



OPEN

Bow-tie structures of twitter discursive communities

Mattia Mattei^{1,2}, Manuel Pratelli^{1,3}, Guido Caldarelli^{1,4,5}, Marinella Petrocchi^{1,3} & Fabio Saracco^{1,6,7}✉

Bow-tie structures were introduced to describe the World Wide Web (WWW): in the direct network in which the nodes are the websites and the edges are the hyperlinks connecting them, the greatest number of nodes takes part to a *bow-tie*, i.e. a Weakly Connected Component (WCC) composed of 3 main sectors: IN, OUT and SCC. SCC is the main Strongly Connected Component of WCC, i.e. the greatest subgraph in which each node is reachable by any other one. The IN and OUT sectors are the set of nodes not included in SCC that, respectively, can access and are accessible to nodes in SCC. In the WWW, the greatest part of the websites can be found in the SCC, while the search engines belong to IN and the authorities, as Wikipedia, are in OUT. In the analysis of Twitter debate, the recent literature focused on discursive communities, i.e. clusters of accounts interacting among themselves via retweets. In the present work, we studied discursive communities in 8 different thematic Twitter datasets in various languages. Surprisingly, we observed that almost all discursive communities therein display a bow-tie structure during political or societal debates. Instead, they are absent when the argument of the discussion is different as sport events, as in the case of Euro2020 Turkish and Italian datasets. We furthermore analysed the quality of the content created in the various sectors of the different discursive communities, using the domain annotation from the fact-checking website Newsguard: we observe that, when the discursive community is affected by m/disinformation, the content with the lowest quality is the one produced and shared in SCC and, in particular, a strong incidence of low- or non-reputable messages is present in the flow of retweets between the SCC and the OUT sectors. In this sense, in discursive communities affected by m/disinformation, the greatest part of the accounts has access to a great variety of contents, but whose quality is, in general, quite low; such a situation perfectly describes the phenomenon of infodemic, i.e. the access to “*an excessive amount of information about a problem, which makes it difficult to identify a solution*”, according to WHO.

Since their first introduction, Online Social Networks (OSN) have been deeply investigated for possible implications of the online public debate on political processes¹. In the last decade, the centrality of OSN for political communications and debates has steady increased: OSN represent one of the most used tool for citizens to get an opinion². It is not surprising, then, that political parties use them extensively to carry out a sort of never-ending propaganda.

Although in the literature there are different opinions on the impact that a particular grouping of users in OSN can have on their offline behavior^{3–5}, it is undeniable that the online social environment is strongly polarized. The origin of such a polarization has been deeply discussed in the sociological literature^{6,7} and seems to be extremely dependent on country's party systems⁸. In particular, the concepts of *selective exposure*, *confirmation bias*, *echo chambers* and *filter bubbles* have had a great relevance in the literature.

Selective exposure leads people to prefer information that confirms their preexisting beliefs^{9,10}, while confirmation bias makes information consistent with one's preexisting beliefs more persuasive¹¹.

Such phenomena imply the formation of groups of users, characterised by following the same information in terms, e.g., of news outlets and personal opinions. These groups are thus closed in so called *echo chambers*: “a bounded, enclosed media space that has the potential to both magnify the messages delivered within it and

¹IMT School For Advanced Studies Lucca, p.zza San Francesco 19, 55100 Lucca, Italy. ²Alephsys Lab, Universitat Rovira i Virgili, Av. Paisos Catalans 26, 43007 Tarragona, Catalonia, Spain. ³Institute of Informatics and Telematics, National Research Council, via Moruzzi 1, 56124 Pisa, Italy. ⁴Department of Molecular Sciences and Nanosystems, Ca' Foscari University of Venice, Ed. Alfa, Via Torino 155, 30170 Venezia Mestre, Italy. ⁵European Centre for Living Technology (ECLT), Ca' Bottacin, 3911 Dorsoduro Calle Crosera, 30123 Venice, Italy. ⁶Institute for Applied Mathematics “Mauro Picone”, National Research Council, via dei Taurini 19, 00185 Rome, Italy. ⁷“Enrico Fermi” Research Center, via Panisperna 89 A, 00184 Rome, Italy. ✉email: fabio.saracco@cref.it

insulate them from rebuttal”^{11–13}. Echo chambers, by being impervious to information coming from outside that may contradict the pre-existing views of the chamber members, are believed to strongly contribute to the polarization of the online debate¹⁴.

Polarization may be also fomented by *filter bubbles*. This paradigm was first introduced by the activist Eli Pariser¹⁵: personalised results provided by search engines and shown in social media feeds can make users be trapped in a bubble of information they like and away from data and viewpoints considered less valuable, but that could challenge their beliefs. Although the user may not be affected in real life by the virtual bubble, due to the various communication channels he or she can take advantage of (see Ref.¹⁶), the customization of algorithms may contribute to the formation of a virtual world apart.

Discourse and discursive communities. Whether circulating within an echo chamber or suggested by recommendation algorithms, the type of information users come across online is fundamental to reinforcing or not the division into ‘closed’ groups. Nevertheless, also the study of the interactions between users is of absolute interest to detect polarization phenomena. The term *discourse community* was coined in 1982 and it indicates ‘groups that have goals or purposes, and use communication to achieve these goals’¹⁷. A discourse community is itself immaterial, and this tends to project it onto the forum on which it operates¹⁸. Thus, with the advent of OSN, discourse communities were projected onto the platforms themselves¹⁹: ‘A discourse community can be viewed as a social network, built from participants who share some set of communicative purposes’. According to Berkenkotter²⁰, ‘just as the digital world is constantly evolving, discourse communities continually define and redefine themselves through communications among members’.

In the discourse community definition, we implicitly know the identities of the individuals forming the community. Actually, in the case of Twitter, it is just partially true, since we have trustworthy information only about a small minority of accounts. For this reason, we prefer to use the term *discursive communities*, as it was introduced in Ref.²¹ to identify group of users that are connected by non-trivial pattern of discourse, but for which we have limited information about the identity of the group itself. Nevertheless, since we can *infer* the discourse community of the discursive community by looking at a set of non-trivial data characterising the group, as the most frequent keywords used therein, the difference is more formal than substantial. Therefore, in the following we will use the two terms interchangeably.

To detect discursive communities in OSN, the first contributions applied mixed approaches to the political debate on Twitter^{22–24}. The work considered political debate on Twitter about the US presidential election campaign, i.e. a ‘perfectly polarized’ one in which two opposite fronts face each other. The authors manually annotated the most frequent keywords characterizing Republicans and Democrats’ narratives and use them to infer the political orientation of accounts using them. The orientation of accounts not using hashtags was later inferred using a label propagation algorithm²⁵. Remarkable, a clear partition in two distinct groups of users, supporters of the two political parties, was observed in the *retweet network* only (the network of users sharing content created by others). Finally, using a label propagation algorithm on the retweet network, the authors were able to successfully assign to all accounts the proper political orientation, that can be translated in the present context to the correct discursive community.

Every country has a different way in which opinions are polarised. This is due to the various party systems and electoral laws and, in principle, there could be more than just two fronts⁸. A methodology for detecting discourse communities less susceptible to human error should therefore rise from the data directly, rather than being based on a priori manual annotation. The approach firstly proposed in Ref.²⁶ meets the desired property: the idea is to infer the various discursive communities starting from accounts whose identity is certified by the social network itself. In Twitter, these are the so called *verified* accounts. This class of accounts tend to produce new content rather than retweet the one created by others²⁶. Since we trust information regarding their identities, the issue is the identification of the discursive communities anchored to them, something that can be done using their interactions with ‘standard’ users (based on the results by Conover et al.^{22–24}, in terms of retweets).

Let the reader consider a pair of verified Twitter users. If they share a large number of retweeters, it is reasonable to think that they attract ‘similar’ users. In this sense, that pair of accounts are perceived to belong to the same discursive community, sharing similar views and ideas. Nevertheless, it is hard to state *a priori* how many common retweeters two verified users should have in order to be considered as ‘similar’; in this sense, a maximum entropy null-model is used as benchmark²⁷. (More technical details on this construction can be found in the “**Discursive communities**” section.) The labels for verified users are then propagated, following the same approach as in Refs.^{22–24}. In the present paper, we are going to follow this approach, that has great performances on manually annotated datasets²⁸.

The recent literature has extensively analysed online debate and discourse communities, focusing, from time to time, on coordinated activities in discursive communities^{26,29,30}, on the semantic network associated to the various discursive communities^{21,31,32}, on their exposure to disinformation campaigns^{33,34}, and on their dynamical evolution^{35–39}.

In the present paper, we tackle the analysis of the network structure of discursive communities: we collect and study 8 thematic Twitter datasets, on topics ranging from sports, to COVID-19, to political elections and immigration policies. Our main result is that *almost all the discourse communities therein features a bow-tie structure*.

Bow-ties. Bow-tie structures were initially introduced by Broder et al. in order to study the structure of World Wide Web (WWW)⁴⁰. The authors represented WWW as a directed network in which webpages are the nodes and the hyperlinks connecting them are the edges. Broder et al. noticed that the network displays a huge Weakly Connected Component (WCC), i.e., the maximal subgraph in which all nodes can be reached by any

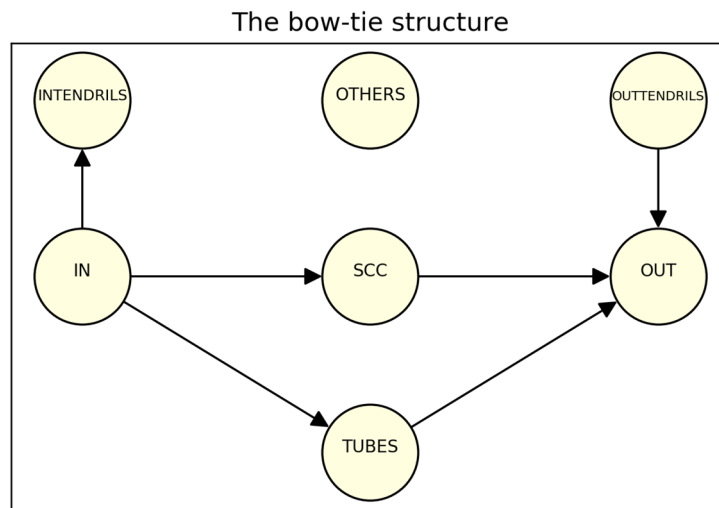


Figure 1. The seven sectors of Yang's bow-tie structure.

other one in the same subgraph, disregarding the direction of the link. This WCC includes more than 75% of all nodes.

WCC breaks into three main pieces: a Strongly Connected Component (SCC), in which each node can be reached by any other one in the same block, following the direction of the links; a group of nodes that can reach SCC, without being reached by it (called IN); a group of nodes that can be reached by SCC, but that cannot reach it (the OUT block). The observation is that SCC is the most populated sector, followed by the IN and the OUT sectors. Most of the websites can be found in the SCC, linking between each other; the IN sector was instead mostly composed by search engines, while the OUT one includes authorities, as Wikipedia.

Yang et al.⁴¹ refined the partition of the structure in⁴⁰, introducing INTENDRILS, OUTTENDRILS, TUBES and OTHERS. The entire situation is pictorially represented in Fig. 1.

Remarkably, the bow-tie structure was detected also in control network of transnational corporations, having deep implications on financial stability⁴².

Results in a nutshell. In the case of our 8 thematic datasets, we find that a bow-tie structure is present in those discourse communities debating (1) about politics, like in the case, e.g., of election campaigns, and (2) about society, e.g., on the proper response to the pandemic or the appropriate management of migration fluxes. Instead, bow-ties are hardly present when the debate is about less socially relevant topics as sports (this confirms what observed in Ref.⁸).

There are two relevant points in observing the presence of bow-tie structures in discursive communities: how big the bow-tie is respect to the entire discursive community (the greatest the accounts in the non-OTHERS blocks, the more informative the bow-tie structure is) and how random the presence of this structure is (i.e., its statistical significance). Regarding the first point, when the bow-tie is informative, even in the worst case, it represents more than 80% of all nodes in the discursive community. Regarding the second point, in order to be sure that the observed bow-ties are not due to a random organization of links only, we compare the observed quantities with a maximum entropy null-model for directed network, conserving the in- and out-degree sequences⁴³. The results show that the dimension of most of the bow-tie sectors are statistically significant, i.e., they carry a signal that cannot be due to the degree sequence only. In this sense, the presence of a bow-tie structure is an extremely non-trivial feature of the system.

We can add more detail to the analysis of this structure. When the bow-tie is informative, we observe two cases: the OUT-dominant and the INTEND-dominant ones, depending on which sector is the largest (respectively, OUT and INTENDRILS). The OUT sector has access to all information produced in the discursive community and, in particular, to the one produced by the most active block, SCC. Instead, in the INTEND-dominant bow-ties, the most crowded sector is the one of INTENDRILS, i.e., the retweeters of IN that are not retweeted by anyone else and that cannot access to all content created by SCC.

In principle, it should be desirable to have an OUT-dominant bow-tie: when the OUT sector is the most populated, there are many accounts that are exposed to information from all other sectors. This should give the accounts a multi-faceted, pluralistic knowledge. However, for the investigated datasets, we carry out an analysis on the production and quality of content in the various sectors of the bow-ties, and the outcome returns a different picture. In fact, regarding content production, SCC is the source of the greatest flux of content. When the discursive community is affected by m/disinformation, the incidence of links from non-reliable sources shared by SCC is much greater than by any other sector and is particularly considerable in the flux between SCC and OUT. In those cases when OUT-sectors dominate the bow-tie, we observe an *infodemic* (According to WHO, “*infodemics are an excessive amount of information about a problem, which makes it difficult to identify a solution*”. Coronavirus disease 2019 (COVID-19) Situation Report—45: <https://www.who.int/docs/default-source/>

[coronaviruse/situation-reports/20200305-sitrep-45-covid-19.pdf?sfvrsn=ed2ba78b_4](https://www.nature.com/scientificreports/20200305-sitrep-45-covid-19.pdf?sfvrsn=ed2ba78b_4)), since OUT, i.e. the widest block, is directly hit by the huge amount of messages of questionable quality produced by SCC.

Summarising, our contribution is twofold:

- almost all the discursive communities in the 8 investigated datasets of Twitter debates, on different topics in different countries, display a bow-tie structure which is statistically significant;
- when the bow-tie is affected by m/disinformation and it is OUT-dominant, the majority of users (i.e. those in the OUT block) is exposed to the flux of m/disinformation. In this sense, the bow-tie structure fuels the phenomenon of infodemic.

