interactions over time and especially so in the last years. For example, the site had almost 57 million visits in November 2020. Concerning demographics, a recent study shows that users that are active on DerStandard tend to be more often male, younger, more highly educated, and more often from Vienna or Upper Austria than respondents of a representative survey in Austria [25]. The advantages of using online interaction data from Birdwatch and DerStandard are that the retrieved information is based on spontaneous behavior instead of elicited reactions to questionnaires. Moreover, the dimensionality of the data allows for both temporal decomposability and evolution of the results as well as stratification in topics of discussion. While our methods to measure alignment, antagonism, cohesiveness, and divisiveness are language-agnostic, it is possible to take advantage of the additional text data to contextualize our results with other NLP analyses.

The signed network data of Birdwatch and DerStandard offer a unique opportunity to directly measure positive and negative relationships, as previous research struggled to infer negative relationship information from unsigned data [26] [27] [28]. This difficulty is particularly pronounced in online social systems, where distinguishing between users not interacting due to animosity versus chance becomes infeasible [13]. Even

the inference of positive interactions from endorsing actions, such as retweets, has been called into question [29]. Some exceptions in very precise contexts exist, such as signed graphs of political elite interactions (e.g. international relations [30] [31] [32] [33] and the US House of Representatives [34]), online platforms with particular functions away from general discussion (e.g. Epinions [35], Slashdot [36] or Wikipedia [37] [38]), and inferred signed interactions from text data in Reddit [39][40]. Our two datasets provide information on general discussions with strong political content and explicit signed interactions of the form of positive and negative ratings resembling likes and dislikes. Both datasets also have temporal information and contextualization features encoded in news tags in the case of DerStandard and text in both datasets. While Birdwatch has been studied in previous research, either aiming at understanding its content [22][41][42], its mechanisms within the platform [43][44], or even the behavior of its users [23], our DerStandard dataset is novel and comprises eight years of signed information interaction between regular users of news discussions.

# Data and Methods

## Datasets

We use two data sources: (a) DerStandard: positive and negative ratings on postings in the forum below articles on the online newspaper page. (b) Birdwatch: agreement and disagreement between raters and their notes, which we treat as positive and negative interactions. Both datasets exhibit an exceptional combination of features: the explicit sign (+ or -) of the user interactions and temporal information (timestamp of postings or note). We differentiate between: (i) *Interactions*: directed pairwise interactions based on the reaction of a user (rater) to the content posted by another user (author), with the timestamp corresponding to the posting of that piece of content, and (ii) *Edges*: undirected and signed relations between users of the platform, based on aggregated interactions exchanged between them through their postings or notes.

### Network Creation

Both datasets contain pairwise interactions between users. Considering a dataset of $n$ users, we model each relation between user $i$ and user $j$ from such interactions as a random variable that follows a Bernoulli distribution with parameter $p_{ij}$. We follow a Bayesian model using a beta prior for estimating $p_{ij}$ with parameters $\alpha_0, \beta_0$. After observing all the interactions between $i$ and $j$ in the dataset, the posterior probability also follows a beta distribution, in this case parametrized by $\alpha_0 + \text{pos}, \beta_0 + \text{neg}$, where pos and neg correspond to the number of positive and negative interactions respectively.

From these posterior probabilities, we build an undirected signed network $G = (V, E, \sigma)$, where $V$ is a set of $n$ nodes, $E$ is a set of $m$ edges, and $\sigma_{ij}$ is the edge sign. Edges are only defined for pairs of users who have a certain bias towards 0 or 1, i.e., $\mathbb{E}[p_{ij}] > 0.6$ or $\mathbb{E}[p_{ij}] < 0.4$, and very low uncertainty, i.e., $\text{Var}[p_{ij}] < 10^{-4}$, where $\mathbb{E}[p_{ij}] = \frac{\alpha}{\alpha+\beta}$ and $\text{Var}(p_{ij}) = \frac{\alpha\beta}{(\alpha+\beta+1)(\alpha+\beta)^2}$. For defined edges, we set their sign according to $\sigma_{ij} = \text{sign}\left(\mathbb{E}[p_{ij}] - \frac{1}{2}\right)$, i.e, two users have a positive (negative) edge if their expected posterior is above 0.6 (below 0.4) with high certainty.

### Birdwatch

Twitter regularly publishes updated and publicly available datasets containing metadata of notes (text, tweet ID, note timestamp, tweet reliability classification, note ratings, rating timestamp) and anonymized birdwatchers data. We retrieved all data covering the time span between the start of Birdwatch in January 2021 and August 2022. In this period, the platform was in a pilot stage and had limited user participation (within the US). Moreover, we re-hydrated the the content and metadata of the original tweet that the note was attached to with the academic access to the Twitter API and computed an ideology score with *tweetscores* [15]. This package provides scores based on Bayesian Ideal Point Estimation: a model that considers ideology as a latent variable that can be inferred by examining which political actors a user is following. Therefore, we can only retrieve a score for tweet authors that have connections (follower or followee relation) to political actors, which leaves us with about 60% of the users.

During the time span covered by our data, the platform changed their rating procedure: first, ratings were cast as a simple `agree` versus `disagree` (Jan 2021 - Jun 2021). Then, Birdwatch introduced a new scheme with differentiation between `helpful`, `somewhat helpful` and `not helpful` for the remaining months in our data set (Jul 2021 - Dec 2022). Moreover, the platform launched a new algorithm to compute note statuses in February 2022, which searched for agreement across different viewpoints [43]. Since these are substantial platform changes, we split the dataset into two parts accordingly: BW1 and BW2, and center our study mostly on BW1, leaving BW2 as comparison only since it comprises a series of platform changes. For the second part, we consider `helpful` and `somewhat helpful` as positive interactions and `not helpful` as negative interactions. On this platform, positive and negative interactions are present in similar proportions (see Table 1). Both possible interactions in Birdwatch, agreement and disagreement, can be considered to be more meaningful than a simple rating in DerStandard because they require an argumentation. Consequently, we use a uniform prior for the beta distribution that characterizes the user relations on Birdwatch.

**DerStandard**

With permission from DerStandard, we automatically retrieved all publicly available postings in the discussion forums below each news piece on DerStandard between Jan 2014 and Dec 2021. DerStandard allows users to positively and negatively rate postings of other users in a similar way as how likes and dislikes operate on other media. Compared to other platforms with similar features, DerStandard uniquely provides information on which users rated a posting in addition to the sign of the rating (see SI for an example of the interface). In addition to postings and ratings, we also retrieved tags that classify news pieces into topics according to the platform (e.g. sports, refugees in Austria, Op-Ed columns, etc).

To avoid influence due to fluctuations (strong influxes or losses of users), we consider only users that voted at least once yearly in our observation period (begin of 2014 - end of 2021). This allows us to identify roughly 14,827 users that we track over 8 years. Our observation period includes a number of major events including the highly contentious European refugee crisis (2015-16), a notoriously turbulent year regarding corruption scandals that led to the dissolution of the Austrian government coalition

(2019), and the two years comprising the COVID-19 pandemic (2020-21). However, being the COVID-19 pandemic an exceptional and unusual circumstance globally, we proceed to use the partitions obtained from the 6 previous years only for the analysis. To test this assumption, we run the partitioning methods using our data with and without the years comprising COVID-19. We find an overlap of above 80% between the two different partitions.

On DerStandard, negative interactions are underrepresented (see Table 1), thus they encode a stronger signal than positive interactions. To account for that, we use a prior distribution that slightly favors negative interactions, especially when the volume of interactions is low, i.e., a beta distribution with $\alpha = 1$ and $\beta = 2$. The resulting network contains a similar number of negative and positive edges.

To give an overview of all of the data we use in our study, Table 1 summarizes the basic descriptive statistics of all datasets (DerStandard, BW1, BW2).

|  | Timespan | Users | Edges | Interactions |
|---|---|---|---|---|
| **BW1** | ∼5 months | 2,676 | 25,562 (28% negative) | 32,323 (28% negative) |
| **BW2** | ∼12 months | 10,662 | 235,493 (38% negative) | 301,041 (38% negative) |
| **DerStandard** | 8 years | 14,827 | ∼5.56M (41% negative) | ∼76M (17% negative) |

**Table 1 Summary statistics of the datasets used to build signed relation networks.**
BW1 and BW2 are the two networks obtained from the Birdwatch platform.

## Measuring Partial Balance

**Main optimization problem.** A signed network is balanced if it can be partitioned into $k \leq 2$ groups such that all negative edges fall outside the partitions and all positive edges fall within the partitions. Following [20] notation, given a signed graph $G = (V, E, \sigma)$, and a partition $P = \{X, V \setminus X\}$, the frustration count will be the sum of the frustration state of all edges, $f_G(P) = \sum_{(i,j) \in E} f_{ij}$, where $f_{ij}$ equals 1 for frustrated edges and 0 otherwise. Frustrated edges correspond to the edges that violate the assumptions of the optimal partition model, i.e. negative edges between members of the same partition or positive edges between members of different partitions. The problem thus is stated as finding the optimal partition $P^*$ with the minimum number of frustrated edges $L_G^* = \min_P f_G(P)$. The value of $L_G^*$ can be used to compute partial balance.

**Computational methods.** The computation of $L_G^*$ is known to be NP-hard [45]. For small scale networks, however, exact computation of the frustration index is feasible using the binary linear programming formulation [45]. Several approximate methods have been proposed that are applicable to large scale networks. For example, Doreian and Mvar apply blockmodeling [46], in which they optimize the criterion function $P(X) = E_{f,p} + E_{f,n}$ via a relocation algorithm, with $E_{f,p}$ defined as the frustrated positive edges and $E_{f,n}$ the frustrated negative edges. In practice, this method, in combination with simulated annealing, provides approximate values of $L_G^*$ that correspond to robust partitions (see SI for details). We use the *Signnet* implementation [47]. Any approximated value for $L_G^*$ will necessarily be equal or higher than its

exact value, given that there is no not-optimal partition that can provide a smaller number of frustrated edges, thus it will be an upper bound.

