# Information diffusion assumptions can distort our understanding of social network dynamics

Matthew R. DeVerna[a,b], Francesco Pierri[c], Rachith Aiyappa[a,b], Diogo Pacheco[a,d], John Bryden[a], and Filippo Menczer[a,b]

[a]Observatory on Social Media, Indiana University
[b]Luddy School of Informatics, Computing, and Engineering, Indiana University
[c]Department of Electronics, Information and Bioengineering, Politecnico di Milano
[d]Department of Computer Science, University of Exeter

This document has not been peer reviewed.
November 7, 2024

## Abstract

To analyze the flow of information online, experts often rely on platform-provided data from social media companies, which typically attribute all resharing actions to an original poster. This obscures the true dynamics of how information spreads online, as users can be exposed to content in various ways. While most researchers analyze data as it is provided by the platform and overlook this issue, some attempt to infer the structure of these information cascades. However, the absence of ground truth about actual diffusion cascades makes verifying the efficacy of these efforts impossible. This study investigates the implications of the common practice of ignoring reconstruction all together. Two case studies involving data from Twitter and Bluesky reveal that reconstructing cascades significantly alters the identification of influential users, therefore affecting downstream analyses in general. We also propose a novel reconstruction approach that allows us to evaluate the effects of different assumptions made during the cascade inference procedure. Analysis of the diffusion of over 40,000 true and false news stories on Twitter reveals that the assumptions made during the reconstruction procedure drastically distort both microscopic and macroscopic properties of cascade networks. This work highlights the challenges of studying information spreading processes on complex networks and has significant implications for the broader study of digital platforms.

# Introduction

The digital age has woven technology into nearly every aspect of daily life, generating an unprecedented volume of data on human behavior and societal trends (1, 2, 3). This abundance of data has catalyzed the rise of computational social science, a field that leverages computational techniques to analyze and interpret digital trace data, providing novel insights into human behavior and society at large (4, 5, 6, 7, 8). Over the past two decades, social media platforms have become a primary source of data, typically accessed through Application Programming Interfaces (APIs) available to the public (9). This access has been vital not only for academic research but also for government and industry sectors, highlighting the pivotal role of digital data in modern research landscapes (10, 11).

The importance of this data is reflected in the rapid growth of the scientific literature on the diffusion of online information, which has surged dramatically over the past decade. A bibliographic analysis (see Supplementary Information) shows that in the last six years alone, over a thousand articles from the Social Sciences, Physics, Engineering, Medicine, and other fields have been published annually on the topic. This vast literature has influenced our understanding of critical societal challenges, such as public health (12, 13, 14, 15), political communication (16, 17, 18, 19), disaster response (20, 21, 22), collective action (23, 24, 25), and human attention (26, 27), demonstrating the interdisciplinary appeal and broad impact of studying information diffusion in the digital age.

Despite significant advances in the field, data from social media platforms present important limitations (5, 26). For example, dynamic socio-technical systems continuously change, influencing user behavior in non-transparent ways even as we attempt to study them (28). The influence of evolving and opaque platform algorithms adds further complexity to analyses (29). Here we focus on an often-overlooked issue arising from the difficulty to reconstruct diffusion cascades. Some platforms, like Facebook and Instagram, only provide aggregated data about cascade sizes. Others, like Twitter (now X), Mastodon, Bluesky, and Threads, provide more data but typically obscure the pathways of information diffusion by attributing all resharing actions to the original poster. This misrepresentation conceals the true dynamics of how information spreads, hindering our ability to fully understand these processes (See Figure 1 (**b, c**)). To address this problem, some researchers have proposed methods to infer the structure of information cascades to predict the actual patterns of online information flow (30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41). However, this problem is notoriously challenging because the reconstruction process relies on numerous assumptions, introducing significant potential for error. Additionally, this potential for error increases as the cascade grows in size, because as the number of nodes increases, it becomes increasingly difficult to accurately identify the parent of a new post (42). In rare cases, researchers may collaborate with platforms, which have access to detailed information about sharing actions (43, 33, 34). However, the general lack of ground-truth data renders validation of reconstruction techniques impossible (11).

Because of this challenge, most studies analyze platform-provided data without reconstructing information diffusion cascades. The few studies that attempt to reconstruct the cascades typically rely on follower networks, which capture social connections between platform users and are leveraged to estimate potential content exposure (30, 44, 38, 40). However, these approaches have not been rigorously validated. Furthermore, gathering this data is typically infeasible for researchers working with large datasets due to platform limitations.

These limitations raise important concerns about the consequences of overlooking the complex dynamics of information diffusion on social media platforms. To what extent could these choices distort our understanding of online information dynamics, such as the role of influential users, the emergence of polarization, and the vulnerability of platforms and users to various forms of manipulation? In this study, we illustrate these concerns by first quantifying how bypassing the cascade reconstruction process altogether impacts measures of social influence in two case studies on Twitter and Bluesky. After determining that this omission dramatically affects assessments of node influence, we investigate the structural effects of different reconstruction approaches. Leveraging a widely studied dataset of over 100,000 Twitter news cascades (40, 45, 46, 47, 48, 49, 50), we uncover substantial discrepancies in cascades at both the micro and macro level.
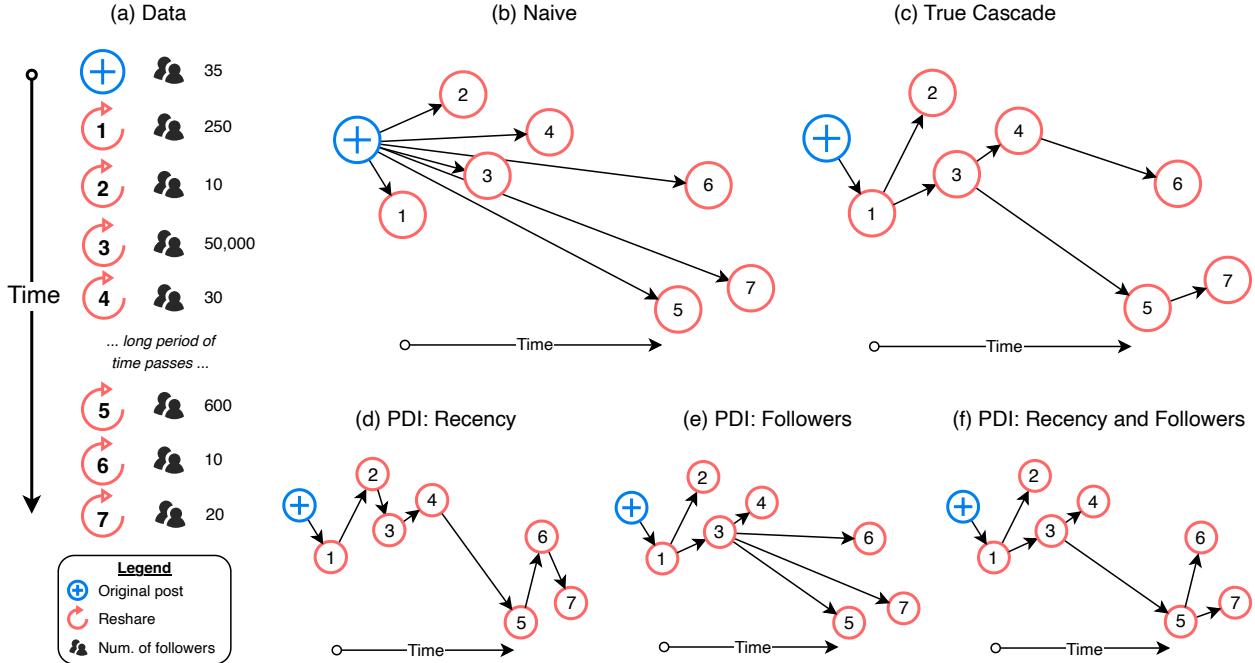
Figure 1: **Cascade reconstruction with Probabilistic Diffusion Inference.** **(a)**: Hypothetical empirical data of a message cascade with an original post (blue cross) and a sequence of resharing actions (red circles) over time. Each post is associated with a timestamp (represented by the time sequence) and the number of followers of the resharing user (next to the user icon). **(b)**: The naive cascade constructed from platform-provided data, which assumes that every user directly reshared the original post. **(c)**: The true cascade, reflecting the actual parent-child relationships. Panels **(d, e, f)** demonstrate different cascade reconstructions when applying various PDI assumptions. The recency assumption **(d)** prioritizes users who reshared the content more recently, capturing temporal dynamics. The followers assumption **(e)** gives higher resharing likelihood to users with more followers, emphasizing popularity. Incorporating both assumptions **(f)** captures both temporal activity and popularity into the cascade reconstruction.