We would like to remark that the results in this manuscript do not represent the only contribution that connects the diffusion of m/disinformation to the network structure (see, for instance, work in^{44–46}, just to consider some of the most recent contributions). However, to the best of our knowledge, for the first time the bow-tie structure emerges in online discursive communities. Moreover, its presence and its peculiarity permit do have a proper description of the phenomenon of infodemic.

Results

Datasets. In order to make our analysis as general as possible, we consider several Twitter datasets across different countries and about different topics. The data collected using the Twitter Streaming API are publicly available for further research and reproducibility and can be found at the following URL: <https://toffee.imtlucca.it/datasets>. In detail:

- **COVID-19 datasets:** we explore Twitter posts containing keywords related to the COVID-19 pandemic in different languages and therefore diffused in different countries (in particular, the keywords for tweets collection were “coronavirus”, “ncov”, “covid”, “SARS-CoV2”, “#coronavirus”, “#coronaviruses”, “#WuhanCoronavirus”, “#CoronavirusOutbreak”, “#coronaviruschina”, “#coronaviruswuhan”, “#ChinaCoronaVirus”, “#nCoV”, “#ChinaWuHan”, “#nCoV2020”, “#nCov2019”, “#covid2019”, “#covid-19”, “#SARS_CoV_2”, “#SARSCoV2”, “#COVID19”. The subset of Italian messages has been matter of investigation in Ref.³³ too). In particular, we consider the **Italian**, **German** and **French** debates about the pandemic, in the period between February and April 2020. The Italian dataset consists of 4,470,648 tweets published between February 17 and April 23. The German dataset contains 1,552,106 tweets posted between February 10 and April 23, the French one has 3,052,708 posts published between March 23 and April 7. The different time frames for data collection have been chosen according to the intensity of the Twitter traffic.
- **Dutch elections dataset:** we collect Twitter posts about the national elections in the Netherlands in 2021. The keywords used for downloading data were “tweedekamer”, “verkiezingen”, “kabinet”, “coalitie”, “stem”, “stembus”, “verkiezingen2021” (respectively, “House of representatives”, “reconnaissance”, “cabinet”, “coalition”, “vote”, “ballot box”, “explorations”) and only messages in Dutch were selected. The dataset contains 1,002,499 tweets posted between February 2 and March 31, 2021.
- **Italian debate on migrants:** we select Twitter posts shared in Italy with keywords regarding the discussion about the migration flows from Northern Africa to the Italian coasts. The dataset consists in 1,081,780 posts, published between January 23, 2019 and February 22, 2019. The dataset is described in more details in Ref.²⁹.
- **Italian debate on the Astrazeneca vaccine:** we examine 583,236 Twitter posts published in Italian, regarding the discussion about the safety of the Astrazeneca vaccine against COVID-19: the keywords used for the download were “astrazeneca”, “aifa”, “ema”, “trombosi” (respectively, “astrazeneca”, “Italian Medicines Agency”, “European Medicines Agency” and “thrombosis”). The dataset contains posts shared between March 15, 2021 and May 15, 2021.
- **Italian and Turkish EURO2020:** we analyze 144,725 Italian tweets and 430,374 Turkish ones about the European Football Championship EURO2020; the keyword used for the download was simply “#euro2020”. The tweets were published between, respectively, June 11–13 and June 11–23, 2021.

So as not to burden the presentation, in the following we will present the results about the Italian COVID-19, Italian EURO2020 and Turkish EURO2020 datasets. We will show the results related to the other datasets wherever there will be something substantially different, compared with the Italian dataset. However, all graphics and results about the other datasets can be found in the “Supplementary Material”.

Discursive communities. Our analysis focuses on the structure of networks of retweets, for each dataset. Retweeting a post is one of the possible ways in which people can interact on Twitter and it consists in sharing the content of a tweet written by another user. It usually means endorsing the post content as it has also the effect of raising the visibility of the original post. It was also shown that, among all possible interactions, retweets are the best performing to infer the political orientation of the various accounts^{22–24}.

We start by distinguishing between *verified* and *non-verified* accounts. The former ones denote Twitter users whose identity has been verified by the social platform. This procedure is usually adopted to certify the accounts of renowned people and organizations and figures of public interest in general, as politicians, journalists, political parties, newspapers and TV-channels. We represent the interactions between verified and non-verified users as a bipartite network. In a bipartite network, nodes belong to two different sets, called layers, and an edge can exist only between vertices placed on different layers; we place the verified accounts on one layer of a bipartite network and the non-verified ones on the other one, again considering links as retweets between them. (In the present

construction, we disregard the information about the direction of the retweet, since we are interested in the interaction between the two class of users. Nevertheless, as mentioned above, verified users tend more to create new contents (i.e. tweet) than to share it with her followers (retweet)). The main idea is to anchor the definition of discursive communities on verified users since they usually introduce new content and posts: as observed in many other studies^{21,26,29,32–34,47}, verified users are, on average, much more retweeted than common users. Such a procedure obtains great performances, since it can be observed that the various discursive communities are coherent in terms of verified users belonging to the same political front; in a further analysis we are comparing this procedure with annotated datasets, better quantifying our performances²⁸.

Following the methodology introduced in Becatti et al.²⁶, we count the common neighbors of each pair of verified users or, in simpler words, the number of non-verified users that have interacted (by retweeting or being retweeted) with the same pair of verified ones. The aim is projecting the bipartite network into the layer of the verified accounts, establishing an edge between two of them if the number of their common neighbors is significantly higher than what expected by a proper null-model. When this happens, we can assert that the two verified users refer to the same audience and, therefore, they probably share similar content and opinions. The statistical significance of the number of common neighbors can be established only comparing it with the predictions of an accurate benchmark, which, in this case, is represented by the Bipartite Configuration Model (*BiCM*⁴⁸), an entropy-based model suited for bipartite networks. A complete description of the model and the projecting procedure can be found in “[Entropy-based null-models for network analysis and their applications](#)”.

The result of the above procedure is a monopartite network of verified users. We further obtain a partition in communities implementing the Louvain algorithm⁴⁹ for the optimization of the modularity, with a slight modification. In fact, the standard definition of the modularity⁵⁰ implements the Chung-Lu null-model⁵¹, which can be considered as a sparse matrix approximation of the entropy-based null-model defined in⁵² and it is known to return wrong results in the presence of strong hubs²⁷. We thus replaced the Chung-Lu null-model in the modularity with the unipartite configuration model (*UCM*) defined in Ref.⁵². Furthermore, we correct for the node ordering bias that affects Louvain algorithm, independently on the objective function chosen. In fact, we perform multiple runs, each time reshuffling the order of the nodes: we finally select the partition displaying the greatest value of the (*UCM*-modified) modularity. More details can be found in “[Entropy-based null-models for network analysis and their applications](#)”.

For all the datasets, looking at the members of each discursive community, we can *a posteriori* associate the latter to a political wing, using the available information for verified users. We thus obtain clusters of users which represent the main wings of the political scenario of each of the examined countries. In addition, in almost all the datasets, we identify also a Media cluster, with official accounts of newspapers, TV-channels, radio and other media.

In the “[Discursive communities for the Italian COVID-19 dataset](#)”, the interested reader can find a complete description of all the discursive communities for the Italian COVID-19 dataset. For the other datasets, a brief description of their discursive communities is in the “Supplementary Material”.

Political orientation of non-verified users. The next step in our procedure consists in extending the discursive communities to non-verified accounts. More in details, following the approach in Ref.²⁹, we use the membership of verified users as (fixed) seeds for the label propagation algorithm proposed by Raghavan et al.²⁵ on the retweet network. This network is a monopartite and directed one in which nodes represent users and links start from the retweeted users and are directed towards the one who retweets. Let us remind that, in case the algorithm cannot find a dominant label for a specific vertex (i.e., in case of a tie), it randomly removes some of the edges attached to that vertex and repeats the procedure: for this reason, we run the label propagation 500 times and assign to each node the most frequent label (actually, the noise in the assignment of the labels is extremely limited).

Figure 2 shows the percentages of nodes placed in the various discursive communities for the Italian COVID-19 dataset (a detailed description of the various communities can be found in the caption of the figure). Considering also the other datasets, in almost all the cases, the label propagation procedure could assign a label to approximately 90% of the nodes. As we could expect, in the COVID-19 datasets, the Media community is always the most numerous one: updates on the spread of the pandemic, written by the official accounts of various media, received a great amount of retweets.

As highlighted in other works^{21,26,29,30,32–34}, the presence of well-defined discursive communities is the signal that users on Online Social Networks (OSNs) are strongly polarized, i.e., they tend to split into groups, which one with same opinions and political orientation.

The bow-tie structure. The original concept of bow-tie by Broder et al.⁴⁰ sees WWW divided into 3 main sectors: a Strongly Connected Component (SCC), in which each node can be reached by any other one in the same block, following the direction of the links; a group of nodes that can reach SCC, without being reached by it (called IN); a group of nodes that can be reached by SCC, but that cannot reach it (the OUT block).

The description by Broder et al. was subsequently refined by Yang et al.⁴¹, who split the network in seven distinct parts:

- the greatest Strongly Connected Component (SCC);
- the IN block;
- the OUT block;
- the TUBES sector, including nodes reachable from IN and accessing OUT, but not being part of SCC;
- the INTENDRILS group, collecting all those nodes pointed by IN that cannot reach the OUT block;
- the OUTTENDRILS sector, containing all those nodes pointing to OUT that cannot reach nodes in IN;

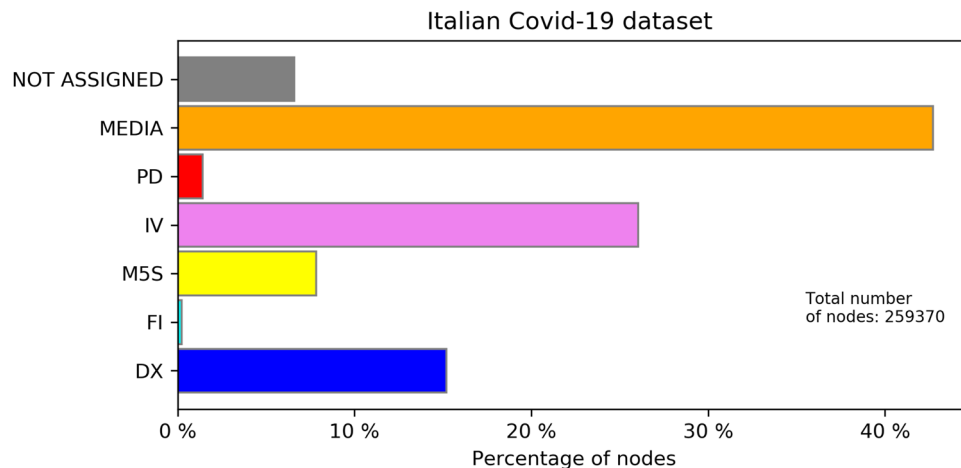


Figure 2. Percentages of nodes in each discursive community, Italian COVID-19 dataset. Due to the presence of politicians and political parties from a specific political area, the various discursive communities are called following their political alignment. “PD” stays for the Italian Democratic Party (*Partito Democratico*); *Italia Viva* (“IV”) is the political party of the former prime Minister and former PD secretary Matteo Renzi, while M5S is the “Movimento 5 Stelle”, a political movement born on the web and being the most represented party in the Italian parliament at the time of the data collection. “FI” stays for *Forza Italia*, the political party of the former Prime Minister Silvio Berlusconi, while the “DX” (*Destra*) community includes right wing parties as Lega and Fratelli d’Italia. The most crowded discursive community is the one of Media in which there are most of the online news outcasts and newspapers. The accounts for which it was not possible to assign a discursive community are in grey.

- the **OTHERS** group, including all those nodes that cannot be placed in one of the previous six sectors.

In Fig. 1 there is a schematic representation of the bow-tie structure defined in Ref.⁴¹. The seven groups of nodes are mutually disjointed.

We remark that every directed network can be divided in blocks using the bow-tie decomposition. Nevertheless, as a rule of thumb, the bow-tie representation is *informative* about the network structure if the number of nodes in blocks other than OTHERS is greater or of the same order of those in OTHERS: the greatest the impact of the non-OTHERS blocks, the more informative the bow-tie structure is.

The bow-tie structure of the discursive communities. In the present manuscript, we investigate the presence of a bow-tie structure in the discursive communities of the retweet network, i.e., in the network composed by Twitter accounts (the nodes) and retweets (the links connecting the original author to the retweeter).

Results show that, when considering political online debates, a bow-tie structure is informative in almost every discursive community of our datasets, while for non-political debates (as the case of Euro2020), the bow-tie structure is less informative. Euro2020 itself records the extreme case in which more than one half of the nodes are in the OTHERS sector. We state that this bow-tie structure is *uninformative*—see, for example, the case of the Turkish debate during Euro2020 in Fig. 8. We remark that the presence of informative bow-ties in many of the discursive communities here investigated is not a trivial result. Indeed, there are no evident reasons for expecting such distribution of the nodes *a priori*.

When a bow-tie structure is informative, we observe two recurrent situations in the investigated datasets and, according to them, we classify the bow-tie into two different categories:

- When the OTHERS block is smaller than SCC, we will refer to **strong** bow-tie structures;
- When the OTHERS block is greater than SCC, we will refer to **weak** bow-tie structures.

Furthermore, when the bow-tie is informative, may it be weak or strong, we can categorize it in two different ways, that we called respectively **OUT-dominant** and **INTEND-dominant**. In OUT-dominant bow-ties, most of the nodes of the bow-tie are placed in the OUT sector. As a rule of thumb, OUT-dominant bow-ties are more frequent when the bow-tie is strong, but we can find some counter-examples. The INTEND-dominant bow-tie is a bow-tie structure in which instead the most part of the nodes is located in the INTENDRILS sector, i.e., when most part of the users retweets accounts from the IN zone and has little to no interaction with the users in the other sectors. INTEND-dominant bow-ties are in general more frequent in weak bow-ties.

We highlight that it is not so strange that the most crowded blocks in the bow-ties are OUT and INTENDRILS: it was already observed in Ref.⁵³ that the greatest number of users tend to mostly retweet content created by others and limit their production of new messages. The difference between OUT-dominant and INTEND-dominant bow-ties is the *access to information*: OUT-dominant bow-ties are those in which the majority of users

can access almost all messages exchanged over the discursive community, while in the INTEND-dominant ones the majority of users limits their retweets to the content produced by accounts in the IN block. Otherwise stated, the main difference between INTEND- and OUT-dominant bow-tie structures is that the former displays a more ‘hierarchical’ structure, i.e., few accounts (those in the IN sector) introduce new content and many others just share it (the INTENDRILS sector). Instead, in OUT-dominant bow-ties, the greatest part of the users (i.e., the OUT block) not only shares posts by accounts in the IN block, but also it retweets content by users in SCC, OUT-TENDRILS and TUBES. We argue that this behaviour, while more ‘democratic’, is, at the same time, more risky.