All the previous definitions and methods are generalizable to $k > 2$ partitions [48][34], which corresponds to a definition of weak structural balance. In that case, each value of $k$ provides an optimal solution $L_G^*(k)$, and a reasonable selection is to keep $k$ with minimum $L_G^*$. In [46], it is shown that $L_G^*$ follows a concave curve with a unique minimum value of $k$, which we refer to as $k^*$. See SI for details in this multi-partition selection for our data.

## Contribution of our approach

The intuition behind our approach is that in a signed social network, the minimum number of frustrated edges will capture the degree up to which the network can be easily separated into groups. This allows us to find a grouping of users that is informative of the global division lines of the community. Moreover, we understand this level of separation as a structural measure related to polarization under the assumption that the groups are of similar sizes, since otherwise it would be a scenario of fragmentation with a minority group. We verify this assumption after finding the optimal partition into groups.

***Normalization of metrics.*** Since we are interested in quantifying the level of structural polarization in our networks, we look for an index that ranges from low to high balance, such as $1 - \frac{L_G^*}{m/2}$ (the "Normalized Line Index of Balance" [21]) with 1 being the completely balanced case. In this index, the $m/2$ term accounts for different network sizes and is an upper bound on the number of frustrated edges. In our framework, balance and structural polarization are equivalent and both grow in an inverse trend compared to frustration in the system. The more frustration there is, the more blended the groups are, and the less polarized the global system is. To be able to compare balance across subsets of our data, we renormalize it by comparing the empirical estimate of $L^*$ versus its mean value in repeated measurements of a null model. The null model based on graph $G$ randomly re-distributes sign attributes while keeping the partition fixed ($\widetilde{G}$). What we see from comparing the obtained measures to the null model is the amount of balance due to the signed nature of our networks, while keeping fixed the overall structure of the network and the assignment of nodes to groups. The value of $L$ in the null model simulations is consistently higher than the frustrated edges in our datasets, proving to be a tighter bound than only considering the number of edges with the term $m/2$. Thus, we define the *Global Signed Alignment Index* as:

$$SAI_G = 1 - \frac{L_G^*}{\langle L_{\widetilde{G}} \rangle}$$

Note that we only change the normalization factor from $m/2$ to the mean of $L$ in the null model, which is also dependent on size, and thus the result is simply a stricter upper bound.
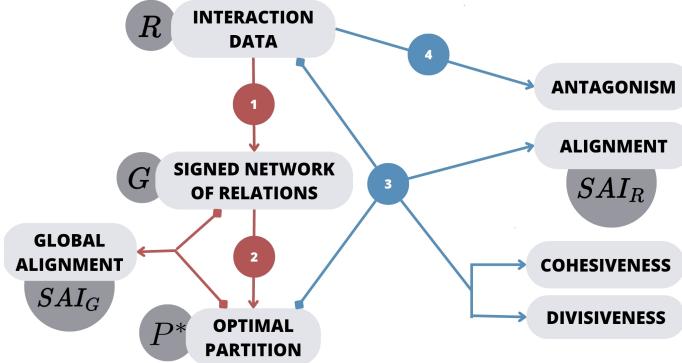
9

**Fig. 1 Schema of our analysis framework for antagonism, alignment, cohesiveness, and divisiveness.** Grey boxes indicate data structures and variables implicated in the pipeline. Step 1 creates the relation network based on an aggregation of interactions through time. Step 2 applies the optimization algorithm, either exact or approximated, to obtain the optimal number of groups and optimal partition. From these two steps we can retrieve a global alignment metric $SAI_R$. Then, by selecting subsets of the interaction data and with the optimal partition information, steps 3 and 4 compute the four metrics of interest: Antagonism, Alignment, Cohesiveness (normalized) and Divisiveness (normalized). See Section Methods for details on each of these metrics.

## Framework Pipeline

After constructing the signed relations based on the interaction data of the full dataset of each network, we calculate the optimal partitions by either using the exact or approximated method. Given that the approximated method involves a stochastic algorithm, we execute it 200 times for different $k$ values and select the number of groups and partitions yielding the minimum $L_G^*$ value. Further details regarding this approach can be found in the SI.

Our normalization approach allows us to obtain a meaningful $SAI$ for sub-sets of the interaction data. To do so, we maintain the optimal partition obtained from the network of relations (i.e. we fix the belonging of each user to a group that is defined by the long-term relation between users), and we proceed to assess how aligned the interactions within that subset of the data are to these partitions. Since this alignment follows the same laws of frustration (e.g. negative interactions within a group are frustrated interactions, and so on), we just have to re-define the $SAI_G$ in the following way:

$$SAI_R = 1 - \frac{L_R}{\langle L_{\widetilde{R}} \rangle}$$

where $R$ accounts for the network of directed interactions within a set or subset of the data, denoted by $R(t)$ in case of a temporal subset, or $(i)$ for a selection based on issue or topic. $L_R$ is then the number of frustrated interactions in that network given the existing assignment of nodes to groups. As in the case of $SAI_G$, $\widetilde{R}$ denotes an instance of the null model applied on $R$, by reshuffling the sign configuration while keeping the network structures and groups. See Figure 1 for the full step-by-step methodology

10

to obtain these measures. We refer to the $SAI_R$ measure as *Alignment* in Section Measurement of Alignment and Antagonism.

Additionally, we formally describe *Antagonism* as the proportion of negative interactions in $R$, which is a simple indication of the amount of conflict or general disagreement. This measure is then not related to the network structure, like Alignment, but it indicates a property of the user-content interaction in terms of the overall presence of disagreement in comparison to agreement. To obtain confidence intervals for our $SAI_R$ and $SAI_G$ measures, we propagate the uncertainty obtained from the reshuffled model. To achieve this, we conduct $10,000$ instances of the null model.

### Divisiveess and Cohesiveness

Similar to earlier work [21], we analysed the two mechanisms that are involved in the alignment of users to the partition: alignment with one's own group (cohesiveness) and alignment against the opposing group (divisiveness). Cohesiveness (divisiveness) is defined by the proportion of internal (external) edges that are positive (negative). Given our optimal partition $P^*$, internal edges are defined by $E_p^e = \{(i,j) \in E | i \in X, j \notin X$ or $j \in X, i \notin X\}$ and external edges are defined by $E_p^i = \{(i,j) \in E | i,j \in X$ or $i,j \notin X\}$.

The measures of cohesiveness and divisiveness defined above cannot be compared between systems with different ratios of negative versus positive interactions. For example, a system with a higher ratio of negative interactions will have by construction a higher divisiveness even if it is not more strongly divided along the division between groups. This can be observed in simulations of our null model, which show that the expected value of divisiveness and cohesiveness is correlated with the fraction of negative interactions, i.e. Antagonism (see SI). To solve this, we use our null model that maintains the network structure and only randomizes interaction signs. We calculate new metrics of Normalized Divisiveness and Normalized Cohesiveness by subtracting the mean of divisiveness and cohesiveness over null model simulations from the original measure of each one. We can compute these normalized measures for both the network of relations, providing a general overview of the cohesiveness and divisiveness in each network, and also for subsets of interactions associated with topics or time periods, to provide an insight into how cohesion and divisiveness vary and explains changes in overall alignment. Furthermore, given the groups found for the network and the directionality of ratings on the network, we can calculate the contribution from each group to normalized cohesiveness and divisiveness, thus allowing us to examine the drivers of increases and drops in alignment over time.

We assess the uncertainty of our measurements of divisiveness and cohesiveness through bootstrapping. For each measurement of divisiveness and cohesiveness, we create $10,000$ bootstrap samples of the interactions with replacement and of the same size as the original. On each bootstrap sample, we calculate divisiveness and cohesiveness and we use the resulting values to calculate the bootstrapping interval around our original measurement.

11

# Results

## Measurement of Alignment and Antagonism

The metric of *Alignment*, $SAI_R$ (defined in detail in Section Framework Pipeline), captures how interactions follow the division of the network into opposed groups, while our metric of *Antagonism* captures the overall tendency towards negative interaction in the network regardless of groups. By considering both these measures, we can provide a more comprehensive picture of polarization than when these two concepts are not explicitly distinguished. Figure 2 shows how these two metrics capture various polarization scenarios given a partition of the network into groups and the positive and negative interactions in the system. A network with low alignment and low antagonism has few negative interactions and no strong division into groups, corresponding to a situation with the weakest polarization. The lower right part of the space, where alignment is high but antagonism is low, corresponds to an *echo chamber* case in which most interactions are positive but happen between like-minded individuals and not across groups. The upper left cases are networks with high antagonism but low alignment, capturing scenarios where disagreement exists but not necessarily following the division of the network into groups. This can happen when everyone is against everyone or where other divisions exist but do not follow the general ideological separation of the network into groups. And finally, the upper right part of the space corresponds to cases where polarization is high, as both antagonism and alignment are high. In this high-polarization case, there is a strong cleavage between groups such that positive interactions are confined within groups while frequent negative interactions happen mostly across groups.

## Approximating Alignment in Birdwatch

In this section, we evaluate our methodology and its performance based on the results obtained from the two Birdwatch datasets. We use Birdwatch for two key factors. Firstly, as described in Section Measuring Partial Balance, we can run the exact method for small networks but for large networks we have to run the approximate algorithm due to the complexity of the problem. The size of BW1 allows us to run both the exact and approximate algorithms and compare the solutions to estimate the difference in signed networks of this kind. The results of both algorithms are very similar in BW1, with the approximated $SAI_G$ being 84% of the $SAI_G$ obtained with the exact method and an average partition overlap coefficient [49] of 0.89.