# Results

## Cascade reconstruction

To understand how reconstructing information cascades impacts various analyses, we first introduce a general, parametric method that infers information (or message) cascades on microblogging platforms by leveraging empirical data about resharing activities. A *message cascade* is a tree structure where the root is the original poster of the message and a parent node's children are the users who reshared the message because they saw the parent's post. For each node in the cascade, the method infers the *parent* node, i.e., the prior node within the tree (user who previously posted or reposted the same message) that led to the resharing action. Linking all the posters of a message through these parent-child connections forms the message cascade.

Our method, called Probabilistic Diffusion Inference (PDI), relies on assumed probability distributions to weigh the likelihood of potential parents being the true parent within an information cascade. While this approach can flexibly incorporate any researcher-formulated probability distribution to capture the latest knowledge or potential platform changes, we adopt two assumptions based on previous work (40) about which users are more likely to be the parent of a resharer: users with more followers (*followers* assumption) and users who are more recently active in the cascade (*recency* assumption). These assumptions are visually represented in Figure 1 (**d, e, f**).

To model these assumptions, we calculate two probabilities for each potential parent node: one based

on their number of followers and the other taking into account the recency of their activity. A parameter $\alpha$ controls how much emphasis is placed on recency, with higher values giving more importance to recent posts. The relative influence of these two factors is adjusted using a parameter $\gamma$—higher values give more weight to follower counts, while lower values prioritize reshare recency. Further details on PDI and these assumptions can be found in the Methods section.

A set of cascade trees reconstructed from the data with the PDI method can be combined into a weighted *resharing network*. Nodes in this network represent users and edges capture the flow of information. Specifically, a link $(i \rightarrow j, w)$ represents a directed edge from user $i$ to user $j$, weighted by $w$, the number of times user $j$ reshared user $i$'s content. However, unlike reconstruction methods that generate a single cascade in deterministic fashion (40, 38), PDI can stochastically generate many different realizations of each cascade. This allows us to construct many versions of the weighted resharing network.

## Social influence measurement

Pinpointing the most influential individuals within social networks is a critical and widely studied challenge across fields ranging from epidemiology (51, 52) and public health (53) to political communication (54, 55, 19) and marketing (56, 57, 30). These key nodes can determine whether an epidemic will spread or whether a messaging campaign will achieve its intended impact.

To understand the effect of reconstructing information cascades on social influence analyses, we conduct case studies using data from two microblogging platforms: Twitter and Bluesky (see Methods for details). For each platform, we construct two types of resharing networks. The first, referred to as a *naive* network, is constructed directly from API-provided platform data connecting all resharing nodes to the original poster and disregarding any intermediate users in the cascade. The second, referred to as a *reconstructed* network, is generated after applying the PDI method as described above. Specifically, we generate 900 resharing networks—100 for each of the nine parameter settings obtained by combining $\gamma \in 0.25, 0.5, 0.75$ and $\alpha \in 1.1, 2.0, 3.0$. Note that the connection of the first reshare node is deterministic, as there is only one possible parent (the root). Therefore, for cascades with only two nodes (the original post and one reshare), no inference is needed. These cascades are included in all resharing networks.
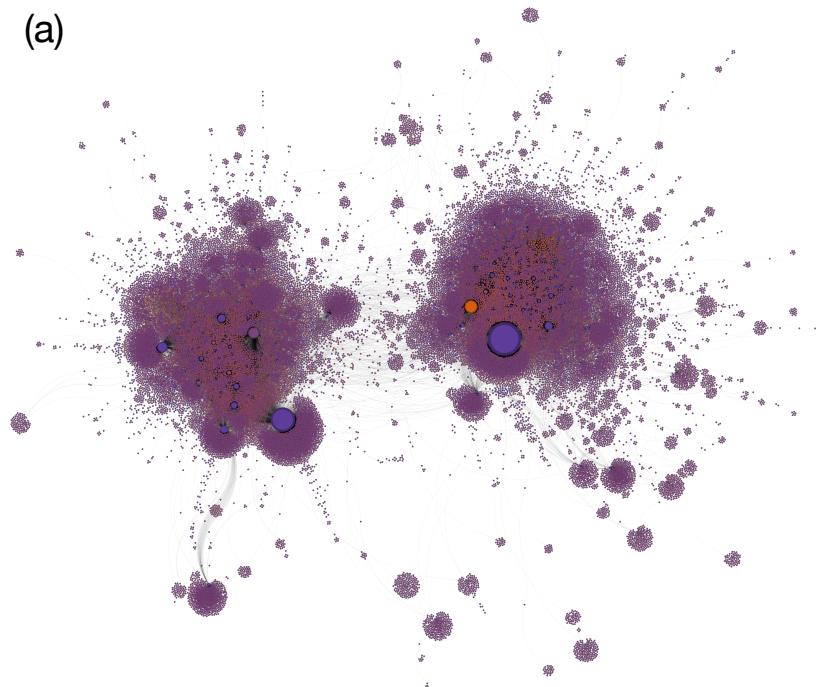
Comparing the two networks lets us determine the effects of the reconstruction method on the analyses of node influence—if the results are very similar, it would indicate that the reconstruction process has minimal impact on these analyses. To measure node influence, we calculate node out-strength, which we refer to as node strength for brevity. This is a widely recognized and intuitive metric, defined as the total number of reshares a node accumulates (58, 59, 60). As shown in Figure 2, there are important differences in node influence before and after reconstruction. In the naive resharing network, influence is concentrated among a few accounts that tend to reshare infrequently. In the reconstructed network, on the other hand, influence is more broadly distributed across many accounts, including amplifiers that tend to reshare content posted by others.

For a more quantitative analysis, our extensive set of reconstructed networks allows us to evaluate both the average impact of the reconstruction process and the robustness of our findings across different parameter settings. We begin by calculating Spearman's rank correlation ($\rho$) between node strength in the naive and reconstructed networks to quantify the changes in relative influence after reconstruction. Here, $\rho = 1$ signifies that the reconstruction process does not affect relative influence, while lower $\rho$ values indicate that node influence is affected.

Figure 3 presents the average correlation values for all tested parameter settings, revealing notable changes in node influence due to the reconstruction process on both platforms. In the Bluesky data, $\rho$ values range from 0.45 to 0.61, indicating a moderate shift in influence. On Twitter, the $\rho$ values are even lower, between 0.19 and 0.33, pointing to a significant reordering of node influence. These low correlations highlight the considerable impact that cascade reconstruction has in altering the perceived influence of nodes on both platforms.

To gain a deeper understanding of how the reconstruction process alters influence, we examine network changes at a single parameter setting ($\gamma = 0.25$ and $\alpha = 3.0$), with results presented in Figure 4. We observe similar trends across all parameter settings. Panels **a, d** compare node strength between a single PDI-reconstructed network and its corresponding naive network, revealing how influence shifts within the network on both platforms. The inclusion of secondary nodes as potential parents rewires network
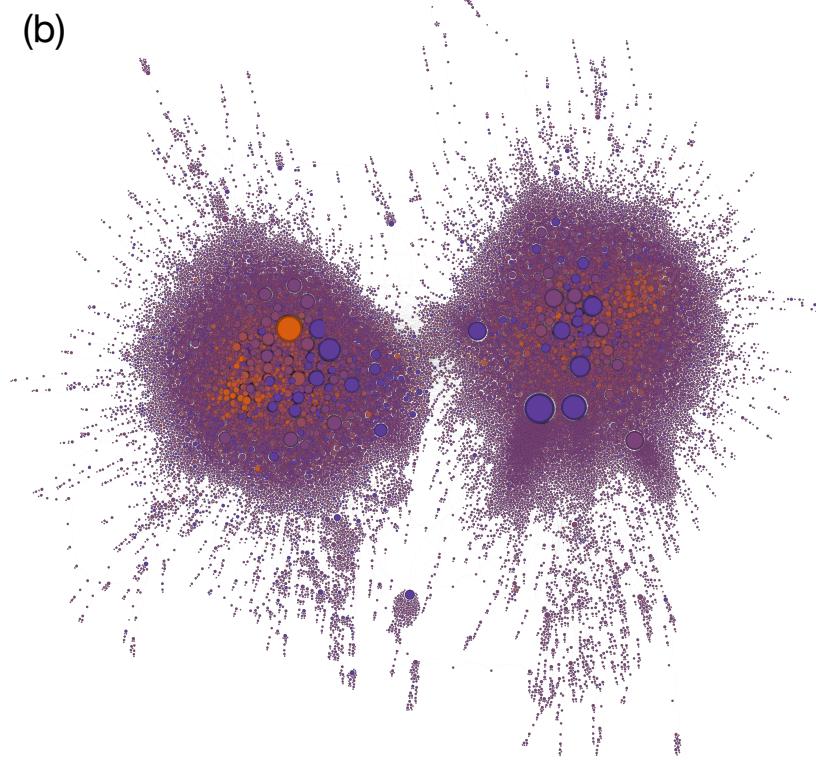
4