In fact, we will see in “[Verified users’ distribution](#)” that users with high visibility and which introduce new content on Twitter can be found mostly in the IN sector: typically, they are verified accounts. As observed in other studies, see, e.g. Ref.³³, verified users tend to limit the spreading of low-quality content. We may argue, then, that users interacting mostly with verified users are safer from m/disinformation campaigns. In the following, we will see that the reputability of information shared confirms our hypothesis and we will come back on the matter.

Figure 3 displays the bow-tie structure of each discursive community for the Italian COVID-19 dataset (analogous plots for the other datasets can be found in the “Supplementary Material”). A single node represents one bow-tie sector and its dimension is proportional to the number of accounts in it. First, according to the definitions given above, the bow-tie structure is informative in all the discursive communities. In the cases of DX and IV, the bow-tie is particularly informative: its blocks include respectively 96.5% and 98.3% of the entire discursive community. Second, different discursive communities display bow-ties with different strengths. For instance, DX and IV discursive communities display strong bow-ties, while, M5S, Media, PD and FI have weak ones, since their SCCs are relatively small (and smaller than OTHERS).

Third, the graph shows that the DX, IV, MEDIA and FI communities display OUT-dominant bow-ties, in which the OUT sector is the biggest one; considering all the investigated datasets, OUT-dominant bow-ties represent the most frequent configuration, being 21 out of 31 communities. Instead, 6 out of 31 discursive communities are INTEND-dominant bow-ties (as PD and M5S in Fig. 3).

We remark that, in all our datasets, all the right wing discursive communities display bow-ties with an OUT-dominant structure; in most of the cases, these bow-ties are also strong. The colours of the nodes in Fig. 3 are going to be explained in the following section.

Statistical significance of the bow-tie structure. It may be argued that the bow-tie structures featured by the discursive communities in our datasets are just an accident, due to the different role of the various users in the debate. In fact, those accounts that have high out-degrees and low in-degrees are naturally in the IN sector; those that, viceversa, have high in-degrees and low out-degrees are in the OUT sector, and so on. To test whether the presence of bow-ties is merely attributable to the behavioral characteristics of the accounts, we compare the dimensions of the different sectors, as observed in the real network, with those in a randomised system in which the in- and out-degree sequences are fixed. If the partition in the various bow-tie sectors were just a matter of the degree sequence, none of the dimensions of the various blocks should be statistically significant. Otherwise, we should observe a significant mismatch with respect to the expectation of the null-model.

In order to have an unbiased benchmark, we build an entropy-based null-model that preserves the in- and out-degree sequences, being maximally random for all the rest (see Ref.²⁷ for a review on the subject). Summarising, starting from a real network, we consider the set of all possible graph realizations (the graph *ensemble*) having the same number of nodes as in the real system. Then, we assign to each representative of the ensemble a different probability of realization by maximising the entropy of the ensemble, but constraining the average value of some topological property of the real network (in our case, the in- and out-degree sequences). In this way, even if the single realization of the ensemble does not display the network properties that we would like to preserve, the entire ensemble, on average, does.

In the last years, such procedure has been adopted to analyse financial and economic systems^{43,48,52,54–70}, biological networks^{71–74} and online social networks^{21,26,29–34,75} and it was shown to be effective to extract the relevant structure from a real network^{76,77}.

Here, we implement the Direct Configuration Model (DCM), firstly introduced in Ref.⁴³ and implemented in the Python module [NEMtropy](#)⁷⁰. More details on the exact derivation of DCM can be found in “[Entropy-based null-models for network analysis and their applications](#)”.

Going back to Fig. 3, the colour of the circles indicates the agreement between the actual size of the bow-tie sectors and the size predicted by the DCM: we are interested in detecting both too “big” and too “small” blocks. In particular, the darker the colour of the sectors in Fig. 3, the larger the $-\log_{10}(\text{p-value})$ (so the lower the p-value) and the greater the disagreement of the real system from the randomization. For each sector, the two-tailed p-value has been calculated looking to a sample of 1000 graphs generated by the DCM.

The p-value tells us about the existence of a disagreement, but not about the direction of the disagreement. For instance, looking at the DX bow-tie in Fig. 3, both the dimensions of OTHERS and SCC have a really small p-value, thus they do not agree with the randomization, but the OTHERS block is smaller than predicted by the DCM, while SCC is larger.

Table 1 reports the exact p-values of the different blocks for the various bow-ties of Fig. 3. The significance of the blocks for each bow-tie can be assessed by using the False Discovery Rate (FDR) correction⁷⁸, setting the statistical significance level to $\alpha = 0.01$. In the present case the correction is limited, due to the small number of blocks in the bow-tie.

It is interesting to observe that, in both strong and weak bow-ties, the OTHERS block is statistically significant in all the discursive communities but PD. In particular, the dimension of the OTHERS block is much smaller than predicted by the null-model and the presence of the bow-tie is not due to the degree sequence only.

ITALIAN COVID-19 DATASET

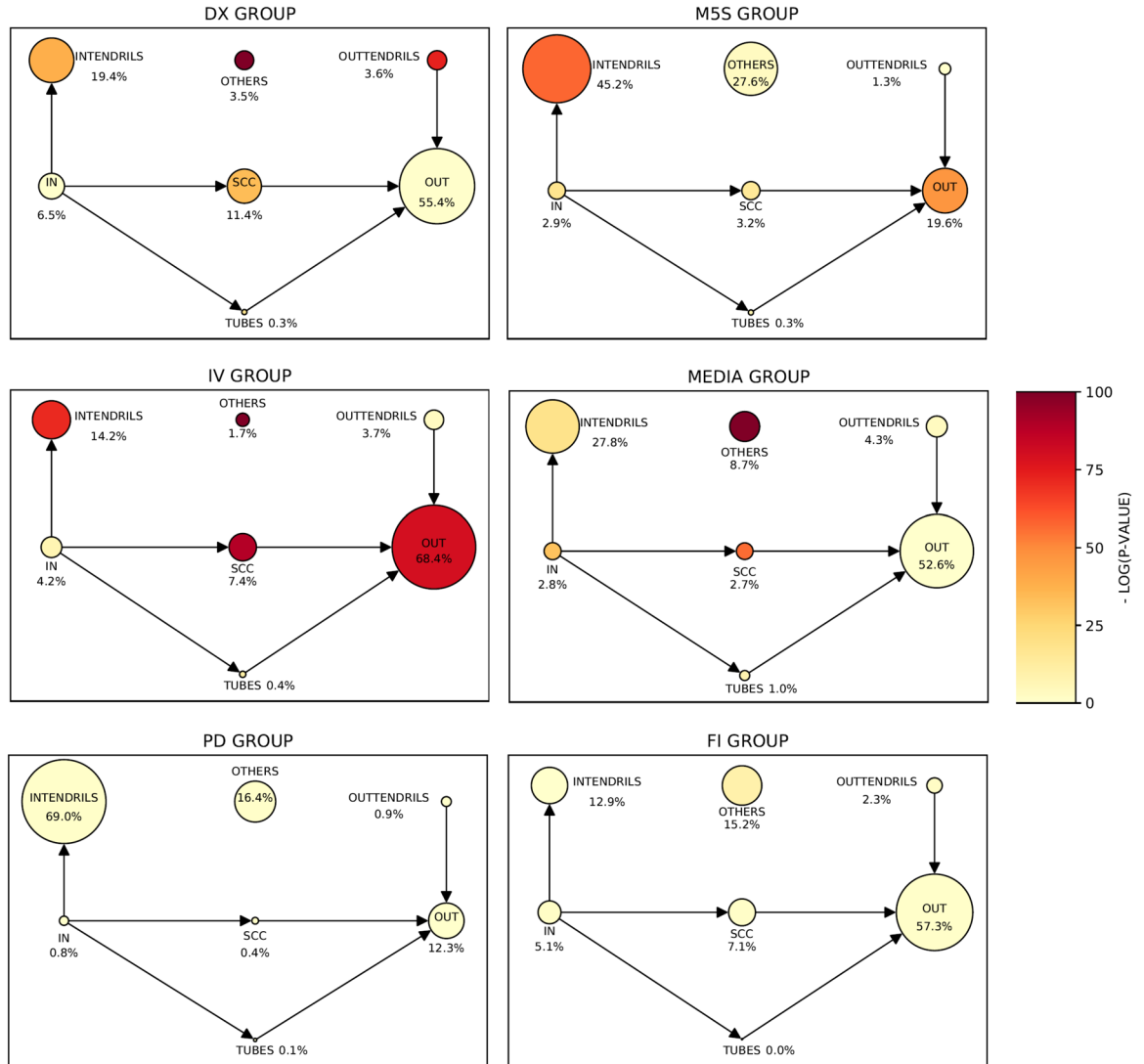


Figure 3. The bow-tie structure of the discursive communities in the Italian COVID-19 dataset. The dimension of the sectors is proportional to the number of nodes: DX and IV discursive communities have strong bow-ties (the OTHERS block is smaller than SCC), while the others are weak (the OTHERS block is greater than SCC, still being smaller than bow-tie WCC). The DX, IV, FI and MEDIA groups display a OUT-dominant bow-tie structure, with the most part of the nodes located in the OUT sector. M5S and PD communities have a INTEND-dominant bow-tie structure, the INTENDRILS sector being the dominant one. The colour of the blocks quantifies the distance between the observed dimensions and those predicted by the Direct Configuration Model (DCM). The observed dimension for the OTHERS sector is significantly less numerous (considering a significance level at $\alpha = 0.01$) for all the communities, but PD. Remarkably, for INTEND-dominant bow-ties, also other sectors, as SCC and INTENDRILS, are usually bigger than what we expect from the model.

SCC is statistically significant (and bigger than expected) for all bow-ties but FI and PD. The IN block is often statistically significant and smaller than expected. We may notice that in the strong bow-tie of IV discursive community the dimensions of all sectors are statistically significant, while none are in the PD bow-tie, which is the smallest discursive community. It is worth noting that also the dimension of the discursive community has a role: due to the limited possible variability, smaller bow-ties feature more agreement with the model.

Verified users' distribution. Usually, verified accounts on Twitter belong to public characters and organizations, such as journalists, politicians, actors, political parties, media, and VIPS in general. Previous studies testify that verified users tend to introduce new content and have high visibility on the platform^{21,26,29,33,53}. Thus, we expect to find them in the IN block. The results in Fig. 4 confirm this intuition: in the case of OUT-dominant bow-ties (leftmost panel), the 33.2% of verified users, on average, are in the IN sector. High percentages of

	SCC	IN	OUT	TUBES	INTE.	OUTTE.	OTHERS
DX •	10^{-35} *	0.7	0.4	10^{-18} *	10^{-39} *	10^{-74} *	0*
M5S •	10^{-15} *	10^{-18} *	10^{-48} *	10^{-9} *	10^{-58} *	0.5	0.0006*
IV •	10^{-90} *	10^{-8} *	10^{-81} *	10^{-12} *	10^{-72} *	0.0004*	0*
PD •	0.04	0.7	0.9	0.08	0.1	0.8	0.1
FI •	0.03	0.01	0.1	0.4	0.5	0.005	10^{-10} *
MEDIA •	10^{-57} *	10^{-33} *	0.9	10^{-9} *	10^{-19} *	0.0002*	0*

Table 1. p-values related to the various bow-tie sectors in the COVID-19 Italian dataset. In orange strong bow-ties, while in teal weak ones; red dots indicate OUT-dominant bow-ties, blue dots indicate INTEND-dominant ones. If we set the statistical significance to $\alpha = 0.01$, then we then have to correct for multiple hypothesis testing per block. In the present case, we used the False Discovery Rate method (FDR⁷⁸). In the table, validated p-values are marked by an asterisk '*'. The OTHERS block is statistically significant (in particular it is smaller than in the randomization) for all discursive communities but the PD one. It is remarkable that the dimension of SCC is significant in all strong bow-ties, while the one of OUT is significant only for IV bow-tie.

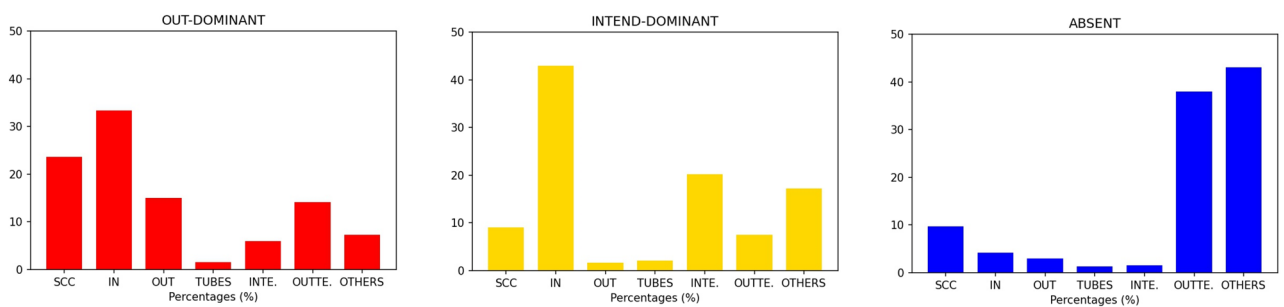


Figure 4. Distribution of the percentage of verified users in each sector of the discursive communities with, respectively, OUT-dominant, INTEND-dominant and not informative bow-ties. Each bar-chart displays the average percentage of verified users in a specific sector, calculated respectively for all the OUT-dominant, INTEND-dominant and not informative bow-ties. In the cases of OUT-dominant and INTEND-dominant bow-ties, the highest percentages of verified accounts can be found in the IN group. Moreover, in OUT-dominant bow-ties, we can find a relevant percentage of verified accounts also in the SCC. Naturally, for those communities with no bow-tie structure the verified accounts are mostly placed in the OTHERS sector and, to less extent, in the OUTTENDRILS one.

verified users are also in the SCC block (23.5%). In the case of INTEND-dominant bow-ties (central panel), the percentage of verified users in the IN group increases to 42.8%; the second block per percentage of verified users is INTENDRILS (20.1%). In those communities where the bow-tie structure is not informative (right panel, Fig. 4), a high percentage (42.9%) of verified users, on average, is in the OTHERS sector. In a few cases of not informative bow-ties, it happens that verified users are mostly in the OUTTENDRILS sector. In this last case, their messages hardly reach a big audience and are simply retweeted by a group of strong retweeters (OUT sector), not catching the interest of the accounts in the SCC. Let us remark that in the case of non-informative bow-ties the dimension of OUT and SCC blocks is nevertheless limited.