For both BW1 and BW2, we find that the optimal number of groups is $k^* = 2$, and the largest groups contain roughly twice the number of users of the second group (see SI for more details). Figure 3 shows the signed network of relations obtained from BW1. Previous literature focused on Birdwatch suggests that the platform is characterized by two opposing factions, corresponding to Republican- and Democrat-leaning users, who attach notes to tweets following behaviors of *counter-partisan policing* and *inner-partisan cheer-leading* [23]. By building on the ideology score extracted from the tweets, we test whether the groups identified through our method reproduce this
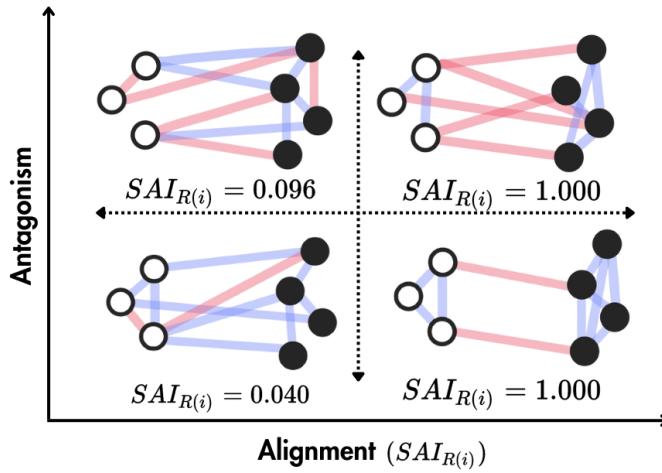
**Fig. 2 Illustration of measures used in this study in a stylized signed network.** The four depicted quadrants represent four different combinations of measures of antagonism and alignment. The four networks have been drawn with the same edge density, number of nodes, and partitioning of nodes for purposes of easier comparison. Negative edges are represented as red, while positive edges are blue. The two upper networks have a higher proportion of negative edges, and thus higher antagonism, than the ones on the lower quadrants. Visual comparison and inspection of the provided numerical $SAI_{R(i)}$ values shows that the two right quadrants exhibit a higher level of alignment, which is due to the lower amount of frustrated edges. Only the right upper quadrant would coincide with a strict definition of polarization in terms of both antagonism and alignment.

behavior, thereby evaluating the coherence of our approach with other metrics of political alignment.

When we retrieve the notes that users from each of these partitions have given to tweets, we find evidence of these policing-cheerleading patterns, as our largest group - which we denote as *inferred Democrats* - is strongly biased towards tagging Republican-leaning tweets as misleading. Contrarily, the smaller partition - *inferred Republicans* - consistently rates like-minded tweets as not misleading (see Figure 3).

## Evolution of Aligntment in Birdwatch

Figure 4 shows the time series of $SAI_{R(t)}$ in BW1. The fluctuations in the measure over time indicate whether the level of alignment among interactions increased or decreased during that particular period. The time series of normalized cohesiveness and divisiveness contextualize these movements, as they show whether peaks are due to higher cohesion within groups or higher division between groups, and what is the contribution of each group to these metrics. In this time series, antagonism and alignment have a low correlation, which emphasizes the need to consider them as two different measures. See SI for details.
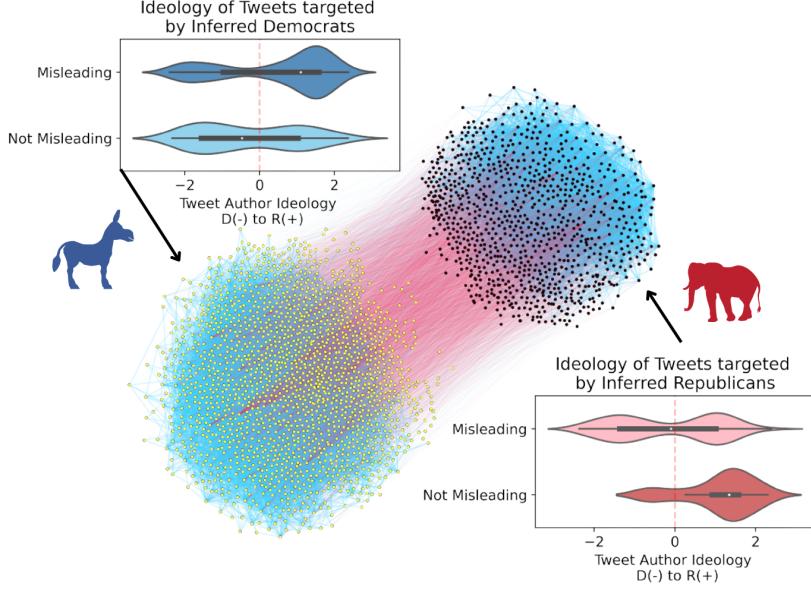
**Fig. 3 Signed network visualization of Brandwatch.** The figure depicts the network of signed relationships for the BW1 dataset, comprising a total of 2,676 users and around 25,562 edges, negative colored red and positive blue. Node color corresponds to their group membership as identified by the exact method. Nodes belonging to the largest group are depicted in yellow, while nodes from the second group are illustrated in black. Negative edges tend to connect different groups, while positive edges predominantly connect nodes within groups, demonstrating a considerable degree of balance. **Insets:** Inferred ideology of the targeted tweet's author separated by which group targeted the tweet and the nature of the note. The larger group gives misleading notes with more probability to tweets authored by Republican users, i.e. counter-partisan policing, with a slightly higher tendency to give not misleading notes to tweets by Democrat users. Thus we identify the larger group as Democrat-leaning. The smaller group is much more likely to give not misleading notes to tweets authored by Democrat users, showing a pattern of cheer-leading within Republicans and thus being identified as Republican-leaning.

We applied a peak detection algorithm and identified five local maxima of $SAI_{R(t)}$ that are marked in Figure 4. To understand the context of the tweets on the day of each peak, we generated wordshift diagrams [50] for each peak in comparison to the rest of the tweets, which can be found in the SI.

Table 2 shows the most important keywords of each peak, illustrating that peaks of alignment happen around controversial topics in the US. For example, we see that the second peak, associated with terms related to COVID-19 vaccination, is driven by an increase in divisiveness, especially from Democrat-leaning users. Alternatively, the third detected peak, which is associated with terms about police shootings, has a stronger contribution of cohesiveness, especially within the Republican-leaning users. The other three peaks (1st, 4th and 5th) are driven by a mix of cohesiveness and divisiveness. The keywords and events at those time periods point towards discussions regarding the US Government and its policies, Donald Trump and 2020 election results,

and other relevant events such as the Capitol insurrection, the Texas Power Crisis and the conversation around banning Critical Race Theory in schools in the state of Florida.
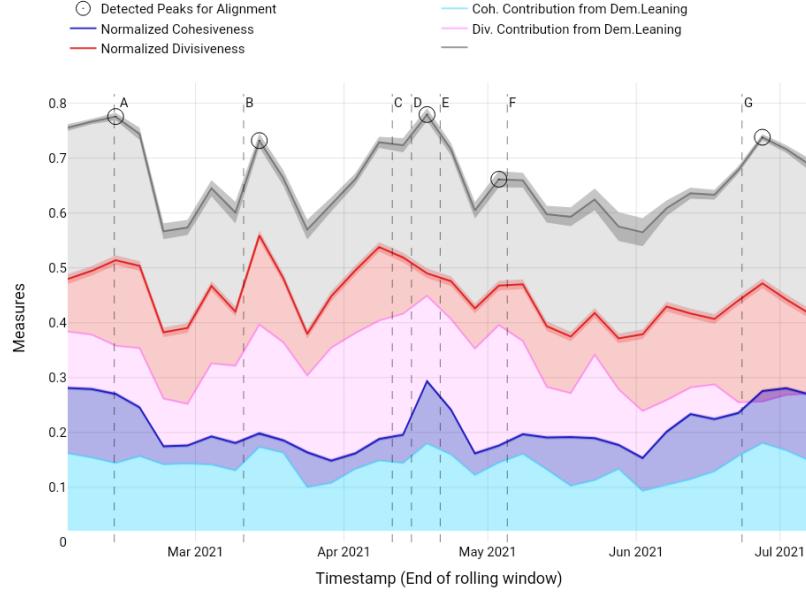


**Fig. 4 Timeline of Alignment, Normalized Cohesiveness and Normalized Divisiveness for BW1.** The time series of each metric is calculated over a rolling window of ten days with increases of 5 days, with values allocated on the right of each window. Shaded areas around each time series show 95% Confidence Intervals calculated against 10,000 instances of the null model. Normalized divisiveness is shown in red and normalized cohesiveness is shown in blue, with lighter areas showing the contribution of Democrat-leaning users to each metric and the remaining area above showing the contribution of Republican-leaning users. Bootstrapping intervals in normalized divisiveness and cohesiveness are obtained for 10,000 re-sampling instances. The alignment measure, $SAI_{R(t)}$, oscillates around a mean value of 0.65. Normalized divisiveness stays consistently above normalized cohesiveness, showing that negative interactions are the main driver of alignment. Detected peaks in $SAI_{R(t)}$ are marked with circles and notable political events in the US are marked with vertical dashed lines for reference. For each peak, a summary text analysis of tweets in that period is shown in Table 2, which can be further contextualized as increases in cohesiveness, divisiveness, or both. An interactive version of this plot can be found at https://emmafrax.github.io/BW1.html.