Figure 2: **Effects of cascade reconstruction on a Twitter resharing network.** (**a**) shows the naive network, while (**b**) displays a version of the same network reconstructed using PDI parameters $\gamma = 0.5$ and $\alpha = 2.0$. For illustration purposes, only nodes from the two largest communities are included. Node size reflects the number of retweets received by an account, with larger nodes representing more influential accounts. Node color represents the number of retweets an account has made, where red nodes indicate amplifiers that extensively retweet others' content.
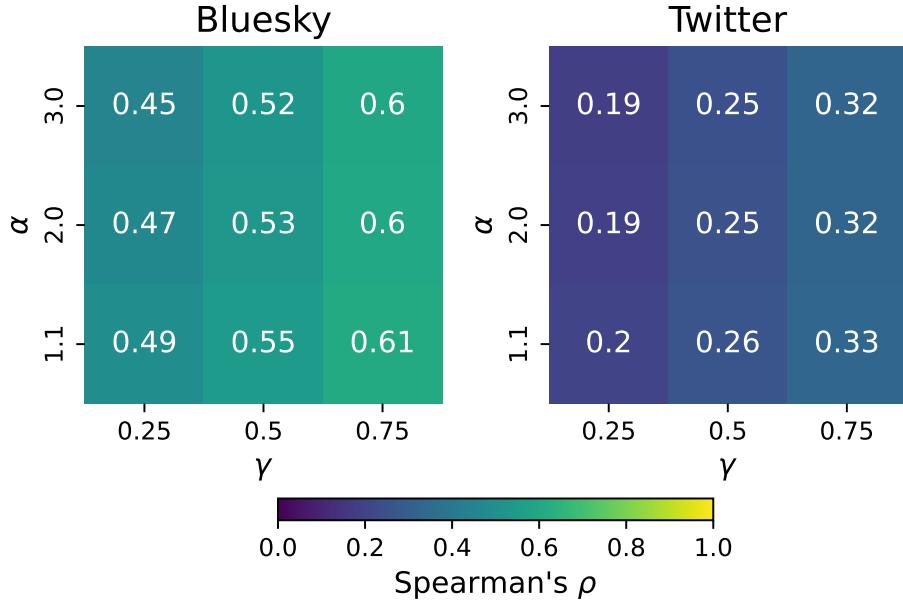
Figure 3: **Node influence is substantially affected by cascade reconstruction.** Heat map cells display the mean Spearman's correlation $\rho$ between node strength values in naive and PDI-reconstructed networks, averaged over 100 versions of the reconstructed network at the specified parameter settings. A $\rho$ value of one means the reconstruction doesn't alter node influence, while values closer to zero suggest significant changes. The maximum standard deviation of correlation values for any parameter setting is 0.001 for Twitter and 0.003 for Bluesky (see the Supplementary Information for full statistics).

connections, causing some to gain influence while others lose it.

Which nodes gain influence through the reconstruction process, and which ones see it diminish? Panels **b, e** show that, on both platforms, nodes with low strength in the naive network tend to experience a modest increase in strength, while nodes with high initial strength undergo a significant decrease. Specifically, 56% of Bluesky users and 91% of Twitter users exhibit a small average increase in influence after reconstruction, as the influence of secondary users is no longer ignored. Only nodes with an initial strength below two on Bluesky and three on Twitter display a median increase in average influence. For most nodes with a higher initial strength, the reconstruction process leads to a substantial decrease in average strength.

Finally, we examine the most influential nodes, defined by their total strength (number of reshares). For each of the 100 reconstructed networks, we compare the top 1%, 5%, and 10% of influential nodes to those in the naive network. We measure their similarity using the Jaccard index, which calculates the ratio of the size of the intersection to the size of the union of the two sets (61). A Jaccard index of one indicates that the reconstruction process does not affect the identification of the most influential nodes, while a value of zero indicates no overlap, signaling a substantial influence shift due to the reconstruction process. This analysis reflects a substantial restructuring of the network, leading to highly dissimilar sets of top influential nodes (panels **c, f**). At most, we observe a mean Jaccard similarity of only about 33% between the two network types on Bluesky (panel **c**; $k$=10%), while at the lowest, the overlap decreases to around 10% on Twitter (panel **f**; $k$=5%). This result suggests that analyses of superspreaders of information based on naive resharing networks might be misclassifying substantial portions of influential nodes.

## Information cascade structure

Let us now analyze how decisions made during the reconstruction process affect *individual* information cascades at the microscopic level. We posit that if distinct reconstruction methods generate cascades with different structural properties, they will have a substantial impact on downstream analyses. Given that no platform-provided data exists to validate *any* proposed method, such a finding would raise concerns about
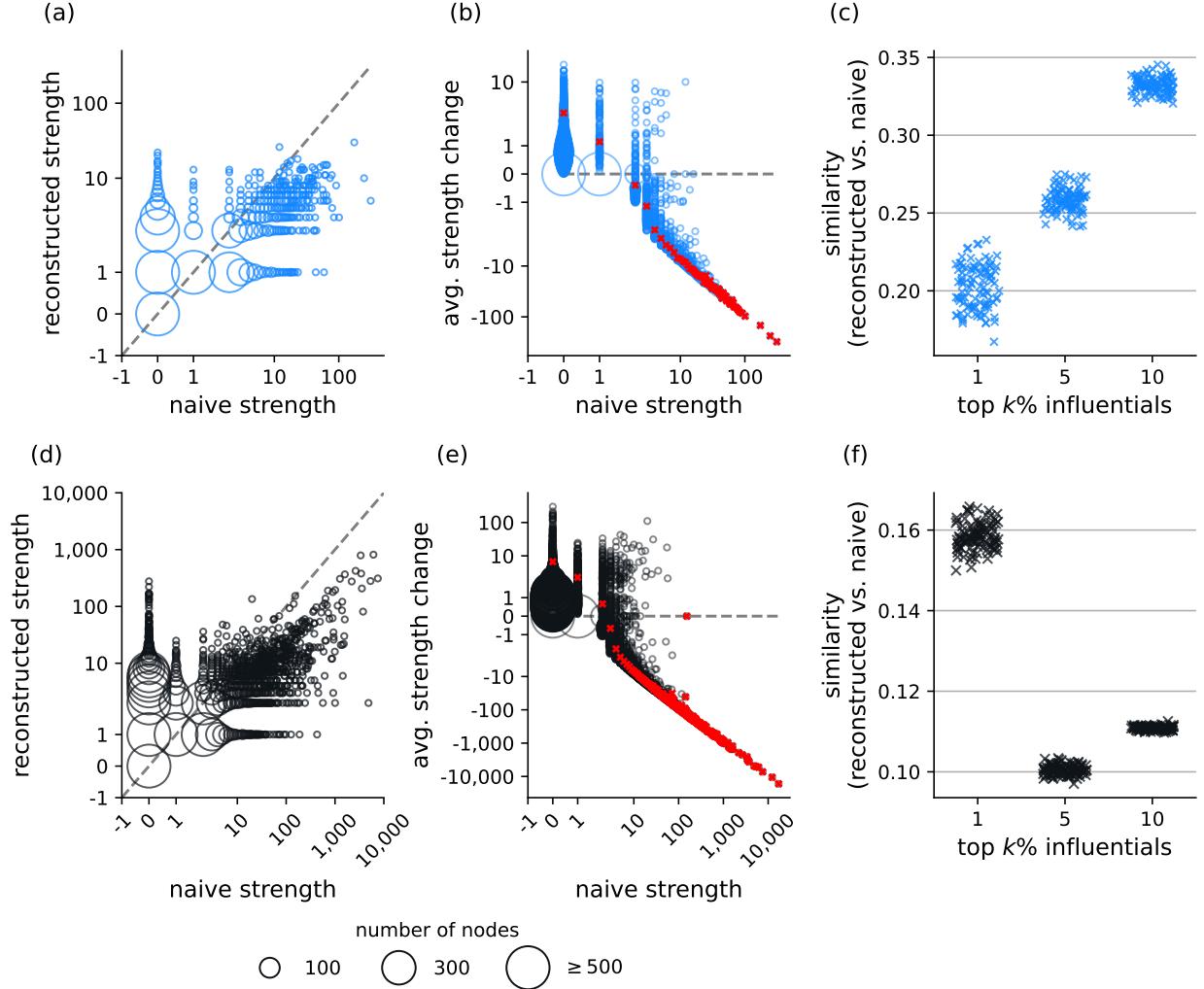
Figure 4: **Resharing networks reconstructed using the PDI method show substantial shifts in node influence compared to those built from naive data, on both Bluesky and Twitter.** Panels **(a, b, c)** present results for Bluesky, while panels **(d, e, f)** show results for Twitter. All panels reflect reconstructions using PDI parameters $\gamma = 0.25$ and $\alpha = 3.0$. **(a, d)**: Comparison of node strength between a single version of the PDI-reconstructed network and the corresponding naive network. **(b, e)**: Average change in node strength relative to naive strength, across all 100 PDI reconstructions. The red crosses show the median values. **(c, f)**: Jaccard similarity between the top $k$% of influential nodes identified based on node strength from reconstructed and naive networks. Each point represents one of the 100 possible comparisons. Circle sizes in panels **(a, b, d, e)** represent the number of nodes at each point. For visualization purposes, we use the same size for all points with 500 or more nodes.

the validity of social network studies that rely on network structure.