Figure 5 reports the same bar-chart, about the presence of verified users, for the bow-ties of the COVID-19 Italian dataset. It is possible to observe that in OUT-dominant bow-ties - i.e., DX, IV, FI and MEDIA - verified users are mainly in IN and SCC sectors. Also, in INTEND-dominant bow-ties, the INTENDRILS sector contains quite a number of verified users. Other user characterizations of the bow-tie blocks can be found in the “Social bots”.

Conservatives groups. The bar-charts in Fig. 6 show the percentage of nodes, the percentage of edges and the number of edges per node in the Strong Connected Component, for each discursive community of the Italian COVID-19 dataset. Not only DX is the one with the greatest number of nodes and the greatest number of links in SCC, but also the link density of SCC in DX is much greater than that of any other discursive community. Thus, the number of links in SCC of DX is not proportionate to the number of nodes, and it results in a greater average degree per node. We found very similar behaviours also for the right-oriented communities of the other datasets.

In fact, in all our datasets, the discursive communities of conservative groups (i.e., DX in the Italian dataset, AfD in the German one, Conservatories in the Dutch one) are those with the highest percentage of nodes and, especially, of edges within SCC. This peculiar feature signals the presence of a common (self-)organization of accounts in line with conservative ideas on Twitter.

NewsGuard (<https://www.newsguardtech.com/it/>) is an independent software toolkit that monitors the quality and transparency of several news websites worldwide. Through the tags that NewsGuard has assigned to news

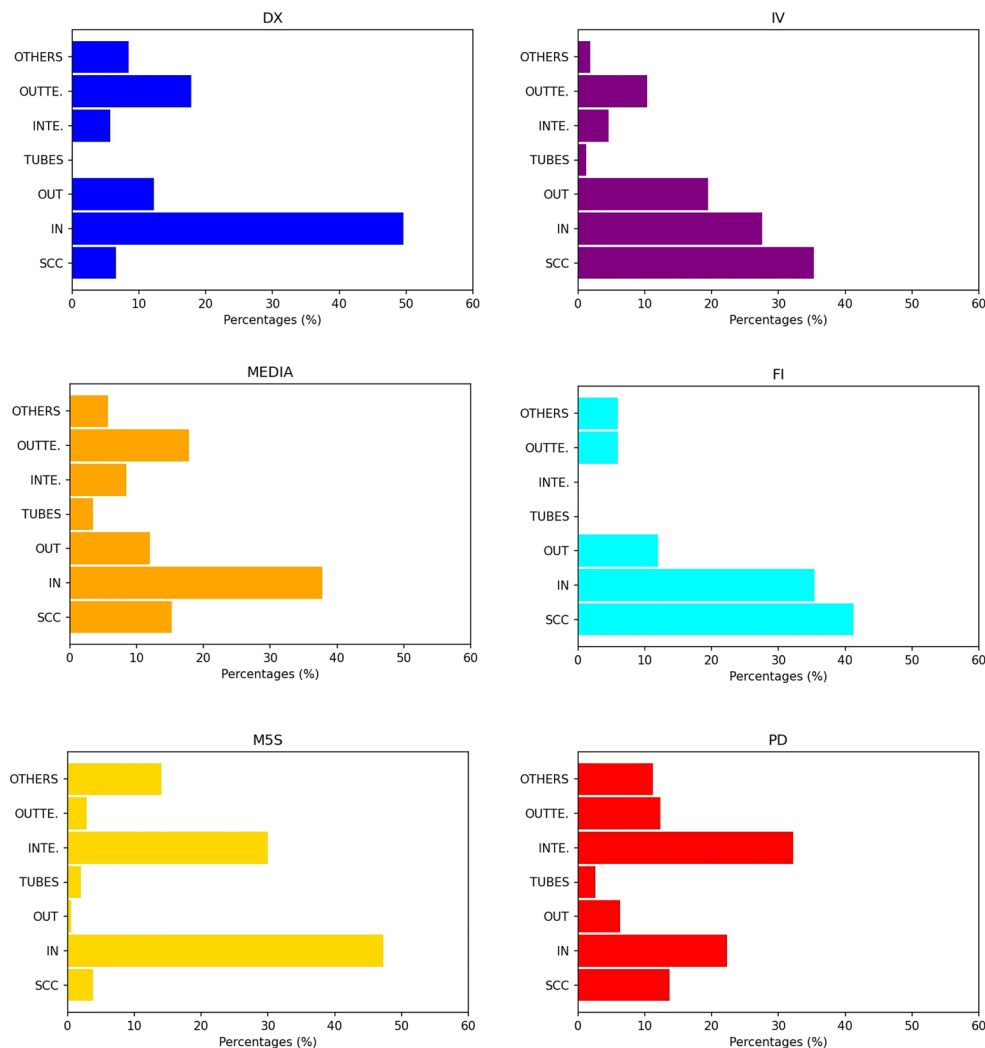


Figure 5. Percentage of verified accounts in the bow-tie sectors for each discursive community of the COVID-19 dataset. The bar-charts confirm that verified accounts are mainly located in the IN sector and, to a less extent, in the SCC one. Only for the PD group, which has an INTEND-dominant bow-tie structure, verified accounts are mostly placed in the INTENDRILS block.

sites whose links appear in the retweets of our communities, we are able to quantify the amount of retweets containing untrustworthy URLs.

The recurrent situation is that almost only the conservative discursive communities display retweets with such URLs. For the Italian COVID-19 dataset, the DX group has 26,318 retweets with links to untrustworthy webpages of news sites, many more than in other communities: 1356 retweets for MSS, 78 retweets for IV, 20 retweets for MEDIA, 9 retweets for FI and 0 for the PD group. A very similar situation has been found for the other datasets, see “Supplementary Material”.

Another interesting aspect is that the most part of retweets containing not reliable URLs has origin in the strongly connected component. Figure 7 shows in red the percentage of retweets containing URLs of untrustworthy news pages within and between the sectors of the bow-tie structure for the DX group. The highest percentage can be found in SCC and between SCC and OUT. Again, this is a recurrent situation also for the conservative communities of the other datasets under investigation.

The case of EURO2020. Here, we devote a specific section to comment about the case of the European football championship (EURO2020) dataset (we do this for academic reasons, and not because Italy won the Euro2020 championship). This dataset features a less divisive, less debated, and less discussed tweets topics. The topics of all the other datasets either have a strong political nature or are debating with sharp different positions. We then analyze whether the fact that topics are less discussed/devated has anything to do with the presence -or absence- of a bow-tie structure in the EURO2020 dataset.

We identified 5 discursive communities for the Italian dataset and 2 discursive communities for the Turkish one. Of these 7, 4 do not have an informative bow-tie structure (in fact, most part of the nodes are in OTHERS),

ITALIAN COVID-19 DATASET

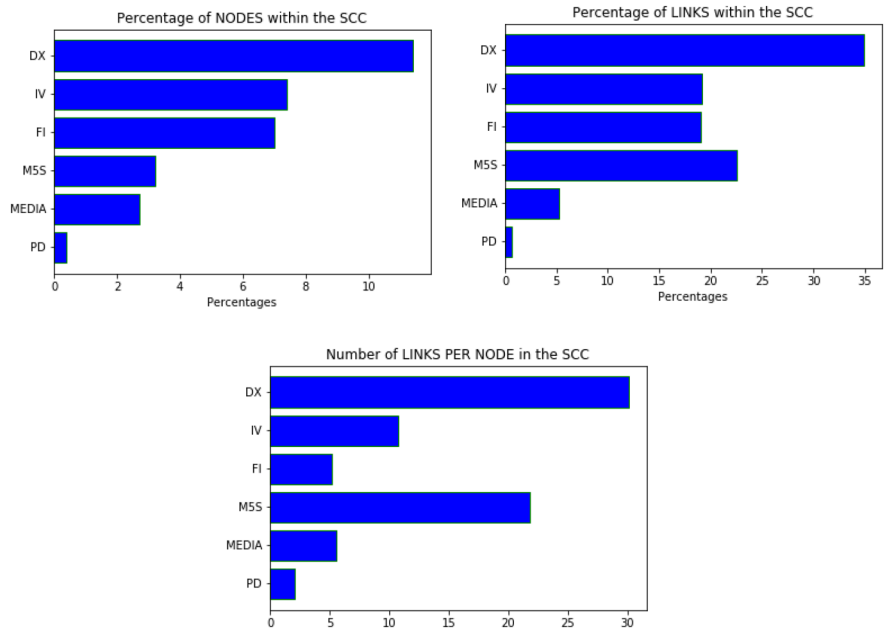


Figure 6. Percentage of nodes and edges in SCC for the communities in the Italian COVID-19 dataset. In the Italian COVID-19 dataset, the conservative and right-oriented discursive community (DX) has more numerous and denser SCCs, as it is displayed in the highest two graphics. In the lowest graphic, it can be seen that, also considering the number of links per node in SCC, DX results again the first discursive community. These results hold for all the conservative groups in all the datasets under investigation.

NEWSGUARD DATA IN DX GROUP

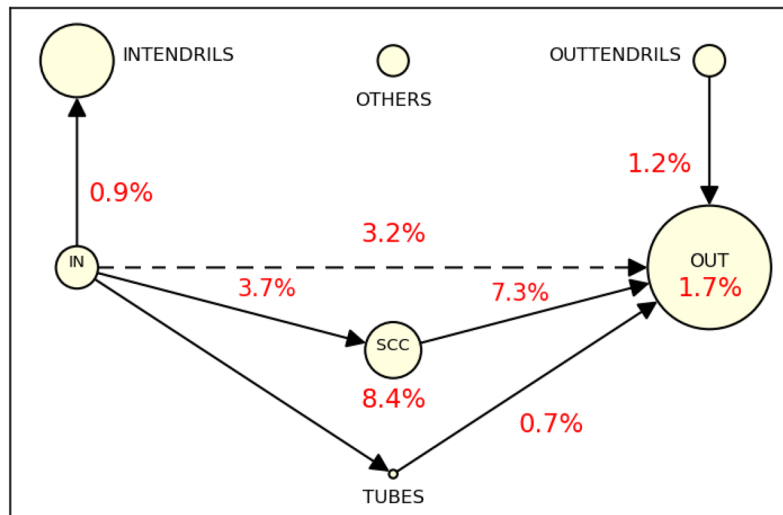


Figure 7. Bow-tie structure of the DX group and percentages of retweets containing URLs of untrustworthy webpages. The DX community in the Italian COVID-19 dataset presents the highest number of retweets containing a link to untrustworthy webpages. Most of them origin from SCC: 8.4% of the retweets in SCC and 7.3% of the retweets between SCC and OUT contain not reliable URLs. In the diagram, we also insert the link between IN and OUT (the dashed line), which, considering the definition of each sector, is not forbidden a priori.

TURKISH EURO2020 DATA-SET

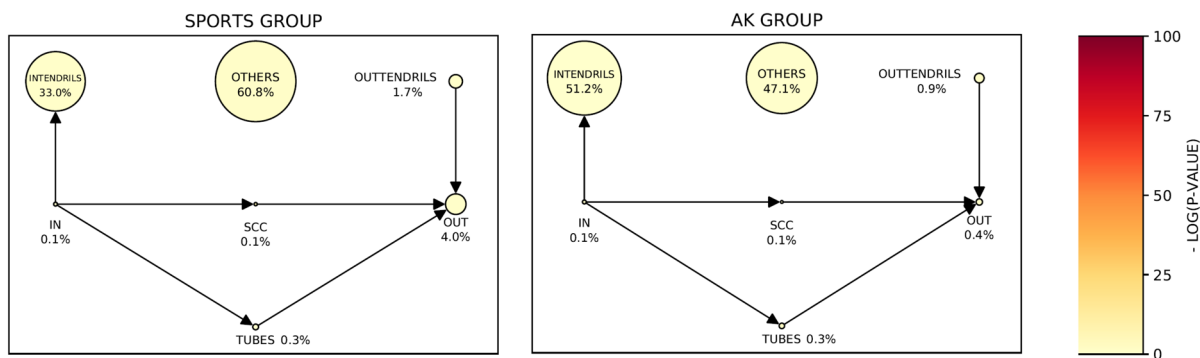


Figure 8. The bow-tie structure of the discursive communities for the Turkish EURO2020 dataset. The SPORTS group contains the official accounts of football players and clubs, and sports newspapers, while AK refers to the Justice and Development Party (Turkish: Adalet ve Kalkınma Partisi, AKP), which is a conservative political party in Turkey, including President Erdogan and his ministries. The SPORTS discursive community does not display an informative bow-tie structure, while the AK one has an extremely weak (INTEND-dominant) bow-tie. The dimension of the sectors is proportional to the number of nodes therein and the color quantifies the distance between the observed and the predicted dimension. Looking to the color of the vertices, it is possible to see that the observed dimensions are not statistically significant.

and the other three have a weak one (OTHERS is smaller than the weakly connected component of the bow-tie, but still greater than the strongly connected one).

Figure 8 reports the bow-tie structures of the two discursive communities in the Turkish dataset. The SPORTS group contains the official accounts of football players and clubs, and of sports newspapers. AK refers to the Justice and Development Party (Turkish: Adalet ve Kalkınma Partisi, AKP), which is a conservative political party in Turkey including President Erdogan and his ministries. While SPORTS does not display any informative bow-tie, AK has a weak one. Following our interpretation, the latter displays a more hierarchical conversation on Twitter, in which the SCC is not numerous. Moreover, the dimensions of the sectors are mostly not statistically significant.

For the Italian case (Fig. 9) the main discursive community is formed by football players, sports newspapers and journalists. There is also a MEDIA community, containing accounts of Italian media, and other three small political communities (DX, IV, M5S). MEDIA, DX and IV does not display an informative bow-tie structure (respectively, 74%, 81.2% and 63.6% of the nodes are in OTHERS), while FOOTBALLERS and M5S show a weak bow-tie (respectively 15.9% and 23.9% of nodes in OTHERS).

Euro2020 dataset is the only, among ours, in which no discursive communities have a strong bow-tie structure.

Discussion

In the present manuscript, we analysed eight thematic Twitter datasets in different languages, related to various debates in Europe. We identified the discursive communities in the retweet networks and we investigated the presence of bow-tie structures in such communities. In previous works, discursive communities were shown to mirror the political orientation of users^{21–24,26,29,30,32–34}, thus the analysis of their structure is of utmost importance to infer the way opinions create and circulate.

