

PAPER

Analysis of flows in social media uncovers a new multi-step model of information spread

To cite this article: Matteo Serafino *et al* *J. Stat. Mech.* (2024) 113402

View the [article online](#) for updates and enhancements.

You may also like

- [A novel translationally invariant supersymmetric chain with inverse-square interactions: partition function, thermodynamics and criticality](#)
Bireswar Basu-Mallick, Federico Finkel and Artemio González-López
- [Spontaneous symmetry breaking in two dimensions under non-equilibrium laminar flows](#)
Yuki Minami and Hiroyoshi Nakano
- [Random sequential adsorption and percolation on discrete substrates](#)
D Dujak, Lj Budinski-Petkovi and I Lonarevi

Analysis of flows in social media uncovers a new multi-step model of information spread

Matteo Serafino^{1,6,*}, Giulio Virginio Clemente^{1,2,6},
James Flamino³, Boleslaw K Szymanski^{3,4}, Omar Lizardo⁵
and Hernán A Makse¹

¹ Levich Institute and Physics Department, City College of New York, New York, NY 10031, United States of America

² IMT School for Advanced Studies, 55100 Lucca, Italy

³ Department of Computer Science and NEST Center, RPI, Rensselaer Polytechnic Institute, Troy, NY, United States of America

⁴ Społeczna Akademia Nauk, Łódź, Poland

⁵ University of California, Los Angeles, CA, United States of America

E-mail: matteo.serafino1991@gmail.com

Received 15 May 2024

Accepted for publication 14 October 2024

Published 8 November 2024



Online at stacks.iop.org/JSTAT/2024/113402

<https://doi.org/10.1088/1742-5468/ad8748>

Abstract. Since the advent of the internet, communication paradigms have continuously evolved, resulting in a present-day landscape where the dynamics of information dissemination have undergone a complete transformation compared to the past. In this study, we challenge the conventional two-step flow communication model, a long-standing paradigm in the field. Our approach introduces a more intricate multi-step and multi-actor model that effectively captures the complexities of modern information spread. We test our hypothesis by examining the spread of information on the Twitter platform. Our findings support the multi-step and multi-actor model hypothesis. In this framework, influencers (individuals with a significant presence in social media) emerge as new central figures and partially take on the role previously attributed to opinion leaders. However, this does not apply to opinion leaders who adapt and reaffirm their influential position on social media, here defined as opinion-leading influencers.

⁶These authors contributed equally to this work.

* Author to whom any correspondence should be addressed.

Additionally, we note a substantial number of adopters directly accessing information sources, suggesting a potential decline in influence in both opinion leaders and influencers. Finally, we found distinctions in the diffusion patterns of left-/right-leaning groups, indicating variations in the underlying structure of information dissemination across different ideologies.

Keywords: information diffusion, influencers, opinion leaders

Contents

1. Introduction	2
2. Results	4
2.1. Modeling information flow	5
2.2. Left vs. right	8
3. Discussion	9
4. Materials and methods	10
4.1. Data	10
4.2. Retweet network	11
4.3. Link validation	11
4.4. Influencer and opinion leader identification	13
4.5. Mapping the information flow: the BFS algorithm	14
Data, materials, and software availability	15
Acknowledgments	15
Appendix A. Validation	15
A.1. Temporal filtering	16
Appendix B. Identifying influencers	16
Appendix C. Further analysis	20
References	21

1. Introduction

The rise of social media has led to a fundamental transformation in how information, opinions, and beliefs propagate in contemporary society. Traditional models of information diffusion developed in sociology and communication in the mid-twentieth century [3, 18, 19, 21, 34] presupposed a linear ‘two-step’ flow of information from sources to the mass public mediated by ‘opinion leaders’ [7, 21, 26, 32], defined as recognized experts or respected public figures with acknowledged credibility in specific fields (see

figure 1(a)). For example, within the two-step model (figure 1(a)), *The New York Times* (the source) releases a news article on the evolution of R_0 , a key epidemiological metric for measuring infectious agent transmissibility. Dr Anthony Fauci, acting as an opinion leader, reads and simplifies the news (S1: first step) before disseminating it to the broader population in a more accessible manner. These individuals, termed adopters, engage with or adopt the idea at various stages (S2: second step).

In the current scenario, by way of contrast, information diffuses via more complex *multi-step* flows, including one-step flows with adopters accessing information directly from the sources [1], without intermediaries, traditionally mediated by two-step, and more complex, longer-path dynamics featuring a heterogeneous set of agents. For instance, adopters may obtain information from other adopters, who, in turn, receive the information from an opinion leader or other mediators (see figure 1(b)). Notably, a new figure has emerged within the intricate structure of contemporary digitally mediated information diffusion: the *influencer*. Unlike traditional opinion leaders, influencers often build their authority through a combination of relatability, engaging content, and a substantial online presence [9].

The rise of influencers has added a new layer to the landscape of information diffusion, introducing a dynamic where individuals with significant followings can swiftly impact trends and opinions. In the context of COVID-19, Elon Musk can be considered a notable influencer who wields considerable social influence that can shape public opinion. Musk's impact is not necessarily rooted in expertise in infectious diseases or pandemics but rather in his extensive reach across online social platforms. Of course, opinion leaders and influencers need not be mutually exclusive groups. Instances exist where traditional opinion leaders have effectively established themselves as influencers. We refer to these individuals as opinion-leading influencers. One notable example is Helen Branswell, a Canadian infectious diseases and global health reporter at *Stat News*. With a 15-year tenure as a medical reporter at *The Canadian Press*, she spearheaded Ebola, Zika, SARS, and swine flu pandemic coverage. Beyond her field of expertise, Branswell maintains a robust online presence, qualifying her as an opinion-leading influencer.

While most observers agree that the traditional opinion-leader-mediated two-step dynamic has certainly been disrupted [1, 10, 12, 17, 31, 33, 35, 36], we know little about the new pathways and actors through which information diffuses in online social media, as well as the extent to which opinion leaders, influencers, or the new hybrid figure of the opinion-leading influencer still serve as key mediators. Or whether individuals have become direct consumers of information from sources, or the extent to which horizontal transmission among adopters, defined here as person-to-person transmission, in contrast to the top-down (vertical) transmission from opinion leaders to individuals, accounts for the bulk of information flow. Despite considerable speculation about how social media has transformed information diffusion, there is still a need for quantitative studies that allow us to clarify from different perspectives the extent to which information flow on digital platforms is mediated through multiple steps. This paper sets out to reconstruct the pathways through which information flows in the era of social media, to characterize how information diffuses through different groups of actors and to ascertain whether the decline in the influence of opinion leaders [1, 2] has been greatly exaggerated or not.