Based on this premise, we compare the PDI method with an alternative reconstruction approach employed in a prominent analysis of verified true and false rumor cascades on Twitter, spanning 2006 to 2017 (40). The latter method, known as Time-Inferred Diffusion (TID), infers a single version of each cascade using heuristics similar to PDI, relying on follower-network data and temporal dynamics to guess the parent of each reshare (38). A key assumption of TID is that the probability of an account resharing a post from someone they don't follow is zero. However, this assumption is problematic in the era of recommendation algorithms, where, for example, half of the content in a user's "For you" feed on X comes from accounts they do not follow (62). While Twitter did not make their algorithmic feed default for all users until 2016, such an assumption also overlooks exposure to content via organic search or off-network exposure. Here, we explore whether these different reconstruction methods alter the resulting cascades.

Our analysis reconstructs over 40,000 cascades from Vosoughi et al. (2018), originally generated by the TID method, using the PDI approach with the same parameter settings from our earlier analysis: $\gamma \in \{0.25, 0.5, 0.75\}$ and $\alpha \in \{1.1, 2.0, 3.0\}$. We focus on cascades with three or more nodes ($n = 28{,}062$), as no inference is required for cascades of size two (where the single resharing user has only one potential parent). For each setting, we generate 100 versions of each cascade using PDI and calculate the similarity between the different versions of the same cascade. We compare the PDI versions of a single cascade against each other ($\binom{100}{2} = 4{,}950$ comparisons) as well as against the TID version (100 comparisons). This allows us to study not only how the PDI and TID reconstruction approaches differ from each other, but also the variety of cascades generated by a specific reconstruction heuristic. We measure the similarity between two cascades using the Jaccard index of their edge sets. A similarity of one indicates that the two cascades are identical, whereas a similarity close to zero suggests significant differences.

Figure 5 shows that, on average, different reconstruction heuristics yield highly dissimilar cascades, regardless of PDI parameter settings. This discrepancy is especially pronounced for larger cascades (size $\gtrsim 100$), with similarity consistently below 0.2 and even dropping below 0.1 when $\gamma = 0.25$. A similar pattern emerges when comparing different PDI versions against each other.

The above results suggest that reconstruction decisions have a substantial impact on the inferred cascades. But how do these differences influence the overall topological structure? To address this question, let us shift our analysis to the macroscopic level. Using all reconstructed cascades from the same dataset, we compare the average distributions of several topological properties based on the 100 cascades produced using each of PDI setting as well as those generated with TID. We examine three key cascade properties: depth, maximum breadth, and structural virality. Depth is defined as the longest chain of unique reshares from the original post in the cascade, whereas maximum breadth captures the largest number of users at any single depth in the cascade. Structural virality (38) is defined as the average shortest-path length between all pairs of nodes in the cascade. It estimates the extent to which content spreads through a single, large broadcast (low structural virality) versus multiple levels, where each individual contributes only a small part to the overall spread (high structural virality).

Figure 6 presents the results of this analysis. For all metrics, we observe that different reconstruction approaches lead to significantly different network distributions, as confirmed by Kolmogorov-Smirnov two-sample tests. We have ten reconstruction heuristics—nine PDI settings plus TID—and three metrics, leading to $3 \times \binom{10}{2} = 135$ possible comparisons. 122 of these (90%) were found to be significantly different after applying Bonferroni's correction ($P < 0.05$; see the Supplementary Information for details). These changes follow expected patterns. For instance, as $\gamma$ decreases, giving more weight to the recency of a potential parent's post, both the depth and structural virality of cascades increase. Reducing $\gamma$ also lowers the maximum breadth, as the influence of individual prominent accounts with many followers diminishes, and longer chains within a cascade are drawn. These findings further emphasize how sensitive the inferred network structure is to the specific reconstruction method used.

## Discussion

This study demonstrates that the reconstruction of information cascades can fundamentally reshape network structures, significantly altering which nodes are identified as influential. In particular, naive network analyses that rely solely on platform-provided data can overestimate the influence of original posters

Figure 5: **Cascades reconstructed in different ways are highly dissimilar, especially for larger cascades.** Each panel shows the mean cascade similarity as a function of cascade size, with similarity measured using the Jaccard index. The panels correspond to different reconstruction parameter settings. Fit lines are generated using locally weighted robust smoothing of the ~28k mean values, while points represent the means in 500 equally-sized x-axis bins. Error bars show 95% confidence intervals calculated from 1,000 bootstraps.

Figure 6: **The structural properties of cascades are significantly altered by different reconstruction methods.** Panels **(a)**, **(b)**, and **(c)** show the complementary cumulative distribution functions (CCDF) for cascade depth, structural virality, and maximum breadth, respectively. Cascades are reconstructed with the TID (purple) and PDI (other lines) methods. 100 versions of each PDI cascade are generated for each parameter setting. Lines represent CCDFs based on the mean values across these versions.
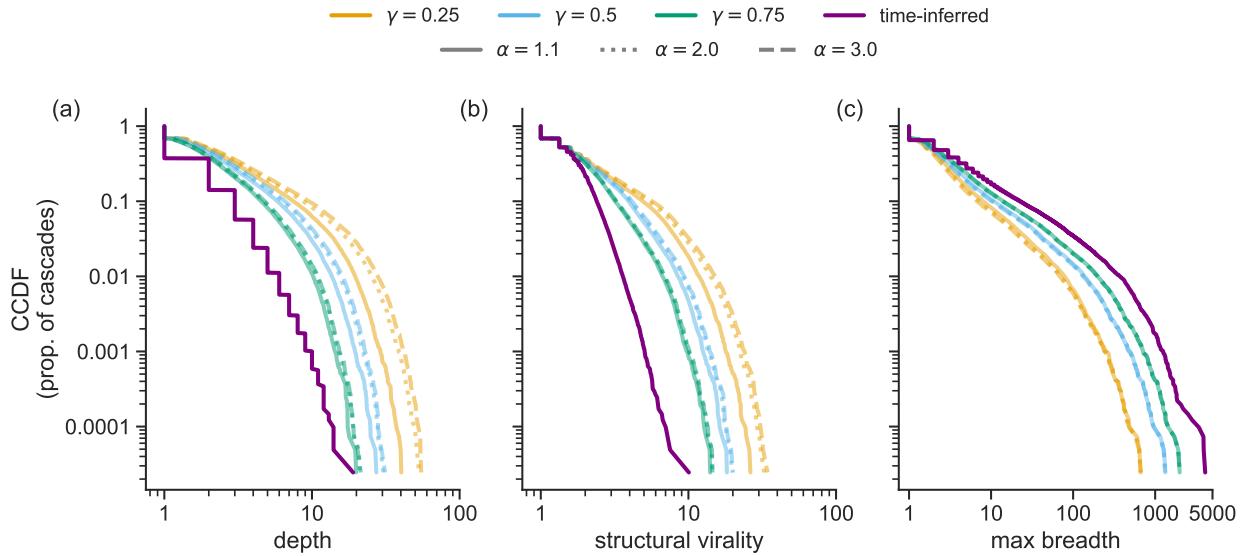
and underestimate the amplification role of intermediate resharers. Future work might examine how other structural features of nodes in the resharing network, like eigenvector and $k$-core centrality, are impacted by the reconstruction process. Furthermore, we observe that the assumptions embedded within different reconstruction methods significantly affect how we interpret the structure of individual cascades and collective resharing networks.

These findings were enabled by Probabilistic Diffusion Inference, a novel and flexible approach for reconstructing information cascades. By combining stochastically generated realizations of each cascade, we are able to construct many versions of weighted resharing (influence) networks. This allows us to explore the variance in outcomes of interest. An extension that we have not explored in this manuscript is probabilistically reflecting the influence of multiple parents on a node within a *single* cascade. By combining multiple reconstruction realizations, we can represent a cascade as a probability-weighted acyclic graph rather than a simple tree. This would make it possible to causally attribute a reshare action not only to one previous action, but to multiple prior exposures (33, 34).