Discursive communities and bow-ties. We found that a bow-tie structure is present in those discursive communities debating about politics, like in the case, e.g., of election campaigns (it is the case of the Dutch elections dataset) or debating about Society, e.g., ‘how to handle a pandemic?’ (it is the case of the Italian, German and French datasets about COVID-19) or ‘how to manage migration fluxes?’ (it is the case of the Italian online debate on migrants). Instead, a bow-tie structure is absent when the topics of the discussion are sportive ones, as in the case of Euro2020 Turkish and Italian datasets.

More in details, we state that the bow-tie is informative if the corresponding WCC includes more than one half of the nodes of the entire discursive community. In the present datasets, we found that bow-ties are informative in all the discursive communities debating about politics. In the case of the Euro2020 dataset, bow-ties are not informative, or, if present, they are extremely weak. When the bow-tie is informative, we found essentially 2 cases: (1) the most crowded block is the OUT one; (2) the most crowded block is the INTENDRILS one. The former is typical of the discursive communities of right wing parties in all European political/societal debates of our datasets, while the latter is more common in less active political discursive communities in many political/societal datasets.

Which users in which bow-tie sectors and the exposure to m/disinformation. A closer inspection of the nodes in the various blocks and the quality of the shared content permit to better characterise the users in the bow-tie. The first observation is that the greatest part of the verified users, i.e., those accounts for which the identity of the owner has been certified by Twitter, in the IN sector, in each bow-tie. This finding is not surprising: as already observed in previous studies, verified users create content and are less active in shar-

ITALIAN EURO2020 DATA-SET

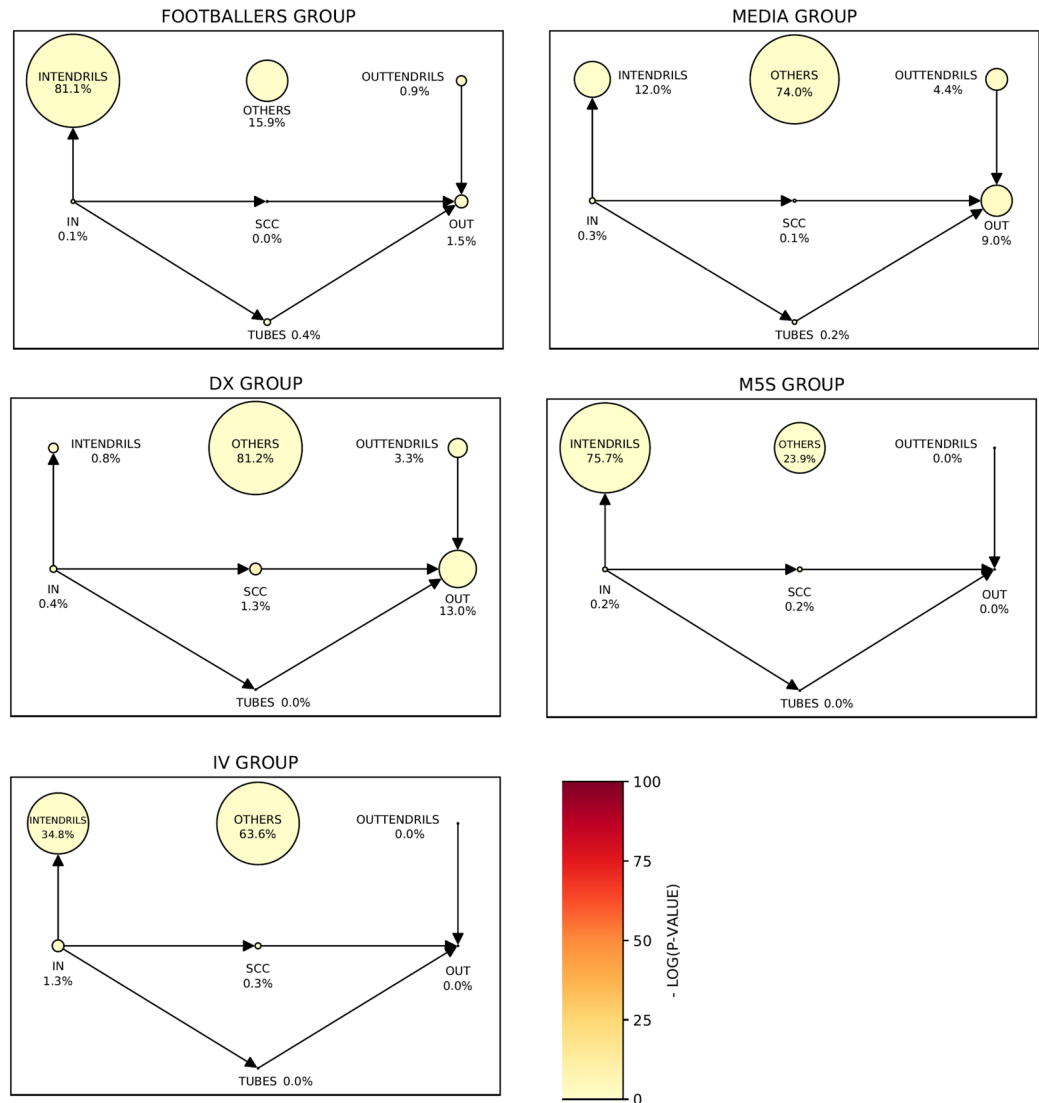


Figure 9. The bow-tie structure of the discursive communities for the Italian EURO2020 dataset. The dimension of the sectors is proportional to the number of nodes therein and the color quantifies the distance between the observed and the predicted dimension. The main discursive community is formed by football players, sports newspapers and journalists. Then, we identified a MEDIA community, containing accounts of Italian media, and three small political communities (DX, IV, M5S). MEDIA, DX and IV do not display an informative bow-tie structure (respectively 74%, 81.2% and 63.6% of the nodes in OTHERS), while FOOTBALLERS and M5S show a weak bow-tie (respectively 81.1% and 75.7% of nodes in INTENDRILS).

ing messages written by others^{21,26,29,33,47}. Verified users are mostly politicians and official accounts of political parties, as well as journalists and official accounts of their newscasts and newspapers. In this sense, a discursive community displaying a INTEND-dominant bow-tie structure (where INTRENDILS is the most crowded block) may appear, at first sight, as a less democratic group: the content is created by a few accounts and shared by a group of followers that limit their interactions to sharing the messages coming from the IN block. Instead, in a OUT-dominant bow-tie, the greatest block is OUT and it can access the content created by all the other blocks in the bow-tie (with the only exception of INTENDRILS), so having the possibility to intercept every voice in the discursive community.

Actually, the issue is on the quality of the content created in the various blocks, see Fig. 7. Leveraging our ongoing collaboration with the NewsGuard organization (<https://www.newsguardtech.com/it/>), we annotated the URLs that appear in tweets in our datasets, based on the reliability and transparency ratings of the news sites to which those URLs belong (such ratings have been assigned by NewsGuard). It turns out that the lowest reliable URLs, in a strong bow-tie, are the ones shared in SCC. The fact that verified accounts are not responsible for the vast majority of m/disinformation sharing was already observed in Ref.³³ and, in the present context, it

reflects the fact that accounts in IN are minimally responsible for the spreading of low quality/untrustworthy content. Otherwise stated, when the source of information is not identifiable, the average quality of the content is lowered down.

An OUT-dominant bow-tie is, in this sense, more exposed to m/disinformation campaigns, as the majority of the accounts, i.e. those in the OUT block, is exposed to a great flow of content, in which the percentage of m/disinformation is quite high. On the other hand, the INTEND-dominant bow-tie is “safer”, since the greatest part of the accounts therein (i.e. the INTENDRILS nodes) accesses the messages from the IN sector that is less prone to m/disinformation campaigns.

It worth to be remarked that, due to the considerations above, the OUT-dominant bow-ties are at risk of infodemic. *Infodemic* is a recently introduced neologism, that became particularly popular during the COVID-19 pandemic. According to the WHO, “*infodemics are an excessive amount of information about a problem, which makes it difficult to identify a solution. Infodemics can spread misinformation, disinformation and rumors during a health emergency. Infodemics can hamper an effective public health response and create confusion and distrust among people*” (Coronavirus disease 2019 (COVID-19) Situation Report—45. https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200305-sitrep-45-covid-19.pdf?sfvrsn=ed2ba78b_4). The effects of the present COVID-19 infodemic, even if debated^{4,5}, may put at risk the countermeasures to the spread of an epidemic and it is worrisome for policy makers (see, for instance, the Joint Communication titled “Tackling COVID-19 disinformation—Getting the facts right” (June 10th, 2020), available at the following link. https://ec.europa.eu/info/sites/default/files/communication-tackling-covid-19-disinformation-getting-facts-right_en.pdf).

Finally, let us consider also the peculiarity of right-wing discursive communities: for all those, the bow-tie is strong (i.e., the dimension of the OTHERS block is smaller than the SCC one) and it is neatly OUT-dominant. The structural exposure of the OUT-dominant bow-tie to infodemic is even more emphasized by the extreme activity of the SCC: for instance, in the COVID-19 Italian dataset the link density in the right-wing bow-tie is at least 3 times greater than any other OUT-dominant strong bow-ties.

Statistical significance of the analysis. Here, we remark an important aspect of our analysis, of utmost importance. In the analysis of a complex network, it is necessary to consider what is being measured, and what is its baseline. A typical example is the modularity, i.e. one of the most used target function for community detection. The problem resides in stating what is the number of links inside a group of nodes that is enough to form a community. In this case, we build a null-model, i.e., a model that shows part of the properties of the original system, being random for all the rest, to have a proper benchmark for our observations. We then compare the number of edges inside a group of nodes with the one expected by the null-model. Without the null-model, we could not know whether the number of links that bind a group of nodes are due to the degree sequence, or whether they are instead the genuine signal of the presence of a community.

In the present study, we used an entropy-based null-model as a benchmark for our analysis^{27,79}. An entropy-based null-model allows to have a benchmark that is tailored to the system under analysis. It fixes (on average) some topological quantities to the values observed in the real network and leaves all the rest completely random. Being based on the (Shannon) entropy maximisation, it guarantees that it uniformly considers all the possible configurations (it is ‘ergodic’, using Statistical Physics jargon), thus it does not introduce any bias in the analysis.

To strengthen the analysis, we study if the bow-tie structures are due to the degree sequence of the nodes in the various discursive communities. In fact, the size of IN and OUTTENDRILS could simply be due to the presence of many nodes with zero in-degree (an analogous consideration could be done for the OUT and the INTENDRILS blocks, considering, instead the out-degree). Thus, strong, weak and not informative bow-ties could be due to degree sequence only, and do not carry any kind of information on their own.

We thus used the Directed Configuration Model defined in Ref.⁵⁴ and implemented by the Python module *NEMtropy*⁷⁰. Our results show that the dimensions of the blocks in the bow-tie are very often statistically significant: the p-value of the observed dimensions of the various blocks against the null-model expected distribution are extremely small, such that they are not compatible with the degree sequence, or, otherwise stated, the dimension of the various blocks cannot be explained using the degree sequences only.

Limitations. Even if we have obtained strong results (see the null-model validation check on the dimension of the bow-tie sectors), we have nevertheless to remark few aspects of our analysis that can limit its generalization. First, the analysis is related to eight different thematic datasets in different languages, all referring to European debates, some of them of political nature. Indeed, while the total amount of messages analysed is quite impressive, we are aware that, even if the spectrum of the arguments covered is various, our findings may be valid on our datasets only. In the near future, we are going to expand the countries covered by our analyses and expand the list of debated arguments.

Following our jargon, OUT-dominant bow-ties expose the majority of their accounts to the risk of infodemic. Nevertheless, it is not a causal relation: the presence of OUT-dominant bow-ties does not imply the presence of an infodemic or of a disinformation campaign. In fact, if the sources shared by SCC are reputable, we will not observe any infodemic or m/disinformation signal. At the same time, it is true that OUT-dominant bow-ties help the diffusion of m/disinformation, when present, since accounts in OUT are exposed to all contents—reliable and not reliable—created by nearly every block in the discursive community.

Final remarks. Let us conclude with some final remarks. First, the bow-tie structure is present in the discursive communities of retweet networks. Let us recall that we build the retweet network by creating a direct link for every retweet, from the author of the original post to the retweeter. Then, from the retweet network we extract the subgraphs relative to the discursive communities obtained through the procedure described in “[Discursive](#)

communities”. With this procedure, we are recovering the flow of information *inside* each discursive community; in this sense, we are disregarding the possible interactions among different discursive communities.

There is another limitation which is unavoidable, due to the nature of Twitter’s data: we have information regarding who retweeted who, but not on the “chain” of retweets, i.e. we cannot distinguish if the retweeter retweeted directly the message of the original author of the post or through one of the retweets given by one of followers of the original author.

Given the structure of the retweet network, it is therefore natural to ask what is the meaning of the bow-tie structure. In particular, what is the sense of the reachability of nodes in the retweet network? (We are thankful to reviewer 2 for suggesting this reflection that permits to give a clearer frame to our results.)

In fact, the retweet network describes an influence flow: users, by retweeting messages, testify that they are influenced by their opinions. In this sense, the opinions expressed in messages written by a user *A* influence the ones of her/his retweeter *B*, that, on turn, influence the one of her/his retweeter *C* and so on, even if *C* has never directly retweeted any of the messages of *A*. Otherwise stated, what the bow-tie is capturing is a division of users in different sets: in *IN* we find the creators of contents, only partially influenced by the content created by others, in *OUT* the big audience of standard users, influenced by the contents created in *IN*, *SCC*, *TUBES* and *OUTTENDRILS*, in *SCC* the active users influencing each other, and so on.

Finally, even if we have not access to the chain of retweets, we still expect the retweet network to be a good proxy of the followership network. Indeed, users either retweet messages from the accounts they follow or they retweet some message when searching for a specific topic. However, we expect the first method to be way more frequent, since in the user’s home the activity (i.e. new messages, likes and retweets) of the followed accounts are present. In this sense, we expect that the bow-tie will be present (and statistically significant) even in the followers network. Such an investigation is going to be part of future research.

Methods

Bow-tie detection. In the following we briefly described the main steps of the detection of the bow-tie structure, following the procedure outlined in Ref.⁴¹. The first step is the identification of the greatest Strongly Connected Component (*SCC*) and then the identification of the nodes in the various sectors, using the bow-tie definition.