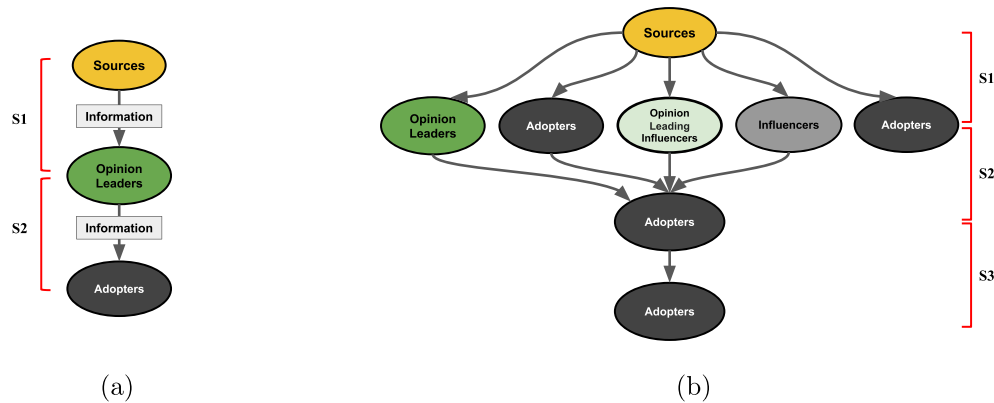


Figure 1. Two information diffusion models. (a) The traditional two-step model of information diffusion involves the flow of information to adopters through the mediation of opinion leaders. (b) The multistep model for information diffusion. Adopters can directly access the original information or obtain it through the more traditional opinion leaders or influencers. Note the possibility of ‘horizontal information flow’, where adopters receive information from other adopters.

For this task, we turn to Twitter content, using a dataset centered around the 2020 US presidential elections [14]. Digital platforms like Twitter provide researchers with the tools to track specific uniform resource locators (URLs) released by the sources. This tracking unveils insights into users who share the same URL, offering a detailed account of user interactions over time. This method enables us to directly assess a proxy for the channels of information flow among various actors, encompassing the source, opinion leaders, influencers, opinion-leading influencers, and ultimately, adopters. We leverage this distinctive capability to disclose the structural characteristics and dynamics initiated by various actors in the digital space.

2. Results

Our analysis proceeds as follows (also refer to appendix, figure A1). Initially, from the raw data, we construct the retweet network by considering tweets that include a URL linking to one of the news outlets listed in the appendix, table A1. Considering only tweets with URLs allows us to trace back the original source of the information, a key aspect in understanding the diffusion information process explained in the introduction. In the retweet network, a connection exists from user i to user j if j has retweeted a tweet from i . The connection is weighted by the number of times j has retweeted a tweet from i . Additionally, we assign an average retweet time to each link, calculated as the average retweet time among the retweets of i by j . However, this process results in the loss of the temporal dimension of the interactions. We then validate these links against an entropy-based null model. The validation process is designed to preserve only significant connections, thereby reducing potential biases in subsequent stages of our methodology. This step is essential due to the original structure of the data, which does not allow for direct measurement of the patterns of interest (section 4). Next,

using the collective influence (CI) algorithm (see appendix B), we identify the top 1000 influencers, representing 0.1% of the users in the network and accounting for more than 65% of the total connections. We classify them into one of the following categories: opinion leaders, influencers, opinion-leading influencers, adopters, and sources (refer to section 4.4 and table 1 for more information on the classification process). It is important to note that, while the choice of considering the top 1000 influencers according to CI is arbitrary, increasing this number does not significantly alter the final results of our analyses. This is because the remaining nodes (which account for more than 99% of the total number of nodes) contribute to less than 35% of the remaining links. Finally, we apply the breadth-first search (BFS) algorithm to the validated network to uncover the underlying structure most likely to facilitate information propagation. In this context, ‘most likely’ refers to the frequency of retweets between groups of nodes, without accounting for the temporal dimension, which is not preserved when constructing the retweet network. We refer to this extracted structure as the ‘backbone’. It represents the skeleton of information diffusion, meaning that if news is shared on Twitter, it would most likely spread through the identified skeleton (or one of its subgraphs). As detailed in appendix A.1, the results remain robust when filtering the links based on their average retweet time. Further analysis, provided in appendix C, contributes to validating the results. Finally, we leverage the tendency of each news outlet to exhibit a bias toward either a left or right ideology. This enables a more in-depth analysis of the diffusion structure and highlights potential differences between diverse ideological perspectives.

2.1. Modeling information flow

Here, our emphasis is on reconstructing the flow of information diffusion originating from the sources documented in the appendix, table A1, regardless of the news outlet bias under consideration. The resulting retweet network consists of 2963 210 nodes and 27 608 480 unique connections. The link-validation procedure validates 51% of the total links, with the validated network consisting of 1 775 194 nodes and 14 258 411 connections. The decrease in the number of nodes is because some nodes are isolated after validation and, therefore, discarded. In this graph, we identify 1173 opinion leaders (politicians or users affiliated with the journal under consideration, see appendix, table A1), 241 sources, 520 opinion-leading influencers, 399 influencers, and 1 772 869 adopters. Refer to section 4.4 for further details of classifications, and to the section 4.6 for a comprehensive list of classified sources and opinion leaders. We consider influencers the users within the top 1000 users by CI who are not opinion leaders. Indeed, the latter are opinion-leading influencers. However, in theory, the sum of these two categories should total 1000, practical deviations occur due to the inclusion of some sources, which also fall within the top 1000 users by CI. Moreover, notice here that despite the number of the initial news outlets consist of 69 elements; here, when we refer to sources we consider all the accounts associated with one of those news outlets. For example, CNN is associated with @CNN but also with @CNNPolitics. Table A2 in the appendix presents the top 10 influencers, opinion-leading influencers and opinion leaders.

After identifying the main actors, we proceed to unveil the skeleton of the information flow using the breadth-first search (BFS) algorithm (see section 4). Apart from the final network structure, to support the implementation of this procedure, we conducted a robustness check to examine the connections directed towards the adopters (see appendix C). These connections can be classified in two ways in relation to the BFS-derived structure. The first classification pertains to connections that contradict the directions identified by the BFS algorithm, and the second relates to connections that may link from step 1 to step 2 in the identified structure. In both cases, however, we can assume that the impact of such connections is negligible, as they account for much less than 1% of the total number of connections. The outcome is illustrated in figure 2. The normalization of connections is computed per step, ensuring that the percentage of connections in each step adds up to 100%. The resulting skeleton comprises 1 718 201 nodes and 5306 961 links.

In the initial step (S1), a substantial group of adopters directly access information from the sources without any mediator. Less than 1% of the connections in this step are directed to opinion leaders, opinion-leading influencers, and influencers (dashed arrows in the figure). Indeed, adopters in this step (S1) account for more than 99% of the nodes accessing the information directly from the sources. Overall, this step accounts for 5% of the total number of nodes and 2.5% of the connections in the skeleton found employing the BFS algorithm. The significant difference between adopters and the other main actors derives from the substantial size gap, with the adopter group being tens of times larger than the other three groups, each of which is of the same order of magnitude. Within the adopters identified in S1, 75% are active adopters, meaning they are directly involved in subsequent steps of the flow and have retweeted in S2, as indicated in the inset of figure 2. Moreover, 25% of the adopters (non-active adopters) defined in S1 serve as information sinks, conforming to a one-step model structure for information diffusion, as illustrated by the orange cloud in the inset.