## Results for DerStandard

Our approach to detecting groups in the DerStandard network shows that this network has an optimal $k^*$ of two groups, as in Birdwatch. The size of these groups is similar, with the largest one comprising 62% of the nodes. Even though the DerStandard dataset spans a much longer period and contains more users than the Birdwatch datasets, the Alignment of the network is substantially high ($SAI_G = 0.3955$), showing

| Peak | Period covered | Wordshift keywords |
|---|---|---|
| **1st** | February, 7th - February, 17th | Trump, Energy, Vote, Impeachment, Trial, Plan, Power, Job, Clear, Start |
| **2nd** | March, 9th - March, 19th | Read, Vaccine, Give, Call, Death, Story, Fact, Make, Stop, Covid |
| **3rd** | April, 13th - April, 23rd | Police, Black, Kill, Shoot, Murder, Justice, Girl, Cop, Name, Veredict |
| **4th** | April, 28th - May, 8th | Election, Want, Trump, Violation, Get, School, Duck, Pandemic, Go, Big |
| **5th** | June, 22nd - July, 2nd | Get, Theory, Crime, Likely, Say, Government, Pay, Right, Voter, Collapse |

| Tag | Date | Event summary |
|---|---|---|
| A | 12th February 2021 | Governor Abbott Issues Disaster Declaration in relation to the Storms and Power Crisis in Texas |
| B | 11th March 2021 | President Biden to Announce All Americans to be Eligible for Vaccinations by May 1 |
| C | 11th April 2021 | Killing of Daunte Wright (20 years old) by the police during a traffic stop for an outstanding warrant |
| D | 15th April 2021 | Release of a relevant body cam video of the killing of Adam Toledo (13 years old) by a CPD Officer |
| E | 21st April 2021 | Killing of Ma'Khia Bryant (16 years old) by a police officer in Columbus, Ohio |
| F | 5th May 2021 | Facebook's Oversight Board upholds ban on Trump |
| G | 17th June 2021 | Biden-Harris Administration Announces Comprehensive Strategy to Prevent and Respond to Gun Crime |

**Table 2 Wordshift keywords of peaks in alignment and notable events during the BW timeline.** Upper table shows the first 10 keywords relevant to the context of each identified peak in Figure 4. We obtain these words by comparing the text of tagged tweets posted in a period surrounding the peaks with the text in the rest of the dataset. In the second part of the table, we collect events that occur close to the time of detected peaks and that help interpret the keywords above.

that alignment can appear across different sizes and time scales. Normalized divisiveness (0.2899) is substantially higher than normalized cohesiveness (0.1409), also mirroring the results for Brandwatch. More details on these results can be found in the Supplementary Information.

Given the classification of news in DerStandard, we can measure Alignment and Antagonism on the full set of user ratings around that topic, thus locating topics in the space of network structures shown in Fig. 2. The scatter plot for Alignment and Antagonism of DerStandard topics is shown in Figure 5, where the spread of values allows for all four combinations outlined by our approach. Alignment and Antagonism have a low correlation across topics (−0.0016), suggesting that these two concepts should not be conflated into a general dimension of polarization. By inspecting the topics falling into each quadrant of the plot, we find their distribution agrees with intuitive expectations. For example, topics with a high conflict potential such as migration, COVID-19 politics, gender politics, climate change, and elections are on the high range of antagonism, whereas lifestyle, sports, and culture topics such as movies, family, travel, art market or international football are located in the low ranges of antagonism. With regard to the dimension of alignment, we find that conflicting topics such as national elections, abortion, military service, or climate change are more aligned than migration or COVID-19 politics. These last two were indeed issues that did not divide the population clearly into left and right. Note that these patterns cannot be explained

16

by the number of ratings, posts, or articles on each topic, as shown more in detail in the SI.

We highlight a few examples within each quadrant of Figure 5 to better illustrate how Alignment and Antagonism relate to each other. While *Refugees* and *COVID-19 politics* are identified as conflicting topics, resulting in higher levels of antagonism, they do not align precisely with the primary division line. During the crucial years for those topics of 2015/16 and 2020/21, we have seen some unexpected political alliances that do not follow from a classical left-right spectrum. These include common platforms between the anti-migration left and right-wing populists or the anti-statist right and anti-vaccine parts of rather left-wing Green parties. These agreements on certain issues between otherwise ideologically distant parties have historically been described by the term "Querfront" ("cross-front") [51]. Conversely, the tag *National elections* exhibits both antagonism and alignment, indicating a combination that favors polarization. This can be explained by federal elections in a representative democracy to lead to more discussion along traditional party lines. Additionally, *Corruption allegations* pertain to specific events involving some of the political parties in Austria. Although it demonstrates alignment, these particular events did not generate substantial conflict within the platform. This could potentially be due to a limited number of defenders of those specific parties that have been covered much in the news in a corruption context (FPÖ and ÖVP, resulting from their joint government coalition), as DerStandard is historically considered a more leftliberal-leaning newspaper. As expected, a more offtopic tag such as *Movies* exhibits low levels of both alignment and antagonism.

While Antagonism and Alignment across topics are weakly correlated, normalized cohesiveness and divisiveness are strongly correlated, as shown on the left panel of Fig 5. This is expected, as the affective component of polarization captured by alignment implies a correlation between out-group animosity and in-group support. Nevertheless, there are topics that deviate from the association between cohesiveness and divisiveness by having substantially higher divisiveness: BVT (Institution), Abortion, Scheuba (Austrian comedian) and ÖVP (Political Party) (see right panel of Fig 5), while this pattern is not mirrored for cohesiveness. As with the time series of Alignment on Birdwatch, measuring cohesiveness and divisiveness is informative even though they both form part of the same phenomenon of alignment.

The time series of alignment in DerStandard reveals how cleavages become salient around politically-relevant events. Figure 4 shows the time series $SAI_{R(t)}$ for all DerStandard discussions in news on three topics: national elections, parties, and the federal president. This highlights political discussions from other, less-contentious topics as identified above. There is a clear change in the trend of alignment at the beginning of 2016, showing steady growth up to the beginning of 2017. This falls into the time period of the so-called "2015 European migrant crisis" [52] when migrants arrived in Europe in numbers that were unprecedented since World War Second. While migration started before 2016, the rise in alignment starts right after the reporting of sexual assaults during New Year's Eve 2015-2016 celebrations in Cologne, Germany[53], which were widely covered in German-speaking media and debated over the following year.

Political events can also drive decreases in alignment, especially if we consider that Austria has a multi-party system. After an election, the political climate changes
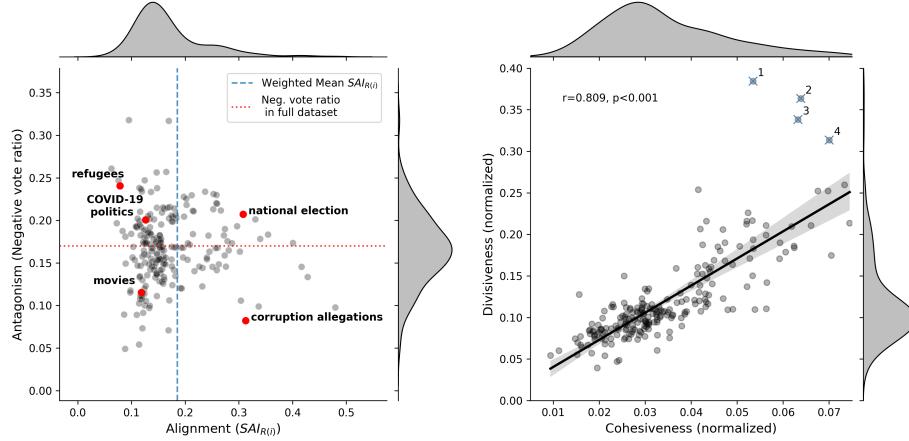
17

**Fig. 5 (Alignment versus Antagonism and Cohesiveness versus Divisivness across Der-Standard topics.** The left panel shows Antagonism and Alignment of the ratings of each news topic in DerStandard. Topics have been selected based on the topic/subtopic tags associated with the articles located above the postings (e.g., sports, climate change, etc.). Dashed lines show the mean values of each metric to identify the quadrants depicted in Figure 2 An interactive version of this figure can be found at https://emmafrax.github.io/scatter.html. The right panel shows the scatterplot of normalized divisiveness versus normalized cohesiveness for DerStandard rating sub-sets based on topics. These two measures, which account for two different mechanisms that define alignment, have a significant correlation across topics of 0.8. The highlighted outliers correspond to: (1) BVT (Institution), (2) Abortion, (3) Scheuba (Austrian comedian) and (4) ÖVP (Political Party)

toward building government coalitions with multiple parties, thus predicting lower alignment as suggested by the case of online networks of Swiss politicians [54]. This can be observed in the time series of alignment in DerStandard if we zoom in to recent elections. Panel A of Figure 4 shows the timeline of alignment during 2016, where the increase in alignment that year accelerates after the result of a presidential election was overturned by the Supreme Court of Justice. This controversial decision lead to a period of increased alignment towards the repetition of the election, to then quickly reset to earlier levels of alignment as soon as the repeated election took place and a candidate won by a large margin.

Panel B of Figure 4 shows a decrease in alignment that happened shortly before the 2017 legislative elections, which was called early since there were clear favorite parties to form a coalition in pre-election polls. The effect of the legislative elections in 2019 showed a sharp decrease in alignment afterward, as the result was not as clearly expected as in 2017 and which led to a new government coalition with a party that was not involved in the previous government.