Let be $G_D(V, E)$ a directed graph where V is the set of nodes and E the set of links, and $G_D^T(V, E)$ its counterpart obtained reversing the direction of the edges. The functions for the identification of the greatest *SCC* and the depth-first searches (*DFS* in the following), used to identify nodes reachable from a given node, are implemented in many python modules, such as [igraph](#) or [networkx](#); for the present analysis, we used the former python module. The algorithm is pretty straightforward and follows the definitions in “[The bow-tie structure](#)”: the pseudo code is presented in Algorithm 1.

Philosophically, the algorithm works as follows. First, consider the greatest strongly connected component and call it *SCC*; then choose randomly a node $v \in SCC$. All nodes that can reach v (identified via *DFS*), but that are not part of *SCC* represent the *IN* sector; an analogous line of reasoning takes to the identification of the *OUT* sector. Regarding the remaining node, the crucial information to be calculated is if they can be reached by nodes in *IN* and if they can reach nodes in *OUT* (again, using *DFS*): based on this final pieces of information, we can identified all remaining sectors.

Algorithm 1 Bow-tie detection algorithm

```

SCC ← all the nodes in the greatest Strongly Connected Component
v ← a randomly chosen node in SCC
DFSGD(v) ← all the nodes that can be reached by v
DFSGDT(v) ← all the nodes that can reach v
OUT ← DFSGD(v)/S
IN ← DFSGDT(v)/S
remainingV ← V/SCC/IN/OUT           ▷ All nodes that are not in SCC, neither in IN, neither in OUT.
for w ∈ remainingV do
  IRW ← (IN ∩ DFSGDT(w) ≠ ∅)           ▷ “IN Reaches w?”
  WRO ← (OUT ∩ DFSGD(w) ≠ ∅)         ▷ “w Reaches OUT?”
  if IRW and WRO then
    w ∈ TUBES
  else if IRW then
    w ∈ INTENDRILS
  else if WRO then
    w ∈ OUTTENDRILS
  else
    w ∈ OTHERS
  end if
end for

```

Entropy-based null-models for network analysis and their applications. *The bipartite configuration model.* In order to create the various discursive communities we needed an appropriate null-model as benchmark for identifying those verified users that share the same audience. In this sense, it is necessary to compare the observed quantities with accurate predictions in order to state their significance: actually, the common audience may appear similar just due to the extreme activity of the considered verified users.

We represent the interaction between verified accounts—the ones whose identity is certified by Twitter platform- and unverified ones (i.e. all the others) via a bipartite undirected binary network in which a link connects a verified users to an unverified ones if there is at least a retweet between one and the other, or viceversa. Since the information about the number of different accounts interacting -via tweet or retweet- with a user is encoded, in this representation, in the degree sequence for nodes of both layers, we need a benchmark discounting it. The natural choice is to choose an entropy-based null-model, since it provides, by definition an unbiased framework²⁷: the null-model is maximally random, but for the constraints imposed on the system. The bipartite null-model discounting the degree sequence is the Bipartite Configuration Model (BiCM⁴⁸). In the present section we will briefly revise the steps of its definition.

Let us consider a bipartite network in which the two layers \top and \perp have dimension, respectively, N_{\top} and N_{\perp} ; in the following, Latin indices will be used to identify nodes on the \top layer while Greek ones will be used for the \perp layer. Then, the bipartite network can be represented by its biadjacency matrix, i.e. a $N_{\top} \times N_{\perp}$ matrix \mathbf{M} whose generic entry $m_{i\alpha}$ is 1 if the node $i \in \top$ is connected to the node $\alpha \in \perp$ and 0 otherwise.

Let us start from a real bipartite network G_{Bi}^* (in the following, all quantities denoted by a * will indicate those measured on the real network). First, let us define an ensemble of graphs, i.e. the set of all the possible bipartite graphs having the same number of nodes of G_{Bi}^* , but with all different topologies, from the fully connected to the empty ones. Then, we can define the Shannon entropy over the ensemble, by assigning a different probability to each of its elements:

$$S = - \sum_{G_{\text{Bi}} \in \mathcal{G}_{\text{Bi}}} P(G_{\text{Bi}}) \ln P(G_{\text{Bi}});$$

where, $P(G_{\text{Bi}})$ is the probability of the generic element of the graph ensemble G_{Bi} . Let us now maximise the entropy, while constraining the network degrees: in particular, we want that the ensemble average of degrees to match the value observed on the real network, in order to have a null-model tailored to the real system. In term of the biadjacency matrix, the degree sequences of the \top and \perp layers respectively read $k_i = \sum_{\alpha} m_{i\alpha}$ and $h_{\alpha} = \sum_i m_{i\alpha}$. Using the method of the Lagrangian multipliers, the constrained maximisation can be expressed as the maximisation of S' , defined as

$$\begin{aligned}
 S' = & S \\
 & + \sum_i \eta_i \left[k_i^* - \sum_{G_{Bi} \in \mathcal{G}_{Bi}} P(G_{Bi}) k_i(G_{Bi}) \right] + \sum_{\alpha} \theta_{\alpha} \left[h_{\alpha}^* - \sum_{G_{Bi} \in \mathcal{G}_{Bi}} P(G_{Bi}) h_{\alpha}(G_{Bi}) \right] \\
 & + \zeta \left[\sum_{G_{Bi} \in \mathcal{G}_{Bi}} P(G_{Bi}) - 1 \right]
 \end{aligned}$$

where S is the Shannon entropy defined above, η_i, θ_{α} are the Lagrangian multipliers relative to the degree sequences, respectively, on \top and \perp , and ζ is the one relative to the probability normalization.

Maximising S' leads to a probability per graph $G_{Bi} \in \mathcal{G}_{Bi}$ that can be factorised in terms of the probabilities per link $p_{i\alpha}$ ⁸⁰, i.e.

$$P(G_{Bi}) = \prod_{i,\alpha} p_{i\alpha}^{m_{i\alpha}(G_{Bi})} (1 - p_{i\alpha})^{1 - m_{i\alpha}(G_{Bi})}, \tag{1}$$

where $p_{i\alpha} = \frac{e^{-\eta_i - \theta_{\alpha}}}{1 + e^{-\eta_i - \theta_{\alpha}}}$. Nevertheless, at this level the above equation is just formal, since we do not know the numerical value of η_i and θ_{α} . To this aim, we can then maximise the likelihood of the real network^{52,81}; it can be shown that the likelihood maximisation is equivalent to imposing

$$\langle k_i \rangle_{\text{BiCM}} = k_i^*, \forall i \in \top; \quad \langle h_{\alpha} \rangle_{\text{BiCM}} = h_{\alpha}^*, \forall \alpha \in \perp.$$

Validated projection of bipartite networks. We want to infer similarities among nodes on the same layer. We can use as a measure of similarity the number of common neighbours—for each couple of verified users, the number of unverified users that have interacted, via tweet or retweet, with both. Let us assume, without loss of generality, that we want to project the information contained in the bipartite network onto the \top layer and call V_{ij} the number of common neighbors between nodes $i, j \in \top$ (following Ref.⁵⁸, we use the letter V to indicate common neighbours, since this pattern appear in the bipartite network as a “ V ” between the layer).

In terms of the biadjacency matrix, V_{ij} can be expressed as

$$V_{ij} = \sum_{\alpha} V_{ij}^{\alpha} = \sum_{\alpha} m_{i\alpha} m_{j\alpha},$$

where we have defined $V_{ij}^{\alpha} = m_{i\alpha} m_{j\alpha}$; $V_{ij}^{\alpha} = 1$ if both i and j are connected to node $\alpha \in \perp$ and 0 otherwise. Let us now compare the observed V_{ij} for each possible pair of nodes in \top with the prediction of the BiCM. Since link probabilities are independent, the presence of each V -motif V_{ij}^{α} can be regarded as the outcome of a Bernoulli trial:

$$\begin{aligned}
 f_{\text{Ber}}(V_{ij}^{\alpha} = 1) &= p_{i\alpha} p_{j\alpha}, \\
 f_{\text{Ber}}(V_{ij}^{\alpha} = 0) &= 1 - p_{i\alpha} p_{j\alpha}.
 \end{aligned}$$

In general, the probability of observing $V_{ij} = n$ can be expressed as a sum of contributions, running on the n -tuples of considered nodes (in this case, the ones belonging to the layer of users). Indicating with A_n all possible nodes n -tuples among the layer of \perp , this probability amounts at

$$f_{PB}(V_{ij} = n) = \sum_{A_n} \left[\prod_{\alpha \in A_n} p_{i\alpha} p_{j\alpha} \prod_{\alpha' \notin A_n} (1 - p_{i\alpha'} p_{j\alpha'}) \right], \tag{2}$$

where the second product runs over the complement set of A_n . Eq. (2) represent the generalization of the usual Binomial distribution when the single Bernoulli trials have different probabilities, also known as Poisson Binomial distribution⁸².

We can, then, verify the statistical significance of the observed co-occurrences by calculating their p-value according to the distribution in Eq. (2), i.e. the probability of observing a number of co-occurrences greater than, or equal to, the observed one:

$$\text{p-value}(V_{ij}^*) = \sum_{V_{ij} \geq V_{ij}^*} f_{PB}(V_{ij}^*). \tag{3}$$

Repeating this calculation for every pair of nodes, we obtain $\binom{N_{\top}}{2}$ p-values. In order to state the statistical significance of the hypotheses belonging to this group, it is necessary to adopt a multiple hypothesis testing correction; in the present paper, we use the False Discovery Rate (FDR⁸³), since it controls the false positives rate.

Direct configuration model. From the entire retweet network, in which the various accounts are represented as nodes in a direct network in which an arrow points the retweeter of a post, starting from its author, we extracted the various subgraphs of discursive community. Then, in order to compare the observed dimensions of the bow-tie sectors of these subgraphs and state their statistical significance, we adopted the *Direct Configuration Model*

(DCM), which is the entropy-based model suited for direct monopartite networks⁴³. For directed networks, the adjacency matrix is (in general) not symmetric, and each node i is characterized by two degrees: the out-degree $k_i^{\text{out}} = \sum_j a_{ij}$ and the in-degree $k_i^{\text{in}} = \sum_j a_{ji}$, where a_{ij} is the generic entry of the (directed) adjacency matrix \mathbf{A} . The Directed Configuration Model (DCM) is therefore defined as the ensemble of direct networks with given out-degree and in-degree sequences. Using the same machinery as in the previous section “Entropy-based null-models for network analysis and their applications”, it is possible to derive a probability per graph: if G_D is the generic representative of the ensemble of directed graphs \mathcal{G}_D , then the probability per graph $P(G_D)$ reads:

$$P(G_D) = \prod_{i,j \neq i} q_{ij}^{a_{ij}(G_D)} (1 - q_{ij})^{1 - a_{ij}(G_D)}.$$

Thus, again the probability per graph factorises in terms of probabilities per link q_{ij} , which can be expressed in terms of Lagrangian multipliers

$$q_{ij} = \frac{e^{-\gamma_i - \delta_j}}{1 + e^{-\gamma_i - \delta_j}},$$

where γ_i and δ_j are the Lagrangian multipliers associated, respectively to the out-degree of node i and to the in-degree of node j . In order to get the numerical value of γ_i and δ_j , we can use the maximum likelihood as in the above section “Entropy-based null-models for network analysis and their applications”, which is equivalent to impose

$$\langle k_i^{\text{out}} \rangle_{\text{DCM}} = k_i^{*\text{out}}, \quad \langle k_i^{\text{in}} \rangle_{\text{DCM}} = k_i^{*\text{in}}, \quad \forall i.$$

Since the bow-tie decomposition is highly non linear, in order to calculate the statistical significance of the dimension of the various blocks, we generated a sample of 1000 different graphs for each discursive community, using the probabilities provided by the DCM. Then, we obtained a distribution for the dimensions of the bow-tie sectors just looking to the decomposition of each graph in our ensemble. At this point, we could calculate a two-tailed p-value with a significance at $\alpha = 0.01$ for estimating the distance between the dimensions observed with those reproduced by the ensemble.

Modularity and community detection. In the present analysis, we inferred the discursive communities from the communities in the validated network of verified users. In particular, we used the modularity based Louvain algorithm⁴⁹.

The modularity⁸⁴ compares the number of edges within the actual communities with its expectation under a certain null-model. Modularity can be written as

$$Q = \frac{1}{2m} \sum_{ij} (a_{ij} - p_{ij}) \delta(C_i, C_j) \quad (4)$$

where m is the total number of links of the network, a_{ij} are the entries of the adjacency matrix, p_{ij} is the probability to have a link between nodes i and j according to the chosen null-model, C_i and C_j are, respectively, the communities of nodes i and j and the Kronecker delta $\delta(C_i, C_j)$ selects all the pairs of nodes contained in the same community (equal to 1 if $C_i = C_j$ or 0 otherwise). In the original definition in Ref.⁸⁵, the null-model chosen is the Chung-Lu one⁵¹, which conserve the degree sequence, but it is known to be inconsistent for dense networks that present strong hubs²⁷. In the present paper we use instead the entropy-based Undirected Configuration Model (UCM) defined in^{52,81}: it can be shown that in the case of sparse network, the UCM can be approximated by the Chung-Lu null-model.

Furthermore, Louvain algorithm is known to be order dependent, i.e. the resulting configuration depends on the order of the nodes given to the algorithm. In order to avoid this bias, we run 100 times the algorithm reshuffling each time the node order. At the end of 100 runs, we select the final partition displaying the maximum of the objective function (in our case, the modified modularity with the UCM).

The multiple runs approach is quite common^{26,38,39,58,75}, but different approaches are present in the literature for the final choice of the resulting partition in communities: for instance, in Refs.^{38,39} the authors, instead of choosing the partition with the greatest value of the modularity, prefer to choose the most common clusters, the choice being motivated by the specific profile of modularity⁸⁶. While the procedure proposed in Refs.^{38,39} is perfectly acceptable, we prefer ours, since it targets directly the objective function we are considering.

NEMtropy. In the present paper, we implemented the BiCM, the DCM and the Louvain algorithm using UCM null-models via the Python module **NEMtropy**, described in Ref.⁷⁰.

Discursive communities for the Italian COVID-19 dataset. Here, we give a brief description of the discursive communities identified in the Italian COVID-19 dataset (Their dimensions are in Fig. 2):

- **DX:** this community collects the official accounts and the main leaders of two Italian right-oriented political parties, ‘Lega’ and ‘Fratelli d’Italia’;
- **M5S:** this community contains the main politicians and the official accounts of the Italian party ‘Movimento 5 Stelle’ (English: 5 Stars Movement), an anti-establishment political movement;

Algorithm	Precision	Recall	F-Measure
RandomForest	0.838	0.828	0.833
Ripper	0.819	0.823	0.821
Multilayer Perceptron	0.709	0.737	0.723
NaiveBayes	0.554	0.986	0.709 2
IBk	0.708	0.726	0.717

Table 2. Performance results after 10-folds cross validation on *cresci-stock-2018* data-set.