The focus of PDI on the underlying assumptions also makes cascade reconstruction transparent. This helps researchers fine-tune assumptions and assess their impact, enabling a deeper exploration of the human and algorithmic factors driving information diffusion. For example, how might diffusion dynamics shift if a platform, like X, actively promotes certain political actors, as some have suggested (63, 64, 65)? How does this differ from Meta's new microblogging platform, Threads, which has indicated it will not insert unwanted political content into user feeds (66)? Researchers could incorporate node features, like political content, into the probability distributions to explore these and other interesting research questions.

Despite these benefits, PDI should not be considered "more accurate" than other techniques, such as the Time-Inferred Diffusion. The validation of inference methods requires the availability of ground-truth information diffusion data (67, 68, 69). If platforms were to publicly share cascade data with researchers, the PDI framework could be leveraged to refine assumptions and optimize parameter settings for more accurate modeling. However, we note that even platforms have to make assumptions about parent attribution, as users may be exposed to a piece of content in different ways prior to sharing, which are not revealed by their specific sharing action.

The substantial divergence we have found between networks reconstructed with different methods un-

10

derscores the potential risks for researchers who study online phenomena—especially when relying on naive networks provided by platforms. Given the widespread reliance on platform-provided data and the lack of ground truth for diffusion cascades, researchers must approach these analyses with caution. Such data is inherently complex and may be incomplete, exacerbating the challenge of accurately capturing underlying dynamics. Therefore, it is essential to develop methods that can effectively account for these limitations. These issues have far-reaching implications for fields that analyze social media networks, such as conservation science (70), political communication (54, 55, 19), public health (53), and epidemiology (71, 72). Future research should focus on identifying which analyses are most sensitive to reconstruction methods and ensuring their robustness across varying assumptions.

Computational social science must continue to develop innovative analytical approaches that make transparent assumptions and are robust to rigorous methodological scrutiny, transparency in assumptions (73, 74, 75). Such progress is crucial for deepening our understanding of complex digital ecosystems and the social dynamics that unfold within them.

# Methods

## Data

The Twitter dataset from our social influence analysis consists of 10,000 English-language retweet cascades sampled from the Indiana University 2022 U.S. Midterms Multi-Platform Social Media Dataset. This dataset captures online conversations about the 2022 midterm elections. It was gathered using a snowball sampling approach to collect keywords that were relevant to the 2022 U.S. Midterm elections. Please see the original publication for all details (76). The cascades analyzed in this study are a subset of the full corpus. We randomly selected cascades that originated between November 2, 2022, and November 8, 2022 (Election day), while including retweets up to November 15, 2022, to fully capture their diffusion (38). The resulting dataset contains over 187,443 tweets shared by 128,930 unique users.

Bluesky is a decentralized, Twitter-like micro-blogging platform designed to offer a federated social experience (77, 78). We collected all data from Bluesky between March 1, 2024, and March 14, 2024, using the public Firehose endpoint, which streams all posts shared on the platform (79). We then randomly sampled 5,000 repost cascades originating in the first seven days of this period, following the platform's public launch (80). The same sampling procedure used for the U.S. Midterm dataset was applied, capturing reposts up to one week later (March 21, 2024). We excluded 290 cascades from our analysis: 271 due to missing metadata for at least one user's follower count, and 19 because of timestamp discrepancies caused by Bluesky's distributed architecture (81). This resulted in a final dataset of 4,710 cascades consisting of 21,338 posts from 15,550 users.

We analyze topological network properties using a dataset of rumor cascades from Twitter (40), provided by the authors in a pre-processed and anonymized format for replication purposes. The original study gathered retweet cascades of both true and false content, verified by six independent fact-checking organizations. Specifically, the authors collected all English-language replies to tweets containing links to fact-checking articles. The initial dataset included approximately 126,000 English-language rumor cascades shared on Twitter by over 3 million users between 2006 and 2017. We excluded 84,221 cascades without any retweets. Additionally, since the PDI method requires follower counts, we also removed 1,242 cascades where this information was missing for at least one user in the cascade. This resulted in a final dataset of 40,839 cascades for analysis.

## Probabilistic Diffusion Inference

The PDI method estimates the likelihood that each user within a social media cascade is the original source of content for subsequent resharers. Consider a cascade $c$ involving a sequence of $N_c$ users, $U^c = \{u_0^c, u_1^c, \ldots, u_{N_c}^c\}$, where $u_0^c$ is the originator of the content, and each subsequent user $u_i^c$ represents the $i$-th person to reshare it. To determine the parent of $u_i^c$ —the source of $u_i^c$'s reshare— PDI considers the subset of all prior users $U_i^c \subset U^c = \{u_j^c \; \forall j < i\}$ as potential parents, each with a probability $p_{ij}$ of being selected as the parent of $u_i^c$. For all resharing users in the cascade, a potential parent is selected as the parent based on these probabilities.

PDI enables flexible computation of the probabilities $p_{ij}$ using researcher-defined assumptions. In this work, we adopt two common assumptions. First, users with more followers are more likely to be the parents of a resharing user ([82]), which we refer to as the *followers* assumption. Second, users who recently reshared the content are more likely to be the true parents of subsequent users ([83], [84]), referred to as the *recency* assumption.

The probability of a potential parent $u_j \in U_i^c$ according to the followers assumption is given by:

$$p_{ij}^{\mathcal{F}} = \frac{F(u_j)}{\sum_{u_k \in U_i^c} F(u_k)} \tag{1}$$

where $F(u)$ represents the mean number of followers of user $u$ during the observed period.

The recency assumption is modeled using a power-law distribution, which has been shown to describe the timing of resharing behavior on social media platforms ([83], [84]):

$$P(\Delta_{ij}^c) = \frac{\alpha - 1}{\Delta_{\min}} \left( \frac{\Delta_{ij}^c}{\Delta_{\min}} \right)^{-\alpha}, \tag{2}$$

where $\Delta_{ij}^c$ is the time (in seconds) between the post by potential parent $u_j^c$ and the reshare by user $u_i^c$, $\Delta_{\min}$ is a minimum time delay (one second), and $\alpha$ is a parameter that expresses the tendency for reshares to be clustered in time. Then, the probability of potential parent $u_j$ according to the recency assumption is calculated by:

$$p_{ij}^{\mathcal{T}} = \frac{P(\Delta_{ij}^c)}{\sum_{u_k \in U_i^c} P(\Delta_{ik}^c)}. \tag{3}$$

We consider the followers and recency assumptions as independent factors and combine them using a weighting parameter $\gamma$, yielding the overall probability that $u_j^c$ is the true parent of $u_i^c$:

$$p_{ij} = \gamma p_{ij}^{\mathcal{F}} + (1 - \gamma) p_{ij}^{\mathcal{T}}. \tag{4}$$

## Acknowledgments

## Author contributions

MRD and FP conceptualized the study design with guidance from FM. FP, MRD and RA collected and curated data. Code was primarily written by MRD with review and input from FP. MRD and FP conducted the analyses and created the visualizations, with input and guidance from FM. The Probabilistic Diffusion Inference method was developed by MRD with input from RA and JB. The first draft was written by MRD and FP, with input from all authors. FM oversaw the progression of the study. Funding was acquired by FM and FP.

## Competing interests

Authors declare that they have no competing interests.

## Data and code availability

Code and data are available in public repositories on GitHub (github.com/osome-iu/cascade_reconstruction) and Zenodo (doi.org/10.5281/zenodo.13994029).

## Ethical matters

This study, focusing on public data, poses minimal risk to human subjects. Consequently, the Indiana University Institutional Review Board has exempted it from review (protocol numbers 1102004860 and 23757).

## References

[1] Ilya Levin and Dan Mamlok. Culture and Society in the Digital Age. *Information*, 12(2):68, 2021. URL `https://doi.org/10.3390/info12020068`.

[2] Seref Sagiroglu and Duygu Sinanc. Big data: A review. In *2013 International Conference on Collaboration Technologies and Systems*, pages 42–47. IEEE, 2013. URL `https://doi.org/10.1109/CTS.2013.6567202`.

[3] Martin Hilbert and Priscila López. The World's Technological Capacity to Store, Communicate, and Compute Information. *Science*, 332(6025):60–65, 2011. URL `https://doi.org/10.1126/science.1200970`.

[4] Jake M. Hofman, Duncan J. Watts, Susan Athey, Filiz Garip, Thomas L. Griffiths, Jon Kleinberg, Helen Margetts, Sendhil Mullainathan, Matthew J. Salganik, Simine Vazire, Alessandro Vespignani, and Tal Yarkoni. Integrating explanation and prediction in computational social science. *Nature*, 595:181–188, July 2021. doi: 10.1038/s41586-021-03659-0. URL `https://doi.org/10.1038/s41586-021-03659-0`.