# Discussion

### Assumptions and limitations of our methods

A conceptual assumption is that our framework, as well as most work on balance in signed networks, switches from Heider's triangle of two people and one issue, towards
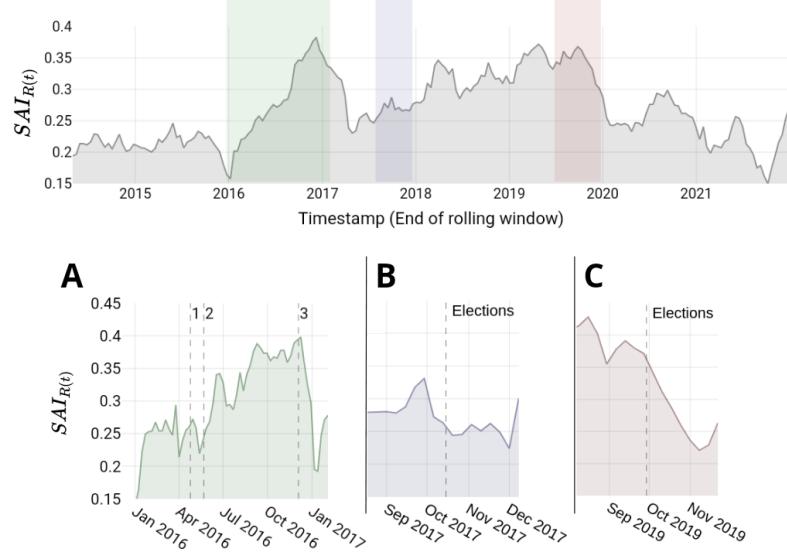
18

**Fig. 6 Alignment timeline in DerStandard voting sub-set of political topics, with detailed fluctuations in election periods.** Upper timeline figure shows the Alignment measure obtained using a rolling window of 120 days of width and a step of 14 days. The features of the rolling window are selected so that the trends in alignment through the eight years are visible, e.g. the change in trend at the start of 2016. In the lower figures we show more detailed changes of alignment, with a rolling window of 30 days of width and a step of 7 days, around the three repetitions of the 2016 Presidential elections (A: 1,2 and 3) and the 2017 and 2019 Legislative Elections (B and C).

three people. Moreover, we also have to consider if users on these platforms are directing their interactions exclusively on the other user, their content, or both. However, since we are using aggregated interactions to construct our network of relations, and ignoring those relations built on an ambiguous set of interactions (high uncertainty), we can infer more meaningful relations between users.

In our methodology pipeline we have a strong assumption on the fixed belonging of users to a partition. We are assuming there is a global clustering to which users are aligned. This is not too far-fetched given the fact that there tends to exist issue alignment in society [55][1] and we control for different numbers of clusters. However, it is true that for time scales as large as the DerStandard dataset, we could be missing relevant changes in the global structure, especially considering COVID-19 times. However, we do test this assumption and control it by not using rating data during the pandemic to construct the network of relations. For future work, it could be interesting to include a tracing system that assesses the partition quality through time and updates it accordingly. This would also be very useful to automatically detect shifts in the main lines of division.

Finally, we are interested in studying the behavior of the platform's groups. We are not aiming at classifying users individually, and therefore there could be users that occupy a more neutral positioning in our strict partitions. This is why our $SAI$ scores

19

are not fully 1 (and we do not expect them to be). If we assume these "moderate" users tend to behave similarly through time and across topics, fluctuations or differences in alignment are still relevant.

The fact that we have to use approximated methods to find (near)-optimal partitions for large scale is also a downside of our methodology. Precisely those users that are not clearly positioned are probably the ones that end up falling on different sides of the division line for different sub-optimal solutions. Even in that situation, we still capture significant values for our metrics and our approximated results are comparable to the exact results for BW1 solutions, which brings us to the conclusion that we are still measuring what we aimed to, even if not at the highest accuracy possible. Moreover, even in the exact solution, it is not possible to ensure a unique single optimal partition, since the method only ensures a unique solution for the minimum amount of frustrated edges, and several partitions can satisfy that requirement [20].

**Advantages and applications**

We define a methodology that is language agnostic and can be applied to other platforms, as long as we can extract or infer positive and negative interactions (like the examples we proposed above). Our framework can be extended to similar use cases, as well as tuned according to platform design choices, for example choosing the prior distribution for user relations or computing the number of optimal groups ($k^*$) in the search of partitions. In the specific cases of DerStandard and Birdwatch, for example, we were able to retrieve a division in the ideological spectrum (left .vs. right). Therefore, it allows us to study the main cleavage in a platform's community without the need of classifying users by their opinions a priori. It can, however, be combined together with methods such as topic modeling or other text analyses to provide a better insight into the discussion themes.

**Conclusions and closing remarks**

We were able to factor online polarization into two dimensions: Antagonism, which represents the degree of conflict, and Alignment, which is determined by the balance in our signed networks of relations. These two measures, although both contributing to polarization, have distinct characteristics and are weakly correlated across topics. We discovered that large-scale online political discussions exhibit an underlying polarized structure based on balance, which becomes more prominent when examining discussions centered around aligned topics. One important conclusion is that online polarization is a dynamic and reactive phenomenon that is heavily influenced by current political and social events. It demonstrates short reaction times, but by examining a sufficiently large time frame, such as in the case of DerStandard, we can observe general trends in addition to specific peaks.

With regard to conclusions drawn from the study of Birdwatch, we found that changes in polarization can arise from different mechanisms within one or both of the groups. Additionally, the identification of Republicans and Democrats provides valuable insights into the status of each topic in online discussions. It is important to note that our findings on Birdwatch, as a platform dedicated to crowd-sourced fact-checking, can be beneficial to understanding the dynamics and effectiveness of using this approach to combat misinformation.

Through our analysis of DerStandard, we discovered that topics such as COVID-19 politics and Refugees, despite being controversial and relevant in discussions, were not aligned with the general Left-Right spectrum in Austria. This finding sheds light on the political divisions within Austria and serves as evidence that our methodology is capable of identifying cross-cutting cleavages. Furthermore, through an analysis of the temporal trends of alignment pertaining to politically relevant topics, our findings demonstrate coherence with expected behaviors given the context of the respective time frames.

Apart from these platform-specific insights, our work contributes in several other ways: Firstly, we have provided the first framework for conducting temporal analysis of structural balance in large-scale online political discussions. This enables a deeper understanding of the dynamics of polarization over time. Secondly, we curated and analyzed a novel dataset obtained from a platform with robust moderation dynamics and an extremely loyal user base. Lastly, we have contributed by comparing approximated and exact partitioning methods for signed networks, which can aid future research in this field.

**Supplementary information.** This article has an accompanying supplementary file.

# Declarations

- Conflict of interest/Competing interests (check journal-specific guidelines for which heading to use). NA
- Ethics approval. NA
- Consent to participate. NA
- Consent for publication. NA
- **Data availability.** In order to ensure the reproducibility of these analysis and contribute to the field with our curated dataset, we intend to release the network data for Derstandard once the article is published. Until that time, we are willing to share the data upon request. The individual responsible for inquiries regarding the data is the corresponding author, EF.
- **Code availability.** In order to ensure the reproducibility of these analyses, we intend to make the code publicly available once the article is published. Until that time, we are willing to share the code upon request. The individual responsible for inquiries regarding the code is the corresponding author, EF.
- **Authors contributions.** E.F. carried out the analyses, code generation, visualizations and main draft for the manuscript. M.P retrieved the data, helped with the analysis and provided contextualization for DerStandard. S.S. provided his insights and expertise in the study of polarization and political systems, and embedded our findings in the relevant literature. V.G. supervised and assessed the methodology approaches and provided original ideas for the formal definition of our metrics . D.G. provided a global guiding and supervision through the project and several of

the ideas behind this framework. All authors contributed to writing the manuscript, and took part in the discussions and decisions which resulted in this work.

# References

[1] DellaPosta, D., Shi, Y., Macy, M.: Why do liberals drink lattes? American Journal of Sociology **120**(5), 1473–1511 (2015)

[2] Rokkan, S.: Geography, religion, and social class: Crosscutting cleavages in norwegian politics. Party systems and voter alignments **367**, 379–86 (1967)

[3] Blau, P.M., Schwartz, J.E.: Crosscutting social circles: Testing a macrostructural theory of intergroup relations (1984)

[4] Mason, L.: A cross-cutting calm: How social sorting drives affective polarization. Public Opinion Quarterly **80**(S1), 351–377 (2016)

[5] Finkel, E.J., Bail, C.A., Cikara, M., Ditto, P.H., Iyengar, S., Klar, S., Mason, L., McGrath, M.C., Nyhan, B., Rand, D.G., *et al.*: Political sectarianism in america. Science **370**(6516), 533–536 (2020)

[6] Lipset, S.M., Lipset, S.M., Rokkan, S.: Party Systems and Voter Alignments: Cross-national Perspectives vol. 7. New York: Free Press, ??? (1967)

[7] Franklin, M.N.: The decline of cleavage politics. Electoral change: Responses to evolving social and attitudinal structures in Western countries, 383–405 (1992)

[8] Kriesi, H., Grande, E., Lachat, R., Dolezal, M., Bornschier, S., Frey, T.: West European Politics in the Age of Globalization. Cambridge University Press, ??? (2008)

[9] Ford, R., Jennings, W.: The changing cleavage politics of western europe. Annual review of political science **23**, 295–314 (2020)

[10] Hooghe, L., Marks, G.: Cleavage theory meets europe's crises: Lipset, rokkan, and the transnational cleavage. Journal of European public policy **25**(1), 109–135 (2018)

[11] Bartolini, S., Mair, P.: Identity, Competition and Electoral Availability: the Stabilisation of European Electorates 1885-1985. ECPR Press, ??? (2007)

[12] Goldberg, A.C.: The evolution of cleavage voting in four western countries: Structural, behavioural or political dealignment? European Journal of Political Research **59**(1), 68–90 (2020)

[13] Guerra, P.C., Meira Jr, W., Cardie, C., Kleinberg, R.: A measure of polarization on social media networks based on community boundaries. In: Seventh

International AAAI Conference on Weblogs and Social Media (2013)

[14] Keuchenius, A., Törnberg, P., Uitermark, J.: Why it is important to consider negative ties when studying polarized debates: A signed network analysis of a dutch cultural controversy on twitter. PloS one **16**(8), 0256696 (2021)