- **IV:** this community is associated to the liberal party of ‘Italia Viva’ (English: Italy Alive) with centre/centre-left political positions;
- **PD:** this cluster contains the politicians of the Italian ‘Partito Democratico’ (English: Democratic Party), the traditional centre-left party;
- **FI:** this group collects the politicians and the official accounts of the Italian centre-right party of ‘Forza Italia’ (English: Go, Italy!);
- **MEDIA:** this type of community is present in almost all the datasets we analyzed. It contains the official accounts of newspapers, journalists, TV-channels, radio channels and in general other Italian media.

Social bots. Social bots are computer algorithms whose behaviour on social platforms is often far from being benign: malicious bots are purposely created to distribute spam, sponsor public characters and, ultimately, induce a bias within the public opinion^{87–89}. Often these agents have the task of increasing the visibility of certain users^{29,33}.

Here, we report the outcome of a study about detection of social bots in the datasets under investigation. For bot detection, we exploit the general-purpose bot detection system based on supervised-learning presented in⁹⁰. Such a system has been shown to be highly accurate, both for unveiling automated accounts that work alone and those that participate in coordinated activities. The bot detector is ‘traditional’, i.e., only one user per time is analyzed during the classification process⁹¹.

The classifier exploits so-called Class A features, i.e., features that can be directly extracted from the user profile. These features were originally introduced in⁸⁷ and, despite their simplicity, proved to be still effective for the detection of novel bots too. Features that are known to be the most expensive to compute (mainly in terms of time needed for data gathering), namely those concerning the account’s relationships (friends and followers) have been disregarded.

Hence, in order to decide about the type of the account (either a bot or not), we (1) train and evaluate different machine-learning algorithms on a dataset where bots and genuine accounts are a priori known, (2) we select the model with the best classification performances and (3) we apply the resulting model to the accounts of the datasets investigated in the main text of the manuscript.

To train and validate the classifier we leverage the publicly available *cresci-stock-2018* (<https://botometer.osome.iu.edu/bot-repository/datasets.html>) dataset. In particular, we use the accounts metadata of the 6842 bots and 5880 human that were still active at the time of data collection; data were crawled on July 2020 through the Tweepy library (<https://www.tweepy.org/>).

To select the best model, we consider five algorithms, each of them belonging to a different category: MLP (Multilayer Perception)⁹², JRip, i.e., a Java-based implementation of the RIPPER algorithm⁹³, Naive Bayes⁹⁴, Random Forest⁹⁵, and the Weka⁹⁶ implementation of the Instance-based Learning Algorithms, i.e., IBK⁹⁷.

The performances of the five different algorithms are evaluated in terms of standard metrics, such as balanced accuracy, precision, and f-measure. The metrics are computed using a 10-fold cross-validation.

For all the experiments, we rely on the open source (Java-based) Weka framework that provide us the implementations of (1) the five machine-learning algorithms (for which we use the default parameter settings, <https://waikato.github.io/weka-wiki/documentation/>), (2) the evaluation metrics and (3) the process of 10-fold cross validation.

In light of our experiments (see Table 2), we select the Random Forest-based model as the classification process since it outperforms the other models.

The resulting model for bot classification is then applied to tag all the accounts involved in our study, giving an average concentration of bots that is around 23.9% in total. In particular, if we focus on specific datasets, we observe percentage of bots around 23.9% for **Italian COVID-19**, 29% for **German COVID-19**, 23.4% for **French COVID-19**, 22.8% for **Dutch elections**, 25.7% for **Italian debate on migrants**, 24% for **Italian debate on the Astrazeneca vaccine** and 18.2% for **Italian and Turkish EURO2020** dataset.

These are quite high values, especially if we take as a baseline measure the one provided by Varol et al.⁹⁸ in a 2017 study which estimated the percentage of active bots on the Twittersphere at between 9 and 15%.

However, in our research, several aspects could motivate both the high values and their variability amongst the datasets. Specifically, (1) we are looking at specific (hot) topics that might involve more significant numbers of bots than the average, (2) we are considering datasets on significantly different topics (thus, the percentage of automated accounts might vary), and (3) we are analyzing data collected in different time intervals, but evaluated with a single classifier (this might further affect the classification performance, due to the possible evolution of bots).

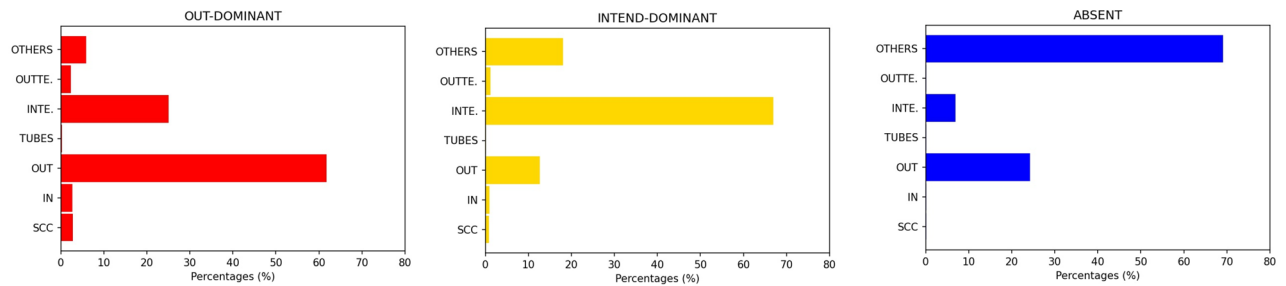


Figure 10. Percentages of social-bots in each sector of the bow-tie structure for OUT-dominant, INTEND-dominant and not informative bow-ties. This figure collects the percentage of bots in every bow-tie's sector for discursive communities with OUT-dominant, INTEND-dominant and not informative bow-tie structure. It is easy to note how the highest percentages can be found in the OUT sector, in the INTENDRILS and in the OTHERS one, respectively for the case of OUT-dominant, INTEND-dominant and not informative bow-tie.

With these premises to keep in mind, we now describe how the potential bots are distributed in the discursive communities. They are equally distributed among the discursive communities, with a slightly higher percentage of bots in the conservative groups: for instance, in the Italian and French COVID-19 datasets, the communities with the highest percentage of bots are DX and RIGHT-WING with, respectively, the 25.5% and the 29.7% of suspicious accounts. In our bow-tie structures, they are basically placed in the OUT sector or in the INTENDRILS one.

In Fig. 10 are shown the percentages of bots in a specific bow-tie sector averaged on all the discursive communities in the usual three categories. Globally, the highest percentages can be found in the OUT sector and in the INTENDRILS one: literally, social bots tend to retweet more than to be retweeted. In particular, in the case of OUT-dominant bow-tie, on average the 60% of bots are placed in the OUT sector and to a lesser extent in INTENDRILS (around 25%). In the case of INTEND-dominant bow-ties, we found even above the 60% of bots in INTENDRILS sector. Instead, when the bow-tie structure is absent, OTHERS is the block that contains the greatest number of bots.

It is worth to be mentioned is that the higher percentages of bots in the strongly connected component can be found in the right-oriented discursive communities. For instance, in the Italian COVID-19 dataset the percentage of bots in the SCC for the DX is the 7%, while for all the others it does not overcome the 2%. Such a situation is particularly dangerous, since the fact that social bots are able of being retweeted by human users (as it is the case for accounts in SCC) means that are able to pass off themselves as genuine accounts.

Data availability

The Twitter datasets used and analysed during the current study available from the corresponding author on reasonable request. The data about the reliability of the various news sources that support the findings of this study are available from Newsguard, but restrictions apply to the availability of these data, which were used under license for the current study, and so are not publicly available. Data are however available from the authors upon reasonable request and with permission of Newsguard.

Received: 7 February 2022; Accepted: 12 July 2022

Published online: 28 July 2022

References

- Adamic, L. A. & Glance, N. S. The political blogosphere and the 2004 U.S. election: Divided they blog. in *3rd International Workshop on Link discovery, LinkKDD 2005, Chicago, Illinois, USA, August 21-25, 2005*, 36–43 (2005).
- Commission, E. & For Communication, D.-G. *Media use in the European Union: Report* (European Commission, 2020).
- Dubois, E. & Blank, G. The echo chamber is overstated: The moderating effect of political interest and diverse media. *Inform. Commun. Society* **21**, 729–745. <https://doi.org/10.1080/1369118X.2018.1428656> (2018).
- Valensise, C. M. *et al.* Lack of evidence for correlation between COVID-19 infodemic and vaccine acceptance (2021).
- Gallotti, R., Pilati, F., Sacco, P. L. & Domenico, M. D. Comment on “The COVID-19 infodemic does not affect vaccine acceptance”. <https://doi.org/10.31219/OSF.IO/M8J32> (OSF Preprints).
- Urman, A. Context matters: Political polarization on twitter from a comparative perspective. **42**, 857–879. <https://doi.org/10.1177/0163443719876541> (2019).
- Yarchi, M., Baden, C. & Kligler-Vilenchik, N. Political polarization on the digital sphere: A cross-platform, over-time analysis of interactional, positional, and affective polarization on social media. 1–42. <https://doi.org/10.1080/10584609.2020.1785067> (2020).
- Barberá, P., Jost, J. T., Nagler, J., Tucker, J. A. & Bonneau, R. Tweeting from left to right: Is online political communication more than an echo chamber?. *Psychol. Sci.* **26**, 1531–1542. <https://doi.org/10.1177/0956797615594620> (2015) (PMID: 26297377).
- Lazer, D. M. J. *et al.* The science of fake news. *Science* **359**, 1094–1096. <https://doi.org/10.1126/science.aao2998> (2018).
- Gangware, C. & Nemr, W. *Weapons of Mass Distraction: Foreign State-Sponsored Disinformation in the Digital Age* (Park Advisors, 2019).
- Del Vicario, M. *et al.* The spreading of misinformation online. *Proc. Natl. Acad. Sci.* <https://doi.org/10.1073/pnas.1517441113> (2016).
- Jamieson, K. & Cappella, O. *Echo Chamber: Rush Limbaugh and The Conservative Media Establishment* (Oxford University Press, 2008).
- Garrett, R. K. Echo chambers online? Politically motivated selective exposure among internet news users. *J. Comput.-Mediated Commun.* **14**, 265–285. <https://doi.org/10.1111/j.1083-6101.2009.01440.X> (2009).
- Zollo, F. *et al.* Debunking in a world of tribes. *PLoS One* <https://doi.org/10.1371/journal.pone.0181821> (2017).