As the diffusion process progresses (S2 in figure 2), 55% of the nodes (80% of the total links) in the skeleton access the information through mediators. Traditional opinion leaders only account for 10% of the connections, indicating a significant lower influence compared to other groups. Opinion-leading influencers, on the other hand, due to their adaptability in the online community, still wield strong influence, accounting for 31% of the connections. Similarly, influencers position themselves with significant influence, mediating 27.5% of the connections. These results are consistent with a two-step model structure for information diffusion, where mediators encompass not only traditional opinion leaders (indicated by the magenta cloud in figure 2) as originally formulated, but also a diverse set of actors, including influencers and opinion-leading influencers. We observe three additional ways to construct a two-step model by substituting opinion leaders with opinion-leading influencers, influencers, or adopters. This forms the basis of a multi-actor model. Additionally, adopters also function as mediators, facilitating the information transfer to other adopters ('horizontal information flow') and accounting for 31.2% of the connections in S2.

As depicted in figure 2, the information diffusion extends beyond S2. The structure presented in the third step (S3) is similar to the one in S2 and accounts for only 39% of the nodes in the skeleton (16% of the links). The cyan cloud in figure 2 represents one of

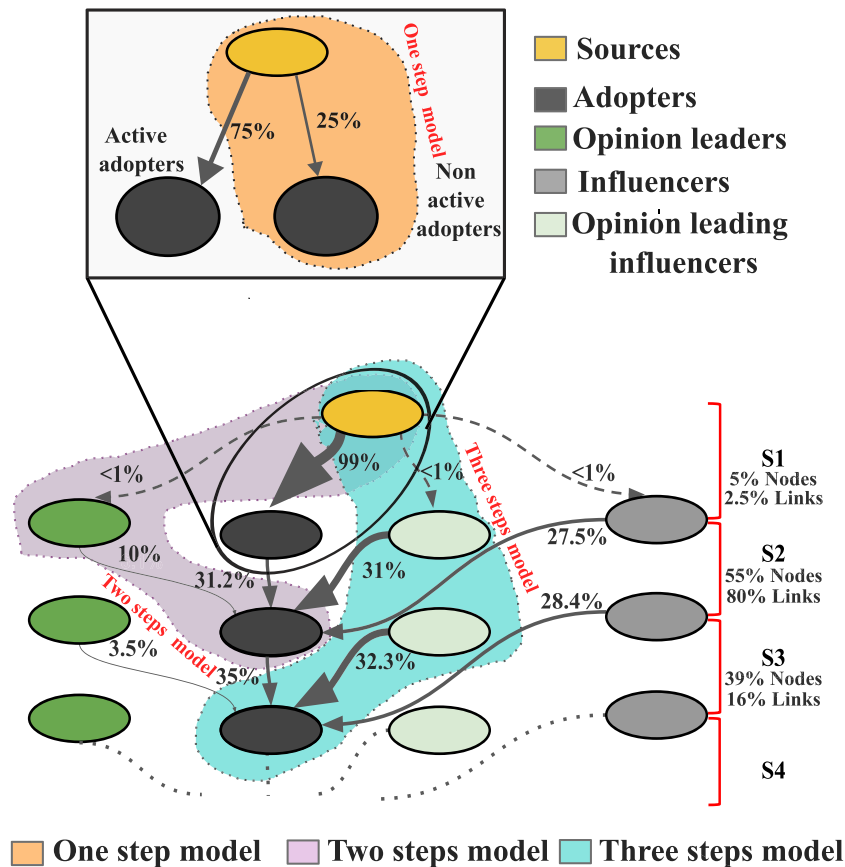


Figure 2. Information flow backbone. Skeleton of the information flow. We define active adopters as individuals directly involved in subsequent steps of the flow, in contrast to non-active adopters. The steps of information flow are highlighted on the right, indicating steps (S1, S2, S3, S4) along with the percentage of corresponding node and connection counts. Overlapping groups are non-existent within the same step and between different steps. Step one (S1) primarily involves mediators (opinion leaders, influencers, opinion-leading influencers, and adopters) retweeting the sources. Step two (S2) consists of adopters retweeting mediators from S1. This pattern is repeated until S9. The final six steps make up the remaining 1% of nodes and 1.5% of the links in the skeleton. Not all steps are displayed for clarity, with an incomplete step S4 alluding to its continuation. Connection normalization per step ensures that the percentage on each layer adds up to 100. Throughout the paper, consistent color associations for the main actors are maintained.

the four potential three-step-like structures for information diffusion. Importantly, from S2 onward, only links representing at least 1% of the connections in each step are displayed for visual clarity. Therefore, connections between the opinion-leading influencers between S2 and S3 in the cyan cloud exist, even if not explicitly shown.

This pattern is repeated until S9, with the remaining six steps comprising the remaining 1% of the nodes and 1.5% of the connections in the skeleton. We do not display all

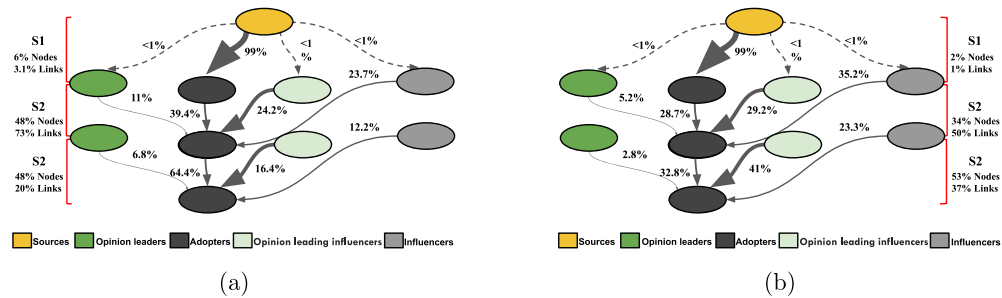


Figure 3. Backbone: left vs. right. (a) Backbone obtained by following the steps depicted in appendix, figure S2 but by only considering news coming from left and left-leaning news outlets. (b) Backbone obtained by following the steps depicted in appendix, figure S2 but by only considering news coming from right and right-leaning news outlets. While stopping at the third step of the information flow, we confirm that the patterns we observed in the general case still hold. Throughout the paper, consistent color associations for the main actors are maintained, as indicated at the bottom of the figure.

the steps in figure 2 for the sake of clarity. However, we leave an incomplete step S4 to allude to its continuation.

2.2. Left vs. right

In this section, we investigate whether the information flow linked to various political media biases displays distinct characteristics. Specifically, we focus on identifying potential differences between content related to the left and right political spectrum. We aim to uncover key disparities in the structure of the information flow and the roles played by primary actors in information dissemination. Our analysis focuses exclusively on left-leaning sources (left and left-leaning in appendix, table A1) for studying the left and right-leaning (right and right-leaning in the appendix, table A1) news outlets for investigating the right.

By following the same steps as in the previous analysis, we obtain a skeleton comprising 475 636 nodes and 902 189 edges for the right-leaning. This represents a reduction of 58% of the nodes and 87% of the edges present in the original retweet network. In the case of left-leaning sources, we end up with a skeleton consisting of 710 432 nodes and 2027 887 edges, indicating a reduction of 71% of the nodes and 90% of the edges present in the original left retweet network.

Figure 3 displays the backbones resulting by considering as sources of information only the left-leaning news outlets (on the left) and the right-leaning news outlets (on the right). The figures stop at the third step of the information flow, which accounts for more than 96% of the nodes and 95% of the links in the skeleton (in both cases). The structure of the consecutive layers is similar to the one shown in figure 2, in agreement with the hypothesis of a multi-step and multi-actor model for the diffusion of information, despite the political polarization of the information circulating.