[15] Barberá, P., Jost, J.T., Nagler, J., Tucker, J.A., Bonneau, R.: Tweeting from left to right: Is online political communication more than an echo chamber? Psychological science **26**(10), 1531–1542 (2015)

[16] Heider, F.: The Psychology of Interpersonal Relations. John Wiley & Sons Inc, Hoboken (1958). https://doi.org/10.1037/10628-000

[17] Cartwright, D., Harary, F.: Structural balance: a generalization of heider's theory. Psychological review **63**(5), 277 (1956)

[18] Leskovec, J., Huttenlocher, D., Kleinberg, J.: Signed networks in social media. In: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, pp. 1361–1370 (2010)

[19] Estrada, E., Benzi, M.: Walk-based measure of balance in signed networks: Detecting lack of balance in social networks. Physical Review E **90**(4), 042802 (2014)

[20] Aref, S., Wilson, M.C.: Measuring partial balance in signed networks. Journal of Complex Networks **6**(4), 566–595 (2018)

[21] Aref, S., Dinh, L., Rezapour, R., Diesner, J.: Multilevel structural evaluation of signed directed social networks based on balance theory. Scientific reports **10**(1), 1–12 (2020)

[22] Pröllochs, N.: Community-based fact-checking on twitter's birdwatch platform. In: Proceedings of the International AAAI Conference on Web and Social Media, vol. 16, pp. 794–805 (2022)

[23] Allen, J., Martel, C., Rand, D.G.: Birds of a feather don't fact-check each other: Partisanship and the evaluation of news in twitter's birdwatch crowdsourced fact-checking program. In: Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems, pp. 1–19 (2022)

[24] Barberá, P.: Birds of the same feather tweet together: Bayesian ideal point estimation using twitter data. Political analysis **23**(1), 76–91 (2015)

[25] Niederkrotenthaler, T., Laido, Z., Kirchner, S., Braun, M., Metzler, H., Waldhör, T., Strauss, M., Garcia, D., Till, B.: Mental health over nine months during the sars-cov2 pandemic: Representative cross-sectional survey in twelve waves between april and december 2020 in austria. Journal of affective disorders **296**,

49–58 (2022)

[26] Andres, G., Casiraghi, G., Vaccario, G., Schweitzer, F.: Reconstructing signed relations from interaction data. arXiv preprint arXiv:2209.03219 (2022)

[27] García, D., Tanase, D.: Measuring cultural dynamics through the eurovision song contest. Advances in Complex Systems **16**(08), 1350037 (2013)

[28] Neal, Z.: The backbone of bipartite projections: Inferring relationships from co-authorship, co-sponsorship, co-attendance and other co-behaviors. Social Networks **39**, 84–97 (2014)

[29] Tufekci, Z.: Big questions for social media big data: Representativeness, validity and other methodological pitfalls. In: Proceedings of the International AAAI Conference on Web and Social Media, vol. 8, pp. 505–514 (2014)

[30] Doreian, P., Mrvar, A.: Structural balance and signed international relations. Journal of Social Structure **16**(1), 1–49 (2015)

[31] Maoz, Z., Terris, L.G., Kuperman, R.D., Talmud, I.: What is the enemy of my enemy? causes and consequences of imbalanced international relations, 1816–2001. The Journal of Politics **69**(1), 100–115 (2007)

[32] Diaz-Diaz, F., Bartesaghi, P., Estrada, E.: Network theory meets history. local balance in global international relations. arXiv preprint arXiv:2303.03774 (2023)

[33] Estrada, E.: Rethinking structural balance in signed social networks. Discrete Applied Mathematics **268**, 70–90 (2019)

[34] Aref, S., Neal, Z.P.: Identifying hidden coalitions in the us house of representatives by optimally partitioning signed networks based on generalized balance. Scientific reports **11**(1), 1–9 (2021)

[35] Guha, R., Kumar, R., Raghavan, P., Tomkins, A.: Propagation of trust and distrust. In: Proceedings of the 13th International Conference on World Wide Web, pp. 403–412 (2004)

[36] Kunegis, J., Lommatzsch, A., Bauckhage, C.: The slashdot zoo: mining a social network with negative edges. In: Proceedings of the 18th International Conference on World Wide Web, pp. 741–750 (2009)

[37] West, R., Paskov, H.S., Leskovec, J., Potts, C.: Exploiting social network structure for person-to-person sentiment analysis. Transactions of the Association for Computational Linguistics **2**, 297–310 (2014)

[38] Maniu, S., Cautis, B., Abdessalem, T.: Building a signed network from interactions in Wikipedia. In: Databases and Social Networks on - DBSocial '11, pp. 19–24. ACM Press, Athens, Greece (2011). https://doi.org/10.1145/1996413.

1996417

[39] Pougué-Biyong, J., Semenova, V., Matton, A., Han, R., Kim, A., Lambiotte, R., Farmer, D.: Debagreement: A comment-reply dataset for (dis) agreement detection in online debates. In: Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2) (2021)

[40] Pougué-Biyong, J., Gupta, A., Haghighi, A., El-Kishky, A.: Learning stance embeddings from signed social graphs. arXiv preprint arXiv:2201.11675 (2022)

[41] Saeed, M., Traub, N., Nicolas, M., Demartini, G., Papotti, P.: Crowdsourced fact-checking at twitter: How does the crowd compare with experts? In: Proceedings of the 31st ACM International Conference on Information & Knowledge Management, pp. 1736–1746 (2022)

[42] Drolsbach, C.P., Pröllochs, N.: Believability and harmfulness shape the virality of misleading social media posts. In: Proceedings of the ACM Web Conference 2023, pp. 4172–4177 (2023)

[43] Wojcik, S., Hilgard, S., Judd, N., Mocanu, D., Ragain, S., Hunzaker, M., Coleman, K., Baxter, J.: Birdwatch: Crowd wisdom and bridging algorithms can inform understanding and reduce the spread of misinformation. arXiv preprint arXiv:2210.15723 (2022)

[44] Drolsbach, C., Pröllochs, N.: Diffusion of community fact-checked misinformation on twitter. arXiv preprint arXiv:2205.13673 (2022)

[45] Aref, S., Mason, A.J., Wilson, M.C.: A modeling and computational study of the frustration index in signed networks. Networks **75**(1), 95–110 (2020)

[46] Doreian, P., Mrvar, A.: Partitioning signed social networks. Social Networks **31**(1), 1–11 (2009)

[47] Schoch, D.: Signnet: An R Package to Analyze Signed Networks. (2020). https://github.com/schochastics/signnet

[48] Davis, J.A.: Clustering and structural balance in graphs. Human relations **20**(2), 181–187 (1967)

[49] Vijaymeena, M., Kavitha, K.: A survey on similarity measures in text mining. Machine Learning and Applications: An International Journal **3**(2), 19–28 (2016)

[50] Gallagher, R.J., Frank, M.R., Mitchell, L., Schwartz, A.J., Reagan, A.J., Danforth, C.M., Dodds, P.S.: Generalized word shift graphs: a method for visualizing and explaining pairwise comparisons between texts. EPJ Data Science **10**(1), 4 (2021)

[51] Third Position. Page Version ID: 1160606245 (2023). https://en.wikipedia.org/

w/index.php?title=Third_Position&oldid=1160606245 Accessed 2023-07-12

[52] 2015 European migrant crisis. Page Version ID: 1159102024 (2023). https://en.wikipedia.org/w/index.php?title=2015_European_migrant_crisis&oldid=1159102024 Accessed 2023-06-14

[53] 2015–16 New Year's Eve sexual assaults in Germany. Page Version ID: 1159999875 (2023). https://en.wikipedia.org/w/index.php?title=2015%E2%80%9316_New_Year%27s_Eve_sexual_assaults_in_Germany&oldid=1159999875 Accessed 2023-06-14

[54] Garcia, D., Abisheva, A., Schweighofer, S., Serdült, U., Schweitzer, F.: Ideological and temporal components of network polarization in online political participatory media. Policy & internet **7**(1), 46–79 (2015)

[55] Baldassarri, D., Gelman, A.: Partisans without constraint: Political polarization and trends in american public opinion. American Journal of Sociology **114**(2), 408–446 (2008)

# Supplemental Information

## Derstandard Contextualization



**Fig. 1 Example of one DerStandard forum post, showing votes.** Each posting in the forum can be up (green) or down votes by other registered user. The bar on the left top of the posting shows the sum for both types of votes. By clicking on that bar, we open a menu that contains the user names of voters and the type of vote cast. (User names have been blurred by the authors in this example.)

### Context: 2015-2016 changes in Austria

It's striking that between 2015 and 2016 a pronounced change occurs in our network measure of polarization. During that time, media discourse was dominated by the events in Cologne (and other German cities) during the New Year's Eve 15/16 celebrations. A substantial number of women was reporting sexual assault by groups in public, an unusual criminal offense in Germany (for a timeline of events and contextualisation see [**?** ] and Wikipedia page). The political discussion starting with those events led to a pronounced shift in public opinion, summarized by the influential German newspaper "Der Spiegel" as such: "The night brought an end to the sense of euphoria that had accompanied the welcoming of hundreds of thousands of refugees into the country earlier that year".

**Assessment of approximated results**

To assess the robustness of the results, given that we work on a lower bound approximated measure, we run the partitioning algorithm 200 times for each network and keep the lowest value of frustrated edges and its respective partition. We provide three robustness checks: (i) To assess the number of iterations of the algorithm, we check if we can find the optimal solution within $it = \frac{1}{2}200$ iterations, and if not we see how different would the final result change. We find that the number of frustrated edges of the optimal solution found in half of the iterations differs less than 1% with the one found in 200 iterations. (ii) We compare the similarities between the partitions within the 3 best solutions found: in the three cases we find almost-identical sizes for the partitions, with a Szymkiewicz–Simpson overlap coefficient of 0.79. (iii) We provide a comparison analysis between the exact solutions and approximated solutions when using a rolling window found for the BW1 dataset, see Figure 2. These partitions have a Szymkiewicz–Simpson overlap coefficient of 0.8.