15. Pariser, E. *The Filter Bubble: What the Internet is Hiding From You* (Penguin Press, 2011).
16. Bruns, A. *Are Filter Bubbles Real?* (Wiley, 2019).
17. Borg, E. Discourse community. *ELT J.* **57**, 398–400. <https://doi.org/10.1093/elt/57.4.398> (2003).
18. Porter, J. *Audience and Rhetoric: An Archaeological Composition of the Discourse Community* (Prentice Hall, 1992).
19. Kehus, M., Kelley, W. & Melanie, S. Definition and genesis of an online discourse community. *Int. J. Learn.* **17**, 67–85 (2010).
20. Berkenkotter, C. A rhetoric for naturalistic inquiry and the question of genre. *Res. Teaching Eng.* **27**, 293–304 (1993).
21. Radicioni, T., Saracco, F., Pavan, E. & Squartini, T. Analysing twitter semantic networks: The case of 2018 Italian elections. *Sci. Rep.* **11**, 1–22. <https://doi.org/10.1038/s41598-021-92337-2> (2021).
22. Conover, M., Ratkiewicz, J. & Francisco, M. Political polarization on twitter. *Icswm* <https://doi.org/10.1021/ja202932e> (2011).
23. Conover, M. D., Gonçalves, B., Ratkiewicz, J., Flammini, A. & Menczer, F. Predicting the political alignment of twitter users. in *Proc.—2011 IEEE Int. Conf. Privacy, Secur. Risk Trust IEEE Int. Conf. Soc. Comput. PASSAT/SocialCom 2011*. <https://doi.org/10.1109/PASSAT/SocialCom.2011.34> (2011).
24. Conover, M. D., Gonçalves, B., Flammini, A. & Menczer, F. Partisan asymmetries in online political activity. *EPJ Data Sci.* <https://doi.org/10.1140/epjds6> (2012).
25. Raghavan, U. N., Albert, R. & Kumara, S. Near linear time algorithm to detect community structures in large-scale networks. *Phys. Rev. E Stat. Nonlinear Soft Matter Phys.* <https://doi.org/10.1103/PhysRevE.76.036106> (2007).
26. Becatti, C., Caldarelli, G., Lambiotte, R. & Saracco, F. Extracting significant signal of news consumption from social networks: The case of Twitter in Italian political elections. *Palgrave Commun.* (2019).
27. Cimini, G. *et al.* The statistical physics of real-world networks. *Nat. Rev. Phys.* **1**, 58–71. <https://doi.org/10.1038/s42254-018-0002-6> (2018).
28. Guarino, S., Mastrostefano, E. & Saracco, F. Discursive community detection on twitter. *In preparation* (2022).
29. Caldarelli, G., De Nicola, R., Del Vigna, F., Petrocchi, M. & Saracco, F. The role of bot squads in the political propaganda on Twitter. *Commun. Phys.* **3**, 1–15. <https://doi.org/10.1038/s42005-020-0340-4> (2020). arXiv:1905.12687.
30. Bruno, M., Lambiotte, R. & Saracco, F. Brexit and bots: characterizing the behaviour of automated accounts on twitter during the UK election. *EPJ Data Sci.* **11**, 1–24. <https://doi.org/10.1140/EPJDS/S13688-022-00330-0> (2022).
31. Patuelli, A., Caldarelli, G., Lattanzi, N. & Saracco, F. Firms' challenges and social responsibilities during COVID-19: A twitter analysis. *PLOS ONE* **16**, e0254748. <https://doi.org/10.1371/JOURNAL.PONE.0254748> (2021).
32. Radicioni, T., Squartini, T., Pavan, E. & Saracco, F. Networked partisanship and framing: A socio-semantic network analysis of the Italian debate on migration. *PLOS ONE* **16**, e0256705. <https://doi.org/10.1371/JOURNAL.PONE.0256705> (2021).
33. Caldarelli, G., Nicola, R. D., Petrocchi, M., Pratelli, M. & Saracco, F. Flow of online misinformation during the peak of the COVID-19 pandemic in Italy. *EPJ Data Sci.* **10**, 1–23. <https://doi.org/10.1140/EPJDS/S13688-021-00289-4> (2021).
34. Mattei, M., Caldarelli, G., Squartini, T. & Saracco, F. Italian twitter semantic network during the COVID-19 epidemic. *EPJ Data Sci.* **10**, 1–27. <https://doi.org/10.1140/EPJDS/S13688-021-00301-X> (2021).
35. Sluban, B., Smailović, J., Battiston, S. & Mozetič, I. Sentiment leaning of influential communities in social networks. *Comput. Social Netw.* **2**, 1–21. <https://doi.org/10.1186/S40649-015-0016-5/TABLES/6> (2015).
36. Cherepnalkoski, D. & Mozetič, I. Retweet networks of the European parliament: Evaluation of the community structure. *Appl. Netw. Sci.* **1**, 1–20. <https://doi.org/10.1007/S41109-016-0001-4/TABLES/3> (2016).
37. Uyheng, J. & Carley, K. M. Characterizing network dynamics of online hate communities around the COVID-19 pandemic. *Appl. Netw. Sci.* **6**, 1–21. <https://doi.org/10.1007/S41109-021-00362-X/FIGURES/6> (2021).
38. Evkoski, B., Mozetic, I., Ljubesic, N. & Kralj Novak, P. Community evolution in retweet networks. *PLOS ONE* **16**, 1–21. <https://doi.org/10.1371/journal.pone.0256175> (2021).
39. Evkoski, B., Pelicon, A., Mozetic, I., Ljubesic, N. & Kralj Novak, P. Retweet communities reveal the main sources of hate speech. *PLOS One* **17** (2022).
40. Broder, A. *et al.* Graph structure in the web. *Comput. Netw.* [https://doi.org/10.1016/S1389-1286\(00\)00083-9](https://doi.org/10.1016/S1389-1286(00)00083-9) (2000).
41. Yang, R., Zhuhadar, L. & Nasraoui, O. Bow-tie decomposition in directed graphs. 1–5 (2011).
42. Vitali, S., Glattfelder, J. B. & Battiston, S. The network of global corporate control. *PLOS ONE* **6**, e25995. <https://doi.org/10.1371/JOURNAL.PONE.0025995> (2011).
43. Mastrandrea, R., Squartini, T., Fagiolo, G. & Garlaschelli, D. Enhanced reconstruction of weighted networks from strengths and degrees. *New J. Phys.* **16**, 043022 (2014).
44. Artime, O., D'Andrea, V., Gallotti, R., Sacco, P. L. & De Domenico, M. Effectiveness of dismantling strategies on moderated vs. unmoderated online social platforms. *Sci. Rep.* **10**, 14392. <https://doi.org/10.1038/s41598-020-71231-3> (2020).
45. Guarino, S., Pierri, F., Giovanni, M. D. & Celestini, A. Information disorders during the COVID-19 infodemic: The case of Italian facebook. *Online Social Netw. Media* **22**, 100124. <https://doi.org/10.1016/J.OSNEM.2021.100124> (2021).
46. Castioni, P., Andrighetto, G., Gallotti, R., Polizzi, E. & Domenico, M. D. The voice of few, the opinions of many: Evidence of social biases in twitter COVID-19 fake news sharing (2021). arXiv:2112.01304.
47. González-Bailón, S. & De Domenico, M. Bots are less central than verified accounts during contentious political events. *Proc. Natl. Acad. Sci.* <https://doi.org/10.1073/pnas.2013443118> (2021).
48. Saracco, F., Di Clemente, R., Gabrielli, A. & Squartini, T. Randomizing bipartite networks: The case of the World Trade Web. *Sci. Rep.* **5**, 10595 (2015).
49. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large networks. *J. Stat. Mech. Theory Exp.* **10008**, 6. <https://doi.org/10.1088/1742-5468/2008/10/P10008> (2008).
50. Newman, M. E. J. & Girvan, M. Finding and evaluating community structure in networks. *Phys. Rev. E* **69**, 026113. <https://doi.org/10.1103/PhysRevE.69.026113> (2004).
51. Chung, F. & Lu, L. Connected components in random graphs with given expected degree sequences. *Ann. Comb.* **6**, 125–145. <https://doi.org/10.1007/PL00012580> (2002).
52. Squartini, T. & Garlaschelli, D. Analytical maximum-likelihood method to detect patterns in real networks. *New J. Phys.* <https://doi.org/10.1088/1367-2630/13/8/083001> (2011).
53. González-Bailón, S., Borge-Holthoefer, J. & Moreno, Y. Broadcasters and hidden influentials in online protest diffusion. *Am. Behav. Sci.* **57**, 943–965. <https://doi.org/10.1177/0002764213479371> (2013).
54. Squartini, T., Picciolo, F., Ruzzenenti, F. & Garlaschelli, D. Reciprocity of weighted networks. *Sci. Rep.* **3**, 1–9. <https://doi.org/10.1038/srep02729> (2013).
55. Squartini, T., van Lelyveld, I. & Garlaschelli, D. Early-warning signals of topological collapse in interbank networks. *Sci. Rep.* **3**, 3357 (2013).
56. Picciolo, F., Squartini, T., Ruzzenenti, F., Basosi, R. & Garlaschelli, D. The role of distances in the world trade web. 784–792, <https://doi.org/10.1109/SITIS.2012.118> (Institute of Electrical and Electronics Engineers (IEEE), 2013).
57. Gualdi, S., Cimini, G., Primicerio, K., Di Clemente, R. & Challet, D. Statistically validated network of portfolio overlaps and systemic risk. *Sci. Rep.* **6**, 39467 (2016).
58. Saracco, F. *et al.* Inferring monopartite projections of bipartite networks: An entropy-based approach. *New J. Phys.* **19**, 16. <https://doi.org/10.1088/1367-2630/aa6b38> (2017).
59. Saracco, F., Di Clemente, R., Gabrielli, A. & Squartini, T. Detecting early signs of the 2007–2008 crisis in the world trade. *Sci. Rep.* **6**, 30286. <https://doi.org/10.1038/srep30286> (2016).

60. Di Gangi, D., Lillo, F. & Pirino, D. Assessing systemic risk due to fire sales spillover through maximum entropy network reconstruction. *J. Econ. Dyn. Control* **94**, 117–141. <https://doi.org/10.1016/j.jedc.2018.07.001> (2018).
61. Squartini, T., Caldarelli, G., Cimini, G., Gabrielli, A. & Garlaschelli, D. Reconstruction methods for networks: The case of economic and financial systems. *Phys. Rep.* **757**, 1–47 (2018) (**Reconstruction methods for networks: The case of economic and financial systems.**).
62. Bardoscia, M. *et al.* The physics of financial networks. *Nat. Rev. Phys.* **3**, 490–507. <https://doi.org/10.1038/s42254-021-00322-5> (2021).
63. Straka, M., Caldarelli, G. & Saracco, F. Grand canonical validation of the bipartite international trade network. *Phys. Rev. E* <https://doi.org/10.1103/PhysRevE.96.022306> (2017).
64. Gabrielli, A., Mastrandrea, R., Caldarelli, G. & Cimini, G. Grand canonical ensemble of weighted networks. *Phys. Rev. E* **99**, 030301. <https://doi.org/10.1103/PhysRevE.99.030301> (2019).
65. Adam, I. *et al.* Maximum entropy approaches for the study of triadic motifs in the mergers & acquisitions network (2019).
66. Bruno, M., Saracco, F., Squartini, T. & Dueñas, M. Colombian export capabilities: Building the firms-products network. *Entropy* <https://doi.org/10.3390/e20100743> (2018).
67. Cimini, G., Mastrandrea, R. & Squartini, T. Reconstructing networks. *Elements Struct. Dyn. Complex Netw.* <https://doi.org/10.1017/9781108771030> (2021).
68. Vece, M. D., Garlaschelli, D. & Squartini, T. Gravity models of networks: Integrating maximum-entropy and econometric approaches (2021).
69. Lin, J.-H., Primicerio, K., Squartini, T., Decker, C. & Tessone, C. J. Lightning network: A second path towards centralisation of the bitcoin economy*. *N. J. Phys.* **22**, 083022. <https://doi.org/10.1088/1367-2630/ABA062> (2020).
70. Vallarano, N. *et al.* Fast and scalable likelihood maximization for exponential random graph models with local constraints. *Sci. Rep.* **11**, 1–33. <https://doi.org/10.1038/s41598-021-93830-4> (2021).
71. Straka, M. M. J., Caldarelli, G., Squartini, T. & Saracco, F. From ecology to finance (and back?): A review on entropy-based null models for the analysis of bipartite networks. *J. Stat. Phys.* **173**, 1252–1285. <https://doi.org/10.1007/s10955-018-2039-4> (2018).
72. Payrató-Borrás, C., Hernández, L. & Moreno, Y. Breaking the spell of nestedness: The entropic origin of nestedness in mutualistic systems. *Phys. Rev. X* **9**, 031024. <https://doi.org/10.1103/PhysRevX.9.031024> (2019).
73. Bruno, M., Saracco, F., Garlaschelli, D., Tessone, C. J. & Caldarelli, G. The ambiguity of nestedness under soft and hard constraints. *Sci. Rep.* **10**, 1–13. <https://doi.org/10.1038/s41598-020-76300-1> (2020).
74. Caruso, T., Rillig, M. C. & Garlaschelli, D. Fluctuating ecological networks: A synthesis of maximum entropy approaches for pattern and perturbation detection (2021).
75. Becatti, C., Caldarelli, G. & Saracco, F. Entropy-based randomization of rating networks. *Phys. Rev. E* **99**, 022306 (2019).
76. Parisi, F., Squartini, T. & Garlaschelli, D. A faster horse on a safer trail: Generalized inference for the efficient reconstruction of weighted networks. *N. J. Phys.* **22**, 053053. <https://doi.org/10.1088/1367-2630/AB74A7> (2020).
77. Neal, Z. P., Domagalski, R. & Sagan, B. Comparing alternatives to the fixed degree sequence model for extracting the backbone of bipartite projections. *Sci. Rep.* **11**, 1–13. <https://doi.org/10.1038/s41598-021-03238-3> (2021).
78. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *J. Royal Stat. Society B* **57**, 289–300 (1995).
79. Squartini, T. & Garlaschelli, D. *Maximum-Entropy Networks. Pattern Detection, Network Reconstruction and Graph Combinatorics* (Springer International Publishing, 2017).
80. Park, J. & Newman, M. E. J. Statistical mechanics of networks. *Phys. Rev. E* **70**, 066117. <https://doi.org/10.1103/PhysRevE.70.066117> (2004).
81. Garlaschelli, D. & Loffredo, M. I. Maximum likelihood: Extracting unbiased information from complex networks. *Phys. Rev. E Stat. Nonlinear Soft Matter Phys.* **78**, 1–5. <https://doi.org/10.1103/PhysRevE.78.015101> (2008).
82. Hong, Y. On computing the distribution function for the Poisson binomial distribution. *Comput. Stat. Data Anal.* **59**, 41–51 (2013).
83. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *J. R. Stat. Soc. B* **57**, 289–300 (1995).
84. Newman, M. *Networks: An Introduction* (Oxford University Press Inc, 2010).
85. Girvan, M. & Newman, M. E. J. Community structure in social and biological networks. *Proc. Natl. Acad. Sci. USA* **99**, 7821–6. <https://doi.org/10.1073/pnas.122653799> (2002). [arXiv:0112110](https://arxiv.org/abs/0112110).
86. Good, B. H., Montjoye, Y. A. D. & Clauset, A. Performance of modularity maximization in practical contexts. *Phys. Rev. E Stat. Nonlinear Soft Matter Phys.* **81**, 046106. <https://doi.org/10.1103/PhysRevE.81.046106> (2010).
87. Cresci, S., Di Pietro, R., Petrocchi, M., Spognardi, A. & Tesconi, M. Fame for sale: Efficient detection of fake twitter followers. *Decis. Support Syst.* **80**, 56–71 (2015).
88. Ferrara, E., Varol, O., Davis, C., Menczer, F. & Flammini, A. The rise of social bots. *Commun. ACM* **59**, 96–104 (2016).
89. Cresci, S., Di Pietro, R., Petrocchi, M., Spognardi, A. & Tesconi, M. The paradigm-shift of social spambots: Evidence, theories, and tools for the arms race. in *26th International Conference on World Wide Web Companion (WWW'17)*, 963–972 ACM, (2017).
90. De Nicola, R., Petrocchi, M. & Pratelli, M. On the efficacy of old features for the detection of new bots. *Inform. Process. Manag.* **58**, 102685 (2021).
91. Cresci, S. A decade of social bot detection. *Commun. ACM* **63**, 72–83 (2020).
92. Pal, S. K. & Mitra, S. Multilayer perceptron, fuzzy sets, and classification. *IEEE Trans. Neural Netw.* **3**, 683–697. <https://doi.org/10.1109/72.159058> (1992).
93. Cohen, W. W. Fast effective rule induction. In *Machine Learning Proceedings 1995* (eds Prieditis, A. & Russell, S.) 115–123 (Morgan Kaufmann, 1995). <https://doi.org/10.1016/B978-1-55860-377-6.50023-2>.
94. John, G. H. & Langley, P. Estimating continuous distributions in bayesian classifiers. in *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence, UAI'95*, 338–345 (Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1995).
95. Breiman, L. Random forests. *Mach. Learn.* **45**, 5–32 (2001).
96. Witten, I. H., Frank, E. & Hall, M. A. *Data Mining: Practical Machine Learning Tools and Techniques* 3rd edn. (Morgan Kaufmann, 2011).
97. Aha, D., Kibler, D. & Albert, M. Instance-based learning algorithms. *Mach. Learn.* **6**, 37–66 (1991).
98. Varol, O., Ferrara, E., Davis, C., Menczer, F. & Flammini, A. Online human-bot interactions: Detection, estimation, and characterization. in *Proceedings of the international AAAI conference on web and social media*, vol. 11 (2017).

Acknowledgements

FS acknowledge Pietro Galgani and Lizanne Dirx for support in both the download and the analysis of the Dutch election dataset; Giulia Andrighetto, Stefano Guarino, Enrico Mastrostefano, Elena Pavan, Eugenia Polizzi and Tiziano Squartini for useful discussions. All authors acknowledge support from IMT PAI project Toffee.

Author contributions

M.Pe. and F.S. planned the research. M.M., M.Pr. and F.S. performed the analyses. M.M. prepared all the figures. M.M., M.Pe., M.Pr. and F.S. wrote the main manuscript text. M.M. and F.S. wrote the Supplementary Materials. All authors reviewed the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-022-16603-7>.

Correspondence and requests for materials should be addressed to F.S.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022