[5] David M. J. Lazer, Alex Pentland, Duncan J. Watts, Sinan Aral, Susan Athey, Noshir Contractor, Deen Freelon, Sandra Gonzalez-Bailon, Gary King, Helen Margetts, Alondra Nelson, Matthew J. Salganik, Markus Strohmaier, Alessandro Vespignani, and Claudia Wagner. Computational social science: Obstacles and opportunities. *Science*, 369(6507):1060–1062, 2020. URL `https://doi.org/10.1126/science.aaz8170`.

[6] Matthew J. Salganik. *Bit by Bit: Social Research in the Digital Age*. Princeton University Press, Princeton, NJ, 2018.

[7] David Lazer, Alex Pentland, Lada Adamic, Sinan Aral, Albert-László Barabási, Devon Brewer, Nicholas Christakis, Noshir Contractor, James Fowler, Myron Gutmann, Tony Jebara, Gary King, Michael Macy, Deb Roy, and Marshall Van Alstyne. Computational Social Science. *Science*, 323(5915):721–723, February 2009. doi: 10.1126/science.1167742. URL `https://doi.org/10.1126/science.1167742`.

[8] Duncan J. Watts. A twenty-first century science. *Nature*, 445:489, 2007. URL `https://doi.org/10.1038/445489a`.

[9] Pascal Jürgens and Andreas Jungherr. A tutorial for using twitter data in the social sciences: Data collection, preparation, and analysis. SSRN, 2016. URL `https://dx.doi.org/10.2139/ssrn.2710146`.

[10] Michelle N. Meyer, John Basl, David Choffnes, Christo Wilson, and David M. J. Lazer. Enhancing the ethics of user-sourced online data collection and sharing. *Nat Comput Sci*, 3:660–664, 2023. URL `https://doi.org/10.1038/s43588-023-00490-7`.

[11] Brittany I. Davidson, Darja Wischerath, Daniel Racek, Douglas A. Parry, Emily Godwin, Joanne Hinds, Dirk van der Linden, Jonathan F. Roscoe, Laura Ayravainen, and Alicia G. Cork. Platform-controlled social media APIs threaten open science. *Nat Hum Behav*, 7:2054–2057, 2023. URL https://doi.org/10.1038/s41562-023-01750-2.

[12] Shu-Feng Tsao, Helen Chen, Therese Tisseverasinghe, Yang Yang, Lianghua Li, and Zahid A. Butt. What social media told us in the time of COVID-19: a scoping review. *Lancet Digital Health*, 3(3): e175–e194, March 2021. URL https://doi.org/10.1016/S2589-7500(20)30315-0.

[13] Dean Schillinger, Deepti Chittamuru, and A. Susana Ramìrez. From "Infodemics" to Health Promotion: A Novel Framework for the Role of Social Media in Public Health. *American Journal of Public Health*, 2020. URL https://ajph.aphapublications.org/doi/abs/10.2105/AJPH.2020.305746.

[14] Taha A. Kass-Hout and Hend Alhinnawi. Social media in public health. *British Medical Bulletin*, 108 (1):5–24, December 2013. URL https://doi.org/10.1093/bmb/ldt028.

[15] Mark Dredze. How Social Media Will Change Public Health. *IEEE Intelligent Systems*, 27(4):81–84, August 2012. URL https://doi.org/10.1109/MIS.2012.76.

[16] Nathaniel Persily, Joshua A Tucker, and Joshua Aaron Tucker. *Social media and democracy: The state of the field, prospects for reform*. Cambridge University Press, 2020.

[17] Joshua A. Tucker, Andrew Guess, Pablo Barbera, Cristian Vaccari, Alexandra Siegel, Sergey Sanovich, Denis Stukal, and Brendan Nyhan. Social Media, Political Polarization, and Political Disinformation: A Review of the Scientific Literature. *SSRN*, 2018. URL https://doi.org/10.2139/ssrn.3144139.

[18] David M. J. Lazer, Matthew A. Baum, Yochai Benkler, Adam J. Berinsky, Kelly M. Greenhill, Filippo Menczer, Miriam J. Metzger, Brendan Nyhan, Gordon Pennycook, David Rothschild, Michael Schudson, Steven A. Sloman, Cass R. Sunstein, Emily A. Thorson, Duncan J. Watts, and Jonathan L. Zittrain. The science of fake news. *Science*, 359(6380):1094–1096, 2018. URL https://doi.org/10.1126/science.aao2998.

[19] Stefan Stieglitz and Linh Dang-Xuan. Social media and political communication: a social media analytics framework. *Soc Netw Anal Min*, 3(4):1277–1291, 2013. URL https://doi.org/10.1007/s13278-012-0079-3.

[20] Christian Reuter, Amanda Lee Hughes, and Marc-André Kaufhold. Social Media in Crisis Management: An Evaluation and Analysis of Crisis Informatics Research. *International Journal of Human–Computer Interaction*, 2018. URL https://doi.org/10.1080/10447318.2018.1427832.

[21] J. Brian Houston, Joshua Hawthorne, Mildred F. Perreault, Eun Hae Park, Marlo Goldstein Hode, Michael R. Halliwell, Sarah E. Turner McGowen, Rachel Davis, Shivani Vaid, Jonathan A. McElderry, and Stanford A. Griffith. Social media and disasters: a functional framework for social media use in disaster planning, response, and research. *Disasters*, 39(1):1–22, 2015. URL https://doi.org/10.1111/disa.12092.

[22] David E. Alexander. Social Media in Disaster Risk Reduction and Crisis Management. *Sci Eng Ethics*, 20(3):717–733, 2014. URL https://doi.org/10.1007/s11948-013-9502-z.

[23] Zachary C. Steinert-Threlkeld. Spontaneous Collective Action: Peripheral Mobilization During the Arab Spring. *American Political Science Review*, 111(2):379–403, 2017. URL https://doi.org/10.1017/S0003055416000769.

[24] Sandra Gonzaalez-Bailón and Ning Wang. Networked discontent: The anatomy of protest campaigns in social media. *Social Networks*, 44:95–104, 2016. URL https://doi.org/10.1016/j.socnet.2015.07.003.

[25] Alexandra Segerberg and W. Social Media and the Organization of Collective Action: Using Twitter to Explore the Ecologies of Two Climate Change Protests. *Communication Review*, 2011. URL https://doi.org/10.1080/10714421.2011.597250.

[26] David Lazer. Studying human attention on the Internet. *Proc Natl Acad Sci U.S.A*, 117(1):21–22, January 2020. URL https://doi.org/10.1073/pnas.1919348117.

[27] Philipp Lorenz-Spreen, Bjarke Mørch Mønsted, Philipp Hövel, and Sune Lehmann. Accelerating dynamics of collective attention. *Nature communications*, 10(1):1759, 2019. URL https://doi.org/10.1038/s41467-019-09311-w.

[28] David Lazer, Eszter Hargittai, Deen Freelon, Sandra Gonzalez-Bailon, Kevin Munger, Katherine Ognyanova, and Jason Radford. Meaningful measures of human society in the twenty-first century. *Nature*, 595:189–196, 2021. URL https://doi.org/10.1038/s41586-021-03660-7.

[29] Claudia Wagner, Markus Strohmaier, Alexandra Olteanu, Emre Kıcıman, Noshir Contractor, and Tina Eliassi-Rad. Measuring algorithmically infused societies. *Nature*, 595:197–204, 2021. URL https://doi.org/10.1038/s41586-021-03666-1.

[30] Eytan Bakshy, Jake M. Hofman, Winter A. Mason, and Duncan J. Watts. Everyone's an influencer: quantifying influence on twitter. In *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining (WSDM)*, page 65–74, 2011. doi: 10.1145/1935826.1935845. URL https://doi.org/10.1145/1935826.1935845.

[31] Manuel Gomez-Rodriguez, Jure Leskovec, and Andreas Krause. Inferring networks of diffusion and influence. *ACM Trans. Knowl. Discov. Data*, 5(4), 2012. doi: 10.1145/2086737.2086741. URL https://doi.org/10.1145/2086737.2086741.

[32] Peter Cogan, Matthew Andrews, Milan Bradonjic, W. Sean Kennedy, Alessandra Sala, and Gabriel Tucci. Reconstruction and analysis of Twitter conversation graphs. In *ACM Conferences*, pages 25–31. Association for Computing Machinery, New York, NY, USA, August 2012. URL https://doi.org/10.1145/2392622.2392626.

[33] P. Alex Dow, Lada Adamic, and Adrien Friggeri. The Anatomy of Large Facebook Cascades. In *Proceedings of the Seventh International AAAI Conference on Web and Social Media (ICWSM)*, pages 145–154, 2013. doi: 10.1609/icwsm.v7i1.14431. URL https://doi.org/10.1609/icwsm.v7i1.14431.