Tables 1 and 2 show the detailed resulting partitions from the exact or approximated methods.

2

| | Method | $K^*$ | Ratio Size Groups | Ratio Internal/ External | % Frus Edges | Coh/ Norm. Coh | Coh Resample Uncertainty | Div/ Norm. Div | Div Resample Uncertainty | $SAI_G$ | $SAI_G$ propagated uncertainty |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **BW1** | Exact | 2 | 65/35 | 67/33 | 14% | 0.929 / 0.201 | 0.004 | 0.72 / 0.43 | 0.01 | 0.669 | 0.001 |
| **BW1** | Approx | 2 | 71/29 | 66/34 | 18% | 0.901 / 0.182 | 0.005 | 0.64 / 0.36 | 0.01 | 0.563 | 0.002 |
| **BW2** | Approx. | 2 | 73/27 | 62/38 | 25% | 0.801 / 0.181 | 0.002 | 0.674 / 0.294 | 0.003 | 0.4746 | 0.0007 |

**Table 1 Birdwatch Global Results.** Summary of the optimal partition results for the two datasets. To evaluate the difference in the use of the methods, for BW1, we show results obtained with both the exact and approximated method. $SAI_G$ propagated uncertainty is the propagation of the standard deviation from the null model (10,000 instances). Divisiveness and cohesiveness resample uncertainties are obtained by bootstrapping for 10,000 instances.
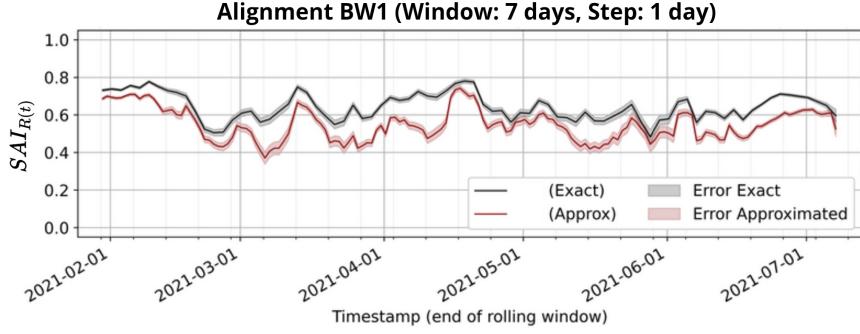
| | Method | $K^*$ | Ratio Size Groups | Ratio Internal/ External | % Frus Edges | Coh/ Norm. Coh | Coh Resample Uncertainty | Div/ Norm. Div | Div Resample Uncertainty | $SAI_G$ | $SAI_G$ propagated uncertainty |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **DS** | Approx. | 2 | 62/37 | 67/33 | 29% | 0.7006 / 0.1409 | 0.0002 | 0.7301 / 0.2899 | 0.0003 | 0.3955 | 0.0002 |

**Table 2 DerStandard Global Results.** Summary of the optimal partition results for the dataset obtained from DerStandard.$SAI_G$ propagated uncertainty is the propagation of the standard deviation from the null model (10,000 instances). Divisiveness and cohesiveness resample uncertainties are obtained by bootstrapping for 10,000 instances.

3

**Fig. 2 Comparison between the timeline results obtained for the approximated and exact methods in the BW1 dataset.** This figure is an analogous of Figure 4 in the main text with different rolling window parameters. It presents the changes in Alignment obtained with the optimal partition of the exact method and the sub-optimal partition obtained through the approximated algorithm with the same data.

## Multipartition study

Below we show the distribution of the partition results for the approximated method. This method is of stochastic nature and therefore we run it several times (i.e. 200 instances) and select the partition that yields the minimum number of frustrated edges. We also use these results to select the optimal number of groups, $k^*$, by selecting $k$ with the best value of frustrated edges. In Figure 3 we show how the trend of results seems to increase with $k$, in agreement with the theorem in [**?** ] which states that the number of minimum frustrated edges is concave when plotted against $k$. We show the results for the Destandard dataset and for the BW2 dataset, which are the datasets that require the use of the approximated method because of their dimensions.

## Metrics normalization

In Fig 4 we show the metrics of divisiveness and cohesiveness before normalization for the time series of BW1. We show both the original data metrics and the null model mean. This figure supports the normalization choice of substracting antagonism (i.e. proportion of negative interactions) from divisiveness to obtain a more meaningful signal on the relevance of sign distribution in a specific time window. Cohesiveness, on the other hand, is perfectly correlated with the proportion of positive interactions and thus should be normalized by substracting such amount.

## Antagonism and Alignment in BW1

In Fig 5 we show the metrics of antagonism and alignment for the time series of BW1. The two time series have a correlation of 0.616. We find this number to be low enough to identify both metrics as different phenomena and thus to emphasize
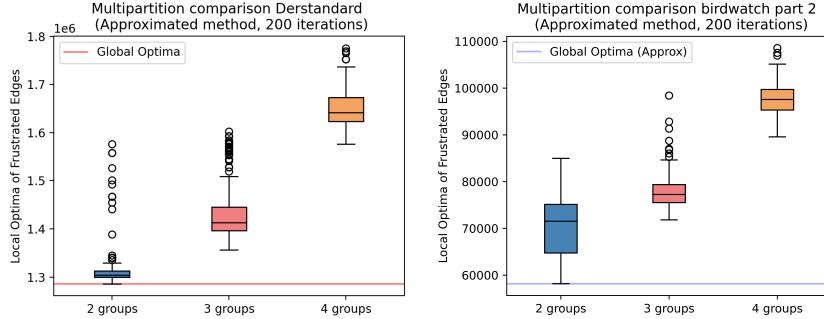
**Fig. 3 Multipartition study for Derstandard (left) and BW2 (right).** We show the distribution of results for the approximated method for $k = 2, 3$ and $4$. In a straight line, we mark the best partition result, which we assume to be the closest to the global optima. All other solutions are sub-optimal and therefore local optima. In both datasets $k = 2$ is the best number of groups.
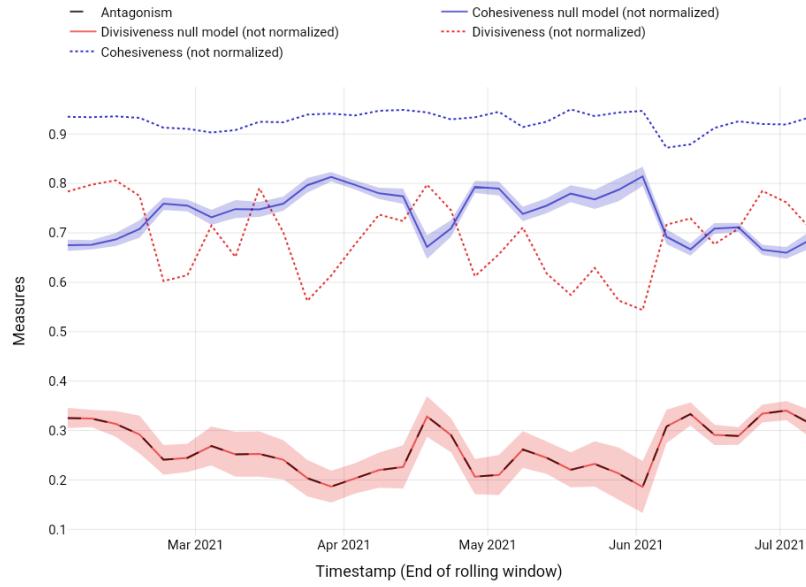


**Fig. 4 Cohesiveness and divisiveness of the original data and null model before normalization for the BW1 time series.** We show the metrics of cohesiveness and divisiveness for the original data (dotted lines) and for the null model. The null model time series is shown with 95% confidence intervals obtained from the $10,000$ instances of re-shuffled sign distributions. The proportion of negative interactions in each time window is represented in a dashed line and is perfectly correlated to the mean divisiveness signal of the null model. Note that it is also inversely correlated to the cohesiveness of the null model.

5

the importance of considering them separately. Moreover, due to the use of a rolling window for the construction of the time series, this correlation measure also contains auto-correlations, and would otherwise be lower.
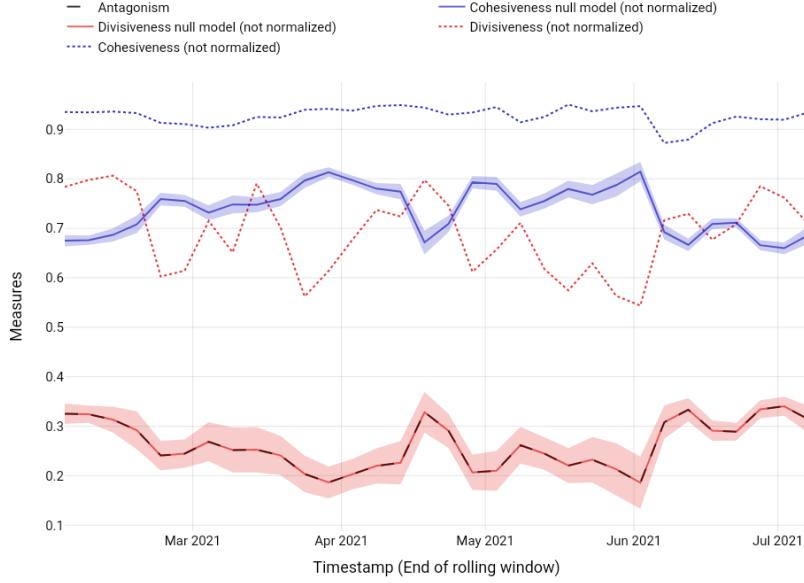


**Fig. 5 Antagonism and alignment of the BW1 time series.** We see that, while fluctuations are similar for both metrics in some time windows, the correlation between the metrics is low enough to consider them as separate measures that provide different insights.