The initial distinction observed between left and right lies in the depth of the information diffusion process. In the case of the left-leaning sources, the first two steps of this

process encompass nearly 55% of the nodes and 76% of the connections. In contrast, these percentages decrease to 36% and 51% for the right-leaning sources, as illustrated in figure 3.

Another intriguing difference highlighted in figure 3 pertains to the role of influencers in the second step (S2) of the information diffusion process. Specifically, our analysis reveals that, in the case of the right-leaning backbone, 35.2% of the connections in S2 originate from adopters who retweeted influencers. In contrast, opinion leaders account for only 5.2% of the connections, while opinion-leading influencers play a substantial role, being retweeted 29.2% of the time. Adopters rank third in terms of the frequency with which they are retweeted, following opinion-leading influencers and influencers.

This hierarchical pattern undergoes a shift in the left-leaning backbone. Here, opinion-leading influencers are retweeted 24.2% of the time in S2, followed by influencers at 23.7%, and adopters contributing for 39.4% to the total retweets. Opinion leaders, in this case as well, emerge as the group with the least impact on the spread of information.

3. Discussion

The rise of social media has fundamentally transformed how information, ideas, and opinions diffuse in contemporary society. Whether traditional models designed in the social sciences to understand this phenomenon, developed in the context of legacy media forms in the middle of the twentieth century still apply, or whether we need a completely different way of thinking about the issue remains a live question.

We reconstructed the multi-step flow of information on a major social media platform (Twitter) at the height of its influence. Our analysis shows that the current structure of information diffusion in the social media era displays characteristics of multiple combined processes. These include both unmediated one-step flows where users connect directly to authoritative information sources, accounting for a small percentage of the total activity. Mediated flows are far from insignificant, particularly those that go via influencers, whether of the traditional opinion leader or the more novel social media influencer kind. This means that the traditional two-step model still helps make sense of this dynamic, indicating that a lot of information flow in social media platforms is curated and mediated by key actors, intervening between sources and potential adopters. Longer multi-step flows can also be observed. Finally, horizontal information flows among adopters, partially independent of accredited or social-media-based mediators, also account for a significant portion of the information flow. This dynamic may be unique to the flow of information in digital platforms. The structure of information flow in the current system is thus closer to a mixed regime, displaying dynamics differing in length, form of intermediation, and the type of actors involved. We also observe intriguing differences in the prevalence of different elements of this hybrid regime across the right/left political divide. In particular, we find that social-media-based influencers, instead of opinion leaders, leave an increased footprint in shaping public opinion regarding right-related news than the left case at earlier steps (see figure 3).

Our work settles the question regarding the death of intermediation and the decline of mediation and opinion leadership in the social media era. While novel ways of accessing information and distinct pathways of mediated information diffusion have indeed opened up, opinion leadership is far from irrelevant. Nevertheless, traditional opinion leaders face significant competition from actors whose source of influence is social media reach. Only opinion leaders who themselves adopt the influencer strategy may be able to become significant intermediators in the social media-driven ecology.

It is worth noting that, by design, the backbone resulting from the BFS algorithm is only an approximation of the information diffusion network, intended to propose a structure through which information is most likely to diffuse. The steps involved in constructing and analyzing this backbone rely on aggregations and averages, aiming to reconstruct a plausible proxy for the diffusion of politics-related news. However, this approach loses the temporal dimension, which is a critical aspect that could reveal interesting differences in diffusion times as information passes through different types of mediators. Moreover, it is important to consider that our analyses are performed exclusively on Twitter data. As a result, the findings may not necessarily be applicable to other social platforms, which could have different network structures, user behaviors, or content-sharing dynamics. Given the increasing heterogeneity of social media systems in the post-Twitter era, future work should adopt a comparative strategy across the plethora of emerging platforms to investigate whether intermediation dynamics of information flow differ systematically, both cross-platform and across language, cultures, topics, and even political divisions, as we observe in this study. The framework we develop in this paper can be readily adapted and scaled for such comparative studies. It is possible that different platforms may encourage their own unique signature combination of one-step, two-step, and multistep information flows, including more horizontal information flows in platforms less dominated by influencer and opinion-leadership dynamics. A deeper understanding of how particular policies and design decisions of different platforms shape the particular structure of information flow within them can provide insights for developing strategies for effective communication and information campaigns at scale.

In the middle of the twentieth century, scholars in sociology and mass communication studies could imagine relatively simple two-step dynamics in which opinion leaders controlled the flow of information to the general public. We are unlikely to ever go back to that world. Nevertheless, the core insight that information mediation is an important phenomenon, one that can co-exist with other ways of both accessing and learning about news, ideas, and opinions, is one that will continue to be central in the era of social media.

4. Materials and methods

4.1. Data

We tracked the spread of political news on Twitter in 2020 by analyzing a dataset containing tweets posted between 1 June and election day (2 November 2020). The data were collected continuously using the Twitter search API with the names of the

two presidential candidates as keywords. The 2020 dataset contains 702 million tweets sent by 20 million users [4, 5, 13, 37].

To control for information polarization [13], we consider tweets containing at least one URL link directing to a news media outlet in a curated list of media outlets. The news outlet classification relies on the website all-sides.com (AS, accessed on 7 January 2021). We classified URL links for outlets that mostly conform to professional standards of fact-based journalism in five news media categories: right, right-leaning, left-leaning, and left. The classifications ('left' and 'right') of media outlets used are subjective and sourced from publicly available datasets by fact-checking organizations. A detailed explanation of the methodologies used by AS for rating news outlets is given in [4, 5, 13, 37]. A full list of the outlets in each category can be found in the appendix, table A1. These news outlets represent the sources of information of the information diffusion model. The dataset under study contains 72.7 million tweets with news links from one of these news outlets sent by 3.7 million users.

4.2. Retweet network

The initial phase of investigating the real-world system involves defining the retweet network, serving as a schematic representation of the diffusion of political opinions (appendix, figure A1(a)). The network is constructed by considering retweets containing a URL leading to one of the news outlets introduced above. Two users (i and j) are connected if one has retweeted the other at least once [25]. Link directions follow the flow of information, with the link going from i to j if j has retweeted a tweet from i . The resulting network is both directed and weighted, with weights denoted by the variable w , representing the number of times user j has retweeted user i . Furthermore, since each tweet is timestamped, we calculate the average retweet time between two nodes. Before proceeding, it is necessary to address critical considerations to set the stage for subsequent steps that we must take to operate on the original retweet network. The Twitter data do not allow us to directly construct the diffusion cascade. Consider this scenario: user 0 posts news, marked by a specific URL. User 1 retweets this post directly from user 0's tweet. Subsequently, user 2 retweets from user 1's post. Ideally, the data for user 2's tweet should cite user 1 as the source. However, the system identifies user 0's original tweet as the source instead. This pattern repeats with subsequent users, resulting in each tweet pointing back to user 0's original post, thereby creating a star graph. Consequently, the final retweet network, composed of multiple star graphs, fails to accurately represent the actual diffusion pathways. To address this, we have developed a strategy that involves a validation and mapping process designed to reconstruct an 'average' cascade structure. This approach allows us to identify the most relevant configuration of the diffusion process.