[34] Adrien Friggeri, Lada Adamic, Dean Eckles, and Justin Cheng. Rumor Cascades. In *Proceedings of the Eight International AAAI Conference on Web and Social Media (ICWSM)*, pages 101–110, 2014. doi: 10.1609/icwsm.v8i1.14559. URL https://doi.org/10.1609/icwsm.v8i1.14559.

[35] Io Taxidou and Peter M. Fischer. Online analysis of information diffusion in twitter. In *ACM Other conferences*, pages 1313–1318. Association for Computing Machinery, New York, NY, USA, April 2014. URL https://doi.org/10.1145/2567948.2580050.

[36] Soroush Vosoughi. *Automatic detection and verification of rumors on Twitter*. PhD thesis, Massachusetts Institute of Technology, 2015. URL https://dspace.mit.edu/handle/1721.1/98553.

[37] Tom De Nies, Io Taxidou, Anastasia Dimou, Ruben Verborgh, Peter M. Fischer, Erik Mannens, and Rik Van de Walle. Towards Multi-level Provenance Reconstruction of Information Diffusion on Social Media. In *ACM Conferences*, pages 1823–1826. Association for Computing Machinery, New York, NY, USA, October 2015. URL https://doi.org/10.1145/2806416.2806642.

[38] Sharad Goel, Ashton Anderson, Jake Hofman, and Duncan J Watts. The structural virality of online diffusion. *Management Science*, 62(1):180–196, 2016. URL https://doi.org/10.1287/mnsc.2015.2158.

[39] Soroush Vosoughi, Mostafa 'Neo' Mohsenvand, and Deb Roy. Rumor Gauge: Predicting the Veracity of Rumors on Twitter. *ACM Trans Knowl Discov Data*, 11(4):1–36, 2017. URL https://doi.org/10.1145/3070644.

[40] Soroush Vosoughi, Deb Roy, and Sinan Aral. The spread of true and false news online. *Science*, 359 (6380):1146–1151, 2018. URL https://doi.org/10.1126/science.aap9559.

[41] Matteo Cinelli, Stefano Cresci, Walter Quattrociocchi, Maurizio Tesconi, and Paola Zola. Coordinated inauthentic behavior and information spreading on Twitter. *Decision Support Systems*, 160:113819, 2022. URL https://doi.org/10.1016/j.dss.2022.113819.

[42] John Bollenbacher, Diogo Pacheco, Pik-Mai Hui, Yong-Yeol Ahn, Alessandro Flammini, and Filippo Menczer. On the challenges of predicting microscopic dynamics of online conversations. *Appl Network Sci*, 6(1):1–21, 2021. URL https://doi.org/10.1007/s41109-021-00357-8.

[43] Eytan Bakshy, Itamar Rosenn, Cameron Marlow, and Lada Adamic. The role of social networks in information diffusion. In *Proceedings of the 21st international conference on World Wide Web*, pages 519–528, 2012.

[44] Sharad Goel, Duncan J Watts, and Daniel G Goldstein. The structure of online diffusion networks. In *Proceedings of the 13th ACM conference on electronic commerce*, pages 623–638, 2012. URL https://doi.org/10.1145/2229012.2229058.

[45] Nicolas Pröllochs and Stefan Feuerriegel. Mechanisms of True and False Rumor Sharing in Social Media: Collective Intelligence or Herd Behavior? *Proc ACM Hum.-Comput Interact*, 7(CSCW2):1–38, 2023. URL https://doi.org/10.1145/3610078.

[46] Christof Naumzik and Stefan Feuerriegel. Detecting False Rumors from Retweet Dynamics on Social Media. In *ACM Conferences*, pages 2798–2809. Association for Computing Machinery, April 2022. URL https://doi.org/10.1145/3485447.3512000.

[47] Jonas L. Juul and Johan Ugander. Comparing information diffusion mechanisms by matching on cascade size. *Proc Natl Acad Sci U.S.A*, 118(46):e2100786118, 2021. URL https://doi.org/10.1073/pnas.2100786118.

[48] Nicolas Pröllochs, Dominik Bär, and Stefan Feuerriegel. Emotions explain differences in the diffusion of true vs. false social media rumors. *Sci Rep*, 11(22721):1–12, 2021. URL https://doi.org/10.1038/s41598-021-01813-2.

[49] Francesco Ducci, Mathias Kraus, and Stefan Feuerriegel. Cascade-LSTM: A Tree-Structured Neural Classifier for Detecting Misinformation Cascades. In *ACM Conferences*, pages 2666–2676. Association for Computing Machinery, August 2020. URL https://doi.org/10.1145/3394486.3403317.

[50] Nir Rosenfeld, Aron Szanto, and David C. Parkes. A Kernel of Truth: Determining Rumor Veracity on Twitter by Diffusion Pattern Alone. In *ACM Conferences*, pages 1018–1028. Association for Computing Machinery, April 2020. URL https://doi.org/10.1145/3366423.3380180.

[51] Sinan Aral and Paramveer S. Dhillon. Social influence maximization under empirical influence models. *Nat Hum Behav*, 2:375–382, 2018. URL https://doi.org/10.1038/s41562-018-0346-z.

[52] Chris T. Bauch and Alison P. Galvani. Social Factors in Epidemiology. *Science*, 342(6154):47–49, 2013. URL https://doi.org/10.1126/science.1244492.

[53] Damon Centola. Social Media and the Science of Health Behavior. *Circulation*, 2013. URL https://www.ahajournals.org/doi/full/10.1161/CIRCULATIONAHA.112.101816.

[54] Kate Starbird, Renée DiResta, and Matt DeButts. Influence and Improvisation: Participatory Disinformation during the 2020 US Election. *Social Media + Society*, 9(2):20563051231177943, 2023. URL https://doi.org/10.1177/20563051231177943.

[55] Alexandre Bovet and Hernán A Makse. Influence of fake news in Twitter during the 2016 US presidential election. *Nature communications*, 10(1):7, 2019. URL https://doi.org/10.1038/s41467-018-07761-2.

[56] David Kempe, Jon Kleinberg, and Éva Tardos. Maximizing the spread of influence through a social network. In *ACM Conferences*, pages 137–146. Association for Computing Machinery, New York, NY, USA, August 2003. doi: 10.1145/956750.956769. URL https://doi.org/10.1145/956750.956769.

[57] Wei Chen, Chi Wang, and Yajun Wang. Scalable influence maximization for prevalent viral marketing in large-scale social networks. In *ACM Conferences*, pages 1029–1038. Association for Computing Machinery, 2010. URL https://doi.org/10.1145/1835804.1835934.

[58] Sarah J. Jackson and Brooke Foucault Welles. #Ferguson is everywhere: initiators in emerging counterpublic networks. *Information, Communication & Society*, 2016. URL https://doi.org/10.1080/1369118X.2015.1106571.

[59] Linyuan Lü, Duanbing Chen, Xiao-Long Ren, Qian-Ming Zhang, Yi-Cheng Zhang, and Tao Zhou. Vital nodes identification in complex networks. *Phys Rep*, 650:1–63, 2016. URL https://doi.org/10.1016/j.physrep.2016.06.007.

[60] Meeyoung Cha, Hamed Haddadi, Fabricio Benevenuto, and Krishna Gummadi. Measuring User Influence in Twitter: The Million Follower Fallacy. *ICWSM*, 4(1):10–17, 2010. URL https://doi.org/10.1609/icwsm.v4i1.14033.

[61] Paul Jaccard. The distribution of the flora in the alpine zone. *New Phytol*, 11(2):37–50, 1912. URL https://doi.org/10.1111/j.1469-8137.1912.tb05611.x.

[62] Twitter Blog. Twitter's Recommendation Algorithm, March 2023. URL https://blog.x.com/engineering/en_us/topics/open-source/2023/twitter-recommendation-algorithm. [Accessed: 25. Aug. 2024].

[63] Barbara Ortutay and Matt O'Brien. How Elon Musk uses his X social media platform to amplify right-wing views, August 2024. URL https://www.pbs.org/newshour/politics/how-elon-musk-uses-his-x-social-media-platform-to-amplify-right-wing-views. Accessed: 2024-08-21.

[64] Tim Murphy. I read everything Elon Musk posted for a week. Send help., May 2024. URL https://www.motherjones.com/politics/2024/05/i-read-everything-elon-musk-posted-for-a-week-send-help. Accessed: 2024-08-21.

[65] Oliver Darcy. Radicalized by the right: Elon Musk puts his conspiratorial thinking on display for the world to see, March 2024. URL https://www.cnn.com/2024/03/19/media/elon-musk-don-lemon-interview-analysis-hnk-intl/index.html. Accessed: 2024-08-21.

[66] Sara Fischer. First look: Meta won't recommend political content on Threads, February 2024. URL https://www.axios.com/2024/02/09/meta-political-content-moderation-threads. Accessed: 2024-08-21.