## Size effects

In Figure 6 we show the correlation between our alignment measures for sub-sets of the data and the size of votes (in the case of a temporal rolling window) or votes, posts and articles (in the case of a topic). These coefficients are computed on the data used for the main text figures. As expected, we see there is no direct correlation between the amount of data we consider for each sub-set and the level of Alignment in the network of interactions.

## Birdwatch wordshift graphs

We apply a word shifts method in order to contextualize the topics of discussion surrounding the peaks we detect in the timeline obtained for BW1. Word shifts extract which words contribute to a difference between two texts and how they do so. We use the tool *Shifterator* [? ], which shows the differences in interpretable horizontal bars that compare two texts. Particularly, we use a Frequency-based proportion shift
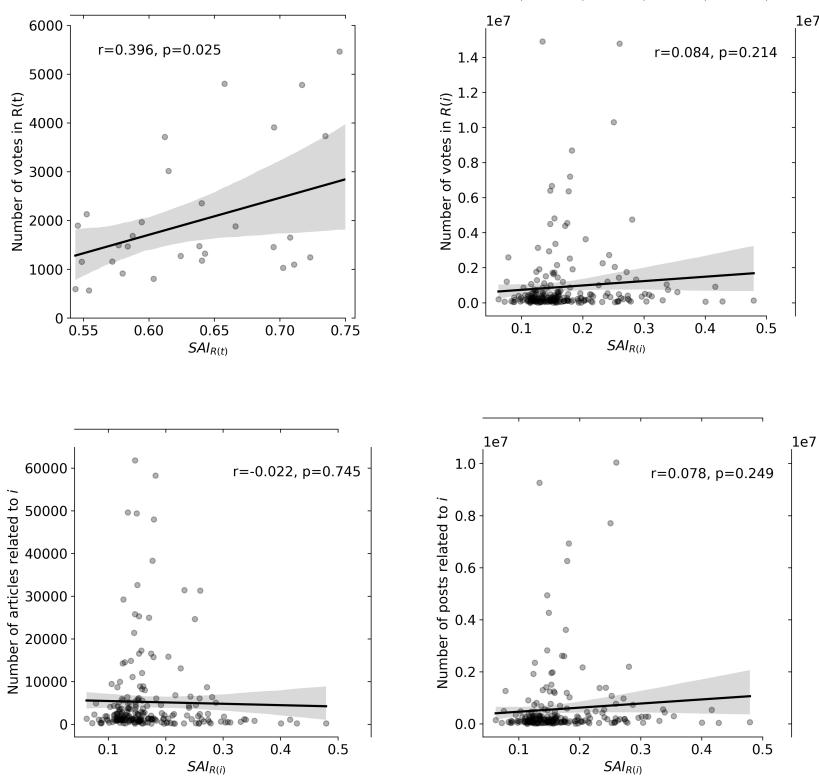
**Fig. 6 Correlation between our Alignment metric and the volume of data of studied subsets.** Scatter plots showing the correlation of the $SAI_R$ measure against the volume of votes for the timeline of BW1 (top left) and the Antagonism-Alignment study for Derstandard (top right). The two lower figures similarly indicate the correlation between the Alginment measures and the volume of articles and posts obtained for each issue for the Derstandard study.

method, that consists on measuring the difference of relative frequencies of a word in each text. In Figures 7, 8 and 9 we show such wordshift graphs for each peak.
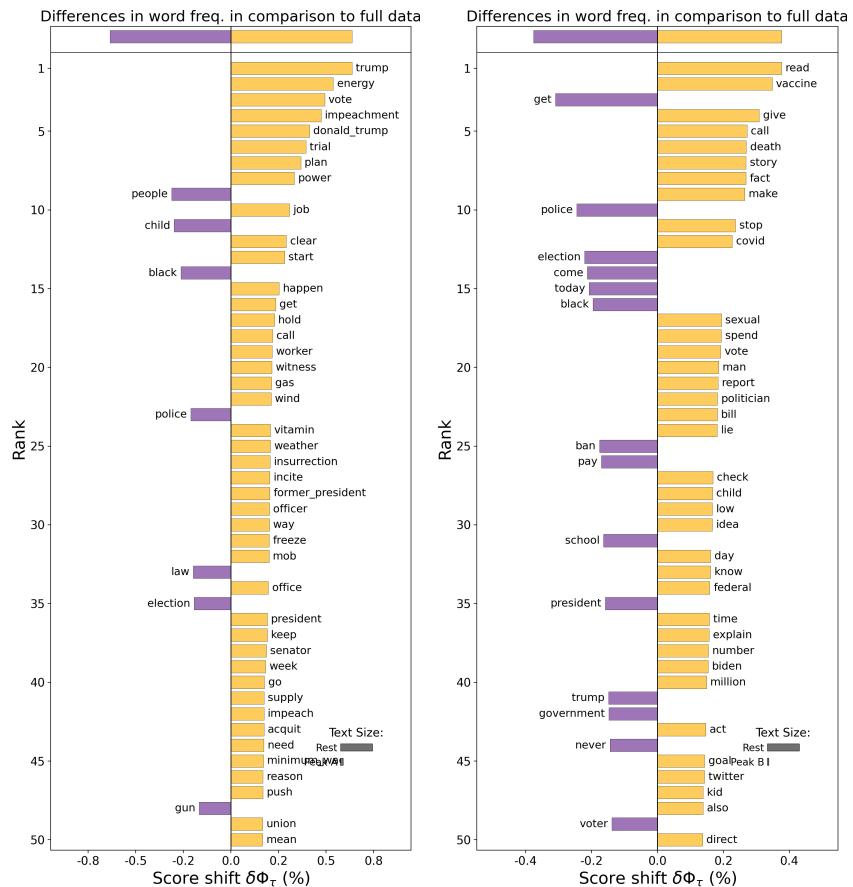
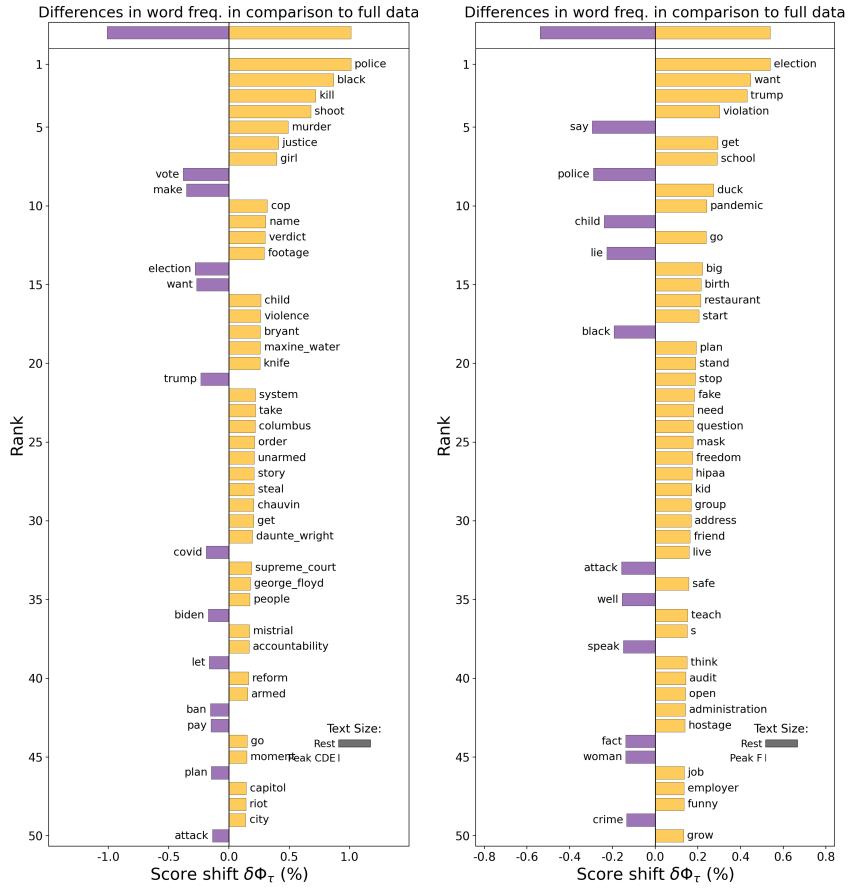**Fig. 7 Wordshift graphs for peak A (left) and peak B (right).**
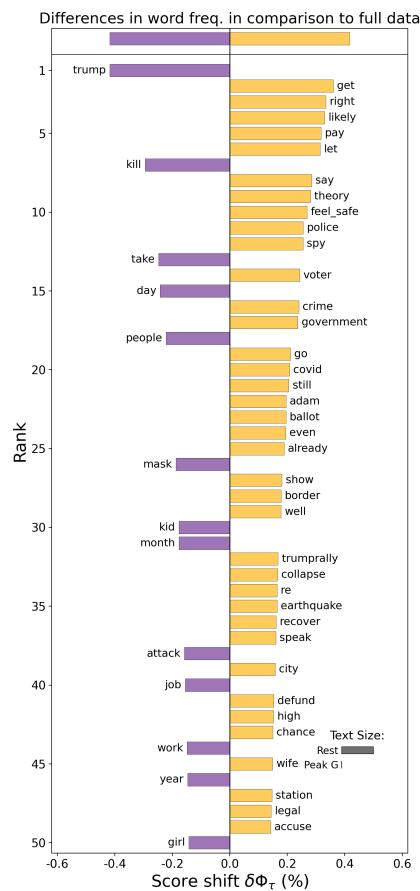
**Fig. 8** Wordshift graphs for peak CDE (left) and peak G (right).

**Fig. 9 Wordshift graph for peak F.**