4.3. Link validation

This step aims to preserve only statistically significant connections [29]. The presence of numerous non-statistically significant links could compromise our results when determining the optimal model for information diffusion. For instance, if many adopters have connections with a source with a weight of one, while only a few connections exist with

opinion leaders with weights greater than one, considering all connections might incorrectly suggest the one step model as the best description of information flow. However, upon statistical validation, the weight-one links would be eliminated, emphasizing the connection to opinion leaders as the most significant one, favoring the two-step model. Validation ensures that the observed connections are not random but are influenced by the shift in the communication paradigm defined by the multi-step model. To validate the structure and assign statistical significance to the observed multi-step model structure, we employ null models. Using null models helps us determine whether a connection between two nodes is unexpected, potentially introducing misleading information, or whether it is expected, indicating a meaningful flow of information between the two nodes. Therefore, selecting the appropriate null model is essential to test the properties considered relevant and to adequately address the research question [15]. Hence, the pertinent question becomes: Given the network we are examining, is it typical for a node i with an out-strength of s_{out} and a node j with an in-strength of s_{in} to be connected? Here, $s_{\text{in}}(i) = \sum_j w_{ij}$ and $s_{\text{out}}(i) = \sum_j w_{ji}$.

To accomplish this task, we employ maximum-entropy models, a versatile class of models that can incorporate fluctuations in measurements [8], thereby enhancing pattern detection quality [6]. These models assume different expressions based on the specific constraints to be reproduced. Although analytical solutions for these models are rarely available, significant progress has been made in addressing this challenge. Various models have been developed, ranging from those suited for bipartite networks [27] to time-varying graphs [11]. Specifically, we leverage the conditional reconstruction method, a maximum-entropy ensemble model [24], for its proficiency in accurately replicating observed system topologies while permitting weight randomization. By doing so, it evenly distributes weights across all available links. Our objective is to determine whether the observed weight of a connection significantly deviates from the average predicted by the ensemble. If the observed weight is markedly lower, we may consider severing that link. We retain the in- and out-strengths during randomization because these metrics could reflect the characteristics of the nodes themselves rather than being inherently linked to the connection. For instance, while some individuals might be more inclined to retweet or be retweeted, we aim to preserve this information. However, we simultaneously control whether the existence of a retweet to or from a specific individual can be justified.

Moreover, beyond validating based on the weight of connections, we also investigate whether the final inferred shape of the network, obtained using the next step of our strategy, changes when we apply additional filtering based on the rapidity of retweets. The rationale is to determine if removing significantly slow links, presumed to result from retweets that occurred with considerable delays and, therefore, likely reached new users through a chain of intermediaries, might distort the link back to the original source of information. To this end, we apply two different filters to the validated networks: the first retains only links that occur within 75% of the time distribution and the second within 20% of the distribution. In both scenarios tested, we observed no significant deviations between the structures with and without temporal filtering. This finding suggests that, for the primary purpose of mapping information diffusion by focusing on groups of nodes rather than individual elements, the BFS algorithm applied to the

original validated networks effectively captures the top fastest connections. Based on these observations, our results represent the general case without the need for temporal filtering.

The validation process, by mainly exploiting the weight of the connections (appendix, section A.1), enables us to concentrate on edges that hold more information or significance within the retweet network, offering a more precise representation of the underlying structure through which information propagates.

4.4. Influencer and opinion leader identification

To determine the most suitable information propagation model for Twitter, we need to identify the actors of the model (appendix, figure A1(c)): sources, opinion leaders (sometimes referred to as traditional opinion leaders), influencers, and the overlap between opinion leaders and influencers, termed opinion-leading influencers.

In order to identify opinion leaders, we examine the URL field in the user's Twitter object. Journalists/reporters associated with a news outlet often include a link redirecting to their outlet in their bio. Therefore, we consider users linking to one of the news outlets classified as sources as opinion leaders. We also consider politicians' profiles as opinion leaders. This set of users undergoes manual verification. Please refer to section 4.6 for a comprehensive list of classified sources and opinion leaders.

To identify the influencers in the validated retweet network, we use the CI algorithm [13, 23, 28] (see appendix B). This widely recognized metric identifies nodes whose removal could disrupt the giant connected component, influencing information diffusion. We select the top 1000 individuals with the highest CI scores among users with non-zero CI values. The top 1000 influencers alone account for more than 85% of the interactions in the network. To check for users indirectly associated with news outlets, a secondary check is performed on the top 1000 influencers identified by CI. Each influencer is classified as an opinion-leading influencer if it is an opinion leader. Otherwise, the user is labeled as an influencer. Users not belonging to sources, opinion leaders, influencers, or opinion-leading influencers are labeled as adopters. By definition, these groups of users have an empty intersection.

As an example, when examining the retweet network derived from utilizing all sources independently, we observe that Donald J. Trump (former US president), Joe Biden (current US president), and Natasha Bertrand (CNN reporter) are identified as opinion-leading influencers. Similarly, Jonathan Landay (Reuters reporter) and Rick Tyler (political analyst at MSNBC) are recognized as opinion leaders, while Donald Trump Jr and Eric Trump are categorized as influencers. Refer to the appendix, table A2 for details. See Table 1 for a synthetic description of the user categories and how they are identified. It is worth noting that an alternative approach for identifying opinion leaders could have been to use the verification badge provided in Twitter's bio information. However, we found that many journalists do not have this verification badge, which would result in the exclusion of many traditional opinion leaders. Additionally, recent changes to Twitter's policies allow users to purchase verification badges, making this distinction less reliable. Despite this, as mentioned earlier, the top 1000 influencers (0.1% of the total users) identified by CI account for more than 65% of the connections

Table 1. Description of each category considered in this study is provided. Accounts described as opinion leaders, influencers, and opinion-leading influencers have been manually reviewed.

Category	Description
Sources	Accounts directly linked to a curated list of news outlets (refer to the appendix table A1).
Influencers	Users who are among the top 1000 most influential nodes in retweet networks, as determined by CI, but are not considered opinion leaders.
Opinion leaders	Users recognized as experts or respected public figures with acknowledged credibility in specific fields, identified here as journalists or other political figures who directly link to one of the considered news outlets in their descriptions.
Opinion-leading influencers	Opinion leaders, as defined above, who are among the top 1000 most influential users according to the CI. Among them, there are also well-recognized politicians or figures that can be directly associated with a political orientation.
Adopters	Users not included in any of the aforementioned categories.

in the network. Since these users were manually checked and labeled, we are confident that the most important actors, in terms of network structure, are included in this list. In other words, while we may miss some opinion leaders or influencers, those omitted have minimal impact on the network structure and, therefore, would not significantly affect our results.

4.5. Mapping the information flow: the BFS algorithm

To identify the information diffusion model that best characterizes the Twitter information diffusion network, we employ a BFS algorithm. The exploration begins with users classified as sources, serving as the root nodes in the BFS algorithm. We examine all the first neighbors of these sources, distinguishing this set of users into influencers, opinion-leading influencers, opinion leaders, and adopters. This initial step identifies the first step (S1) of the information diffusion model (appendix, figure A1(d)). At this stage, we consider only connections among users with a weight above one.

Next, we consider all the first neighbors of the newly identified nodes (the first neighbors of the first neighbors), making sure not to select users already chosen in the previous step. This step identifies the second step (S2) in the information diffusion process. Each iteration of the above procedure adds a new step to the information

diffusion process and our algorithm halts when no additional neighbors for the nodes defined in the earlier steps are available.

Data, materials, and software availability

The Twitter data and codes can be accessed at the following link: <https://osf.io/u9svz/>.

Acknowledgments

HAM and MS were supported by National Science Foundation Grant NSF-SBE award 2214217. BKS was partially supported by DARPA under contract HR001121C0165, as well as from the National Science Foundation Grant NSF-SBE 2214216. OL was partially supported by National Science Foundation Grant NSF-SBE 2214216.

Appendix A. Validation

We adopt a reference model constructed specifically to validate the connections between nodes in our network. This model is designed to preserve the topology of the retweet network, as well as the expected value of the total number of retweets made and received by each node. The resolution to an analogous challenge is detailed in a study where the authors introduce the CReMa (Conditional Reconstruction Method Model A) [24]. The CReMa model allows to define a probability distribution over a set of graphs, that effectively replicates both the topology and the expected values of the network's in and out strength sequences. Additionally, it ensures that all other network observables are maximally random, and can be either analytically or numerically calculated [30]. The main tool employed in defining the model relies on the fundamental principle of entropy maximization [16].

In our case, we estimate node-specific parameters $\vec{\beta}_{\text{in}}$ and $\vec{\beta}_{\text{out}}$ that are intrinsic to the model, correlating directly with the count of incoming and outgoing retweets for each node. This process allows us to calculate the expected weight (w_{ij}) of each link (how many times node i has retweeted node j) as well as its standard deviation. This enables the definition of the probability of observing the actual weight, considering only node-specific characteristics and not those of the individual link.

Specifically, for the employed model, the expected value and the standard deviation are the same and have the following expression:

$$\sigma(w_{ij}) = \langle w_{ij} \rangle = \frac{1}{(\beta_{\text{in}}(i) + \beta_{\text{out}}(j))}, \quad (\text{A.1})$$

from here we can attach a z -score to each link [15] that is:

$$z(w_{ij}) = \frac{w_{ij}^* - \langle w_{ij} \rangle}{\sigma(w_{ij})}. \quad (\text{A.2})$$

The parameters used to define the probability of connection in the referenced null model are obtained by solving the following 2 N set of coupled equations (where N is the number of nodes):

$$\begin{cases} \sum_{j(\neq i)} \frac{a_{ij}^*}{\beta_i^{\text{out}} + \beta_j^{\text{in}}} = s_i^{\text{out}*}, & \forall i \in N \\ \sum_{j(\neq i)} \frac{a_{ji}^*}{\beta_j^{\text{out}} + \beta_i^{\text{in}}} = s_i^{\text{in}*}, & \forall i \in N \end{cases} \quad (\text{A.3})$$

where $s_i^{\text{out}*}$ and $s_i^{\text{in}*}$ represent, respectively, the total number of times node i has been retweeted and the number of times node i has retweeted. a_{ij}^* is the ij entrance of the empirical adjacency matrix, that in our case is kept fixed. We chose to randomize while keeping the network's topology fixed because, in this case, the topology of the retweet network crucially represents the propagation of information among users. Randomizing by introducing connections between users who have never retweeted each other would distort this structure. Therefore, we opted to preserve the integrity of the original connections. For further details on the model and calculations details, see [24].

After calculating the z -score, it is possible to establish a threshold. With this threshold, we can assess whether the observed actual value can be accepted or rejected based on the subjective evaluations. In our analysis, we have chosen a threshold of -1. Therefore, we accept all links that have a z -score within the range $[-1, +\infty]$.

A.1. Temporal filtering

We introduce an additional step to incorporate the temporal dimension into our analysis. After validating the connections by assessing the frequency of retweets among users, this final step involves testing and applying a final filter based on retweet timing before executing the BFS algorithm. For each connected pair of nodes in the retweet network, we calculate an average retweet time from the time differences between the original tweet and its retweets, each timestamped. We suppose that if the average retweet time between two users, i and j , is longer than that involving another user, k , it indicates that user i typically accesses information from user j before user k . Building on this assumption, and given that the BFS algorithm in the subsequent step will consider only one connection per pair of users to deduce the most relevant diffusion pattern, we analyze the distribution of average retweet times. We then apply and evaluate two types of filters: one retaining links within the top 20% of this distribution and another within the top 75%. Links outside these thresholds are removed. The BFS algorithm is then applied to this refined network structure.

Appendix B. Identifying influencers

We identify influencers using the CI method [22], which involves an algorithm designed to find a minimal set of nodes capable of triggering a global cascade within the network,

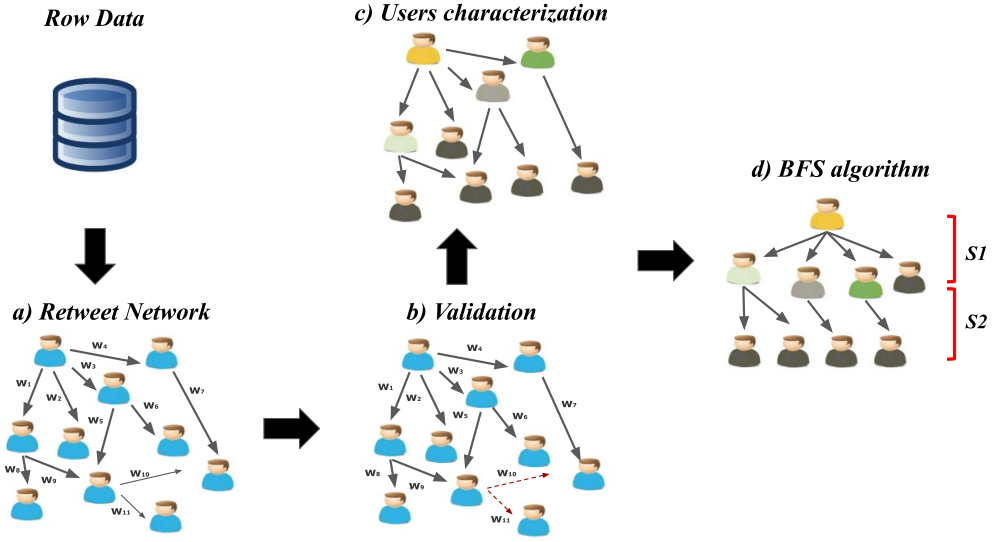


Figure A1. Pipeline. Illustrative representation of the methodology we follow in this work. (a) Starting from the raw data, we build a retweet network by considering all the retweets containing a URL redirecting to one of the news outlets in table A1. To each connection, we associate weights (w) and an average Δt as explained in section 4. (b) The retweet network undergoes link validation, and (c) after this step, users are classified as sources (sandy yellow), opinion leaders (green), opinion-leading influencers (light green), influencers (light gray), or adopters (dark gray). (d) On this network, we perform a BFS algorithm to identify the backbone (or skeleton) of the information diffusion.

following the linear threshold model [20]. For each node i , the CI is defined as follows:

$$CI_{\ell}(i) = (k_i - 1) \sum_{j \in \partial \text{Ball}(i, \ell)} (k_j - 1), \quad (\text{B.1})$$

where $\text{Ball}(i, \ell)$ is the set of nodes inside a ball of radius ℓ around node i , with the radius defined as the shortest path distance, and $\partial \text{Ball}(i, \ell)$ is the frontier (surface) of the ball. Here, k_i is the degree of node i . The value obtained for each node effectively evaluates the node's influence, considering the connectivity of nodes in its neighborhood. For our case, we choose $\ell = 1$. Moreover, since this task is nondeterministic polynomial-time (NP) complete, the algorithm is impractically slow. Therefore, we apply a computationally efficient CI heuristic that provides an approximate solution.

After computing the CI for each node in the network, whose distribution for the general case network is represented in figure B1, we identify the most influential nodes. We select a number that captures, on average, the top 0.1% of the influencers for the three categories studied: general case, left, and right; this resulted in an arbitrary threshold of 1000 elements as influencers.

Table A1. Hostnames in each media category. The tables contain information about the pages related to the news outlets considered in this study.

Left-leaning news		Right-leaning news	
Hostnames	Username	Hostnames	Username
1	nytimes.com	nytimes	nypost.com
2	washingtonpost.com	washingtonpost	WSJ
3	cnn.com	CNN	forbes.com
4	politico.com	politico	WashTimes
5	nbcnews.com	NBCNews	FoxBusiness
6	theguardian.com	guardian	BulwarkOnline
7	theatlantic.com	TheAtlantic	MarketWatch
8	abcnews.go.com	ABC	RealClearNews
9	npr.org	NPR	detroitnews
10	bloomberg.com	business	dallasnews
11	cbsnews.com	CBSNews	Rasmussen_Poll
12	cnbc.com	CNBC	chicagotribune
13	axios.com	axios	Jerusalem_Post
14	msn.com	MSN	
15	news.yahoo.com	YahooNews	
16	independent.co.uk	Independent	
17	latimes.com	latimes	
18	citizensforethics.org	CREWcrew	
19	buzzfeednews.com	BuzzFeed	

Right news		Left news	
Hostnames	Username	Hostnames	Username
1	foxnews.com	FoxNews	rawstory.com
2	dailycaller.com	DailyCaller	RawStory
3	washingtonexaminer.com	dcexaminer	MSNBC
4	justthenews.com	jsolomonReports	thedailybeast.com
5	thefederalist.com	FDRLST	thedailybeast
6	dailywire.com	realDailyWire	HuffPost
7	theepochtimes.com	EpochTimes	politicosusa
8	nationalreview.com	NRO	palmerreport.com
9	saraacarter.com	SaraCarterDC	PalmerReport
10	townhall.com	townhallcom	MotherJones
11	theblaze.com	theblaze	motherjones.com
12	thepostmillennial.com	TPostMillennial	vox.com
13	westernjournal.com	WestJournalism	voxdotcom
14	redstate.com	RedState	vanityfair.com
15	thegreggjarrett.com	GreggJarrett	VanityFair
16	bizpacreview.com	BIZPACReview	nymag.com
17	twitchy.com	TwitchyTeam	NYMag
18	trendingpolitics.com	CKeirns	NewYorker
19	lifenews.com	LifeNewsHQ	dailynos.com
			dailykos
			Slate
			Salon
			RollingStone
			thenation
			AlterNet
			rdevro

Table A2. Example of influencers (I), opinion leaders (OL), and opinion-leading influencers (OLI) from the retweet network obtained by considering all the sources.

Influencers			
	Number of followers	Name	Username
1	664 345	The Lincoln Project	ProjectLincoln
2	879 979	Laurence Tribe	tribelaw
3	2248 593	Mark R. Levin	marklevinshow
4	2428 507	James Woods	RealJamesWoods
5	607 927	Tea Pain	TeaPainUSA
6	5164 374	Donald Trump Jr	DonaldJTrumpJr
7	3685 092	Eric Trump	EricTrump
8	1074 520	60 Minutes	60Minutes
9	31 896	Don Moynihan	donmoyn
10	109 779	Ryan Goodman	rgoodlaw
Opinion-leading influencers			
	Number of followers	Name	Username
1	81 994 886	Donald J. Trump	realDonaldTrump
2	638 314	Natasha Bertrand	NatashaBertrand
3	1357 350	Maggie Haberman	maggieNYT
4	6142 647	Joe Biden	JoeBiden
5	687 572	Bill Kristol	BillKristol
6	2564 714	Jake Tapper	jaketapper
7	272 043	Greg Sargent	ThePlumLineGS
8	804 111	Daniel Dale	ddale8
9	233 827	Jeffrey Goldberg	JeffreyGoldberg
10	432 294	Peter Baker	peterbakernyt
Opinion leaders			
	Number of followers	Name	Username
1	36 665	Rick Tyler-Still Right	rickwtyler
2	17 428	Jonathan Landay	JonathanLanday
3	34 233	jimrutenberg	jimrutenberg
4	19 490	Henry J. Gomez	HenryJGomez
5	8339	Christian Datoc	TocRadio
6	9655	Michael Schwirtz	mschwirtz
7	67 721	Charlie Savage	charlie_savage
8	132 581	Senator Ron Johnson	SenRonJohnson
9	446 379	Sherrod Brown	SenSherrodBrown
10	135 322	Leana Wen, M.D.	DrLeanaWen

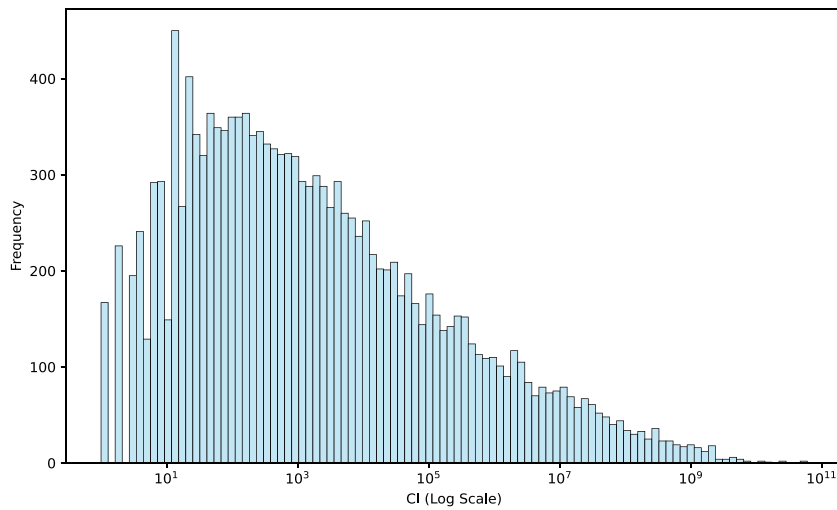


Figure B1. Distribution of collective influence in the retweet network. News related to both left- and right-leaning sources.

Appendix C. Further analysis

To support the results obtained using the BFS algorithm and to assess the potential loss of connections, and by extension, information, we measured the number of connections from adopters to other types of actors within the validated networks. The idea is to assess how many connections go against the direction individuated by the BFS and, therefore, observe the impact of our approximation. The results are shown in table C1.

Table C1. Links from adopters (ADP) to OL, OLI and I.

Link type	Left (%)	Right (%)	Full (%)
<i>from ADP to OL</i>	0.0926	0.0555	0.0559
<i>from ADP to I</i>	0.1531	0.4571	0.0565
<i>from ADP to OLI</i>	0.0949	0.0310	0.0397

The low values observed in table C1 support our main results and corroborate the assumption that connections in other directions are less relevant to our primary hypothesis. These connections can be viewed in two ways. A part of those can be interpreted in relation to the results obtained in the BFS algorithm as the connections that goes from adopters to the other categories in the step 2 of the multi-step model. Another part of these connections can be viewed as noise on the main structure, which is obtained by averaging over many different events. The observation of such connections indicates that, although rare, there are instances where information flows in the opposite direction from that identified by our approach. However, these events are infrequent enough to be disregarded for the purposes of our analysis. Along with the robustness checks that

incorporate temporal information, these results confirm the reliability of the inferred structure.

References

- [1] Bennett W and Manheim J 2006 The one-step flow of communication *Ann. Am. Acad. Pol. Soc. Sci.* **608** 213–32
- [2] Bennett W L and Iyengar S 2008 A new era of minimal effects? The changing foundations of political communication *J. Commun.* **58** 707–31
- [3] Berelson B, Lazarsfeld P and McPhee W 1954 *Voting: a Study of Opinion Formation in a Presidential Campaign* (University of Chicago Press)
- [4] Bovet A and Makse H A 2019 Influence of fake news in Twitter during the 2016 US presidential election *Nat. Commun.* **10** 7
- [5] Bovet A, Morone F and Makse H A 2018 Validation of Twitter opinion trends with national polling aggregates: Hillary Clinton vs Donald Trump *Sci. Rep.* **8** 8673
- [6] Bruno M, Saracco F, Garlaschelli D, Tessone C and Caldarelli G 2020 The ambiguity of nestedness under soft and hard constraints *Sci. Rep.* **10** 11
- [7] Burt R 1999 The social capital of opinion leaders *Ann. Am. Acad. Pol. Soc. Sci.* **566** 37–54
- [8] Caruso T, Clemente G, Rillig M and Garlaschelli D 2022 Fluctuating ecological networks: a synthesis of maximum–entropy approaches for pattern detection and process inference *Methods Ecol. Evol.* **13** 09
- [9] Casalo L V, Flavián C and Ibáñez-Sánchez S 2020 Influencers on Instagram: antecedents and consequences of opinion leadership *J. Bus. Res.* **117** 510–9
- [10] Choi S 2015 The two-step flow of communication in Twitter-based public forums *Soc. Sci. Comput. Rev.* **33** 696–711
- [11] Clemente G V, Tessone C J and Garlaschelli D. 2023 Temporal networks with node-specific memory: unbiased inference of transition probabilities, relaxation times and structural breaks
- [12] Dubois E and Gaffney D 2014 The multiple facets of influence: identifying political influentials and opinion leaders on Twitter *Am. Behav. Sci.* **58** 1260–77
- [13] Flamino J, Galeazzi A, Feldman S, Macy M W, Cross B, Zhou Z, Serafino M, Bovet A, Makse H A and Szymanski B K 2021 Shifting polarization and Twitter news influencers between two U.S (arXiv:2111.02505)
- [14] Flamino J, Galeazzi A, Feldman S, Macy M W, Cross B, Zhou Z, Serafino M, Bovet A, Makse H A and Szymanski B K 2023 Political polarization of news media and influencers on Twitter in the 2016 and 2020 US presidential elections *Nat. Hum. Behav.* **7** 1–13
- [15] Gotelli N 2000 Null model analysis of species co-occurrence patterns *Ecology* **81** 2606–21
- [16] Jaynes E T 1957 Information theory and statistical mechanics *Phys. Rev.* **106** 620–30
- [17] Karlens R 2015 Followers are opinion leaders: the role of people in the flow of political communication on and beyond social networking sites *Eur. J. Commun.* **30** 301–18
- [18] Katz E 1957 The two-step flow of communication: an up-to-date report on an hypothesis *Public Opin. Q.* **21** 61–78
- [19] Katz E and Lazarsfeld P 1955 *Personal Influence: the Part Played by People in the Flow of Mass Communications* (Free Press)
- [20] Kempe D, Kleinberg J and Tardos E 2003 Maximizing the spread of influence through a social network *Proc. 9th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining* pp 137–46
- [21] Lazarsfeld P, Berelson B and Gaudet H 1944 *The People's Choice: How the Voter Makes Up His Mind in a Presidential Campaign* (Columbia University Press)
- [22] Morone F and Makse H A 2015 Influence maximization in complex networks through optimal percolation *Nature* **524** 65–68
- [23] Morone F, Min B, Bo L, Mari R and Makse H A 2016 Collective influence algorithm to find influencers via optimal percolation in massively large social media *Sci. Rep.* **6** 30062
- [24] Parisi F, Squartini T and Garlaschelli D 2020 A faster horse on a safer trail: generalized inference for the efficient reconstruction of weighted networks *New J. Phys.* **22** 053053
- [25] Pei S, Muchnik L, Andrade J S, Zheng Z and Makse H A 2014 Searching for superspreaders of information in real-world social media *Sci. Rep.* **4** 5547
- [26] Rogers E and Cartano D 1962 Methods of measuring opinion leadership *Public Opin. Q.* **26** 435–41
- [27] Saracco F, Di Clemente R, Gabrielli A and Squartini T 2015 Randomizing bipartite networks: the case of the world trade web *Sci. Rep.* **5** 10595

- [28] Serafino M, Monteiro H S, Luo S, Reis S D, Igual C, Lima Neto A S, Travizano M, Andrade J S J and Makse H A 2022 Digital contact tracing and network theory to stop the spread of COVID-19 using big-data on human mobility geolocalization *PLoS Comput. Biol.* **18** e1009865
- [29] Serrano M, Boguñá M and Vespignani A 2009 Extracting the multiscale backbone of complex weighted networks *Proc. Natl Acad. Sci. USA* **106** 6483–8
- [30] Squartini T and Garlaschelli D 2017 *Maximum-Entropy Networks: Pattern Detection, Network Reconstruction and Graph Combinatorics* (Springer)
- [31] Turcotte J, York C, Irving J, Scholl R and Pingree R 2015 News recommendations from social media opinion leaders: effects on media trust and information seeking *J. Comput. Mediat. Commun.* **20** 520–35
- [32] Watts D and Dodds P 2007 Influentials, networks and public opinion formation *J. Consum. Res.* **34** 441–58
- [33] Weeks B, Ardévol-Abreu A and Zúñiga H 2017 Online influence? Social media use, opinion leadership and political persuasion *Int. J. Public Opin. Res.* **29** 214–39
- [34] Weimann G 1982 On the importance of marginality: one more step into the two-step flow of communication *Am. Sociol. Rev.* **47** 764–73
- [35] Winter S, Neubaum G, Stieglitz S and Ross B 2021 Opinionleaders: a comparison of self-reported and observable influence of Twitter users *Inf. Commun. Who says what to whom on TwitterSoc.* **24** 1533–50
- [36] Wu S, Hofman J, Mason W and Watts D 2011 Who says what to whom on Twitter *Proc. 20th Int. Conf. On World Wide Web* pp 705–14
- [37] Zhou Z, Serafino M, Cohan L, Caldarelli G and Makse H A 2021 Why polls fail to predict elections *J. Big Data* **8** 1–28