[67] Platform Transparency: Understanding the Impact of Social Media | United States Senate Committee on the Judiciary, May 2022. URL https://www.judiciary.senate.gov/committee-activity/hearings/platform-transparency-understanding-the-impact-of-social-media. [Accessed: 23. Aug. 2024].

[68] Rebekah Tromble. Where Have All the Data Gone? A Critical Reflection on Academic Digital Research in the Post-API Age. *Social Media + Society*, 7(1):2056305121988929, 2021. URL https://doi.org/10.1177/2056305121988929.

[69] Deen Freelon. Computational Research in the Post-API Age. *Political Communication*, 2018. URL https://www.tandfonline.com/doi/full/10.1080/10584609.2018.1477506.

[70] Tuuli Toivonen, Vuokko Heikinheimo, Christoph Fink, Anna Hausmann, Tuomo Hiippala, Olle Järv, Henrikki Tenkanen, and Enrico Di Minin. Social media data for conservation science: A methodological overview. *Biol Conserv*, 233:298–315, 2019. URL https://doi.org/10.1016/j.biocon.2019.01.023.

[71] Jamie Bedson, Laura A Skrip, Danielle Pedi, Sharon Abramowitz, Simone Carter, Mohamed F Jalloh, Sebastian Funk, Nina Gobat, Tamara Giles-Vernick, Gerardo Chowell, João Rangel, Rania Elessawi, Samuel V Scarpino, Ross A Hammond, Sylvie Briand, Joshua M Epstein, Laurent Hébert-Dufresne,

and Benjamin M Althouse. A review and agenda for integrated disease models including social and behavioural factors. *Nature Human Behaviour*, 5(7):834–846, 2021. URL https://doi.org/10.1038/s41562-021-01136-2.

[72] J Sooknanan and D. M. G. Comissiong. Trending on social media: Integrating social media into infectious disease dynamics. *Bulletin of Mathematical Biology*, 82(7), 2020. URL https://doi.org/10.1007/s11538-020-00757-4.

[73] Timon Elmer. Computational social science is growing up: why puberty consists of embracing measurement validation, theory development, and open science practices. *EPJ Data Sci*, 12(1):1–19, 2023. URL https://doi.org/10.1140/epjds/s13688-023-00434-1.

[74] Derek Ruths and Jürgen Pfeffer. Social media for large studies of behavior. *Science*, 346(6213):1063–1064, 2014. URL https://doi.org/10.1126/science.346.6213.1063.

[75] Carter T. Butts. Revisiting the Foundations of Network Analysis. *Science*, 325(5939):414–416, 2009. URL https://doi.org/10.1126/science.1171022.

[76] Rachith Aiyappa, Matthew R. DeVerna, Manita Pote, Bao Tran Truong, Wanying Zhao, David Axelrod, Aria Pessianzadeh, Zoher Kachwala, Munjung Kim, Ozgur Can Seckin, Minsuk Kim, Sunny Gandhi, Amrutha Manikonda, Francesco Pierri, Filippo Menczer, and Kai-Cheng Yang. A Multi-Platform Collection of Social Media Posts about the 2022 U.S. Midterm Elections. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 17, pages 981–989, June 2023. URL https://ojs.aaai.org/index.php/ICWSM/article/view/22205.

[77] Dorian Quelle and Alexandre Bovet. Bluesky: Network Topology, Polarization, and Algorithmic Curation. arXiv Preprint, 2024. URL https://doi.org/10.48550/arXiv.2405.17571.

[78] Andrea Failla and Giulio Rossetti. "I'm in the Bluesky Tonight"': Insights from a Year Worth of Social Data. arXiv Preprint, 2024. URL https://doi.org/10.48550/arXiv.2404.18984.

[79] Bluesky. Firehose API, March 2024. URL https://docs.bsky.app/docs/advanced-guides/firehose. [Accessed: 8. Mar. 2024].

[80] Erfan Samieyan Sahneh, Gianluca Nogara, Matthew R. DeVerna, Nick Liu, Luca Luceri, Filippo Menczer, Francesco Pierri, and Silvia Giordano. The Dawn of Decentralized Social Media: An Exploration of Bluesky's Public Opening. arXiv Preprint, 2024. URL https://doi.org/10.48550/arXiv.2408.03146.

[81] Timestamps — Bluesky, October 2024. URL https://docs.bsky.app/docs/advanced-guides/timestamps. [Online; accessed 2. Oct. 2024].

[82] Seth A. Myers and Jure Leskovec. The bursty dynamics of the Twitter information network. In *WWW '14: Proceedings of the 23rd International Conference on World Wide Web*, pages 913–924. Association for Computing Machinery, April 2014. URL https://doi.org/10.1145/2566486.2568043.

[83] Riley Crane and Didier Sornette. Robust dynamic classes revealed by measuring the response function of a social system. *Proc Natl Acad Sci U.S.A*, 105(41):15649–15653, 2008. URL https://doi.org/10.1073/pnas.0803685105.

[84] Daniele Notarmuzi, Claudio Castellano, Alessandro Flammini, Dario Mazzilli, and Filippo Radicchi. Universality, criticality and complexity of information propagation in social media. *Nature Communications*, 13(1308):1–8, 2022. URL https://doi.org/10.1038/s41467-022-28964-8.

[85] David Y. Hancock, Jeremy Fischer, John Michael Lowe, Winona Snapp-Childs, Marlon Pierce, Suresh Marru, J. Eric Coulter, Matthew Vaughn, Brian Beck, Nirav Merchant, Edwin Skidmore, and Gwen Jacobs. Jetstream2: Accelerating cloud computing via jetstream. In *Practice and Experience in Advanced Research Computing (PEARC '21)*, pages 1–8, New York, NY, USA, 2021. Association for Computing Machinery, Association for Computing Machinery. doi: 10.1145/3437359.3465565.

[86] Timothy J. Boerner, Stephen Deems, Thomas R. Furlani, Shelley L. Knuth, and John Towns. ACCESS: Advancing Innovation: NSF's Advanced Cyberinfrastructure Coordination Ecosystem: Services & Support. In *Proceedings of the Practice and Experience in Advanced Research Computing (PEARC '23)*, page 4. Association for Computing Machinery (ACM), July 2023. URL https://doi.org/10.1145/3569951.3597559.

[87] J. Priem, H. Piwowar, and R. Orr. OpenAlex: A fully-open index of scholarly works, authors, venues, institutions, and concepts. Preprint[arXiv], 2022. URL https://arxiv.org/abs/2205.01833.

Figure S1: **Research on information diffusion and social media has grown rapidly since the early 2000s across various fields.** The barplot in the top left panel displays the cumulative number of peer-reviewed publications across various academic fields, from 2006 to 2023. The time series in the bottom right panel breaks down publication trends annually over the same period.

# Supplementary Information

## Bibliographic analysis

We analyze bibliographic data from OpenAlex (87) to track growing interest in social media information diffusion, shown in Figure S1. To collect this data we query the `search-works` endpoint of the OpenAlex Application Programming Interface (API).[1] We employ boolean search parameters to match a wide range of publications that are clearly related to information diffusion on popular social media platforms. We employ the search query

> ("information diffusion" OR "diffusion of information" OR "information spread" OR "spread of information") AND ("social media" OR "facebook" OR "twitter" OR "reddit")

to return entities that find exact matches (case insensitive) within titles, abstracts, or full text. Our search was not limited by time, retrieving all works in the database that matched our query. Despite the fact that only a subset of the OpenAlex database contains full text, we obtained 19,294 matching works. For our analysis, we narrowed the focus to peer-reviewed articles and conference publications from 2006 onward, the year Facebook opened to the public.[2] This filtering removed 6,723 works from the dataset initially returned by our query.

The ten fields with the smallest number of publications that were grouped together in Figure S1 to create the "Other" group are: "Neuroscience" ($n = 33$); "Earth and Planetary Sciences" ($n = 16$); "Immunology

and Microbiology" ($n = 15$); "Dentistry" ($n = 14$); "Pharmacology, Toxicology and Pharmaceutics" ($n = 5$); "Energy" ($n = 4$); "Nursing" ($n = 3$); "Materials Science" ($n = 1$); "Chemistry" ($n = 1$); "Veterinary" ($n = 1$). Collectively, these papers account for 0.74% of the publications in our collection.

## Node strength correlations

Table S1 presents the statistics for the mean and standard deviation of node strength correlation values between naive and reconstructed networks. Correlations are calculated using Spearman's $\rho$.

| $\gamma$ | $\alpha$ | Twitter | | Bluesky | |
|---|---|---|---|---|---|
| | | $\bar{\rho}$ | $\sigma$ | $\bar{\rho}$ | $\sigma$ |
| 0.25 | 1.1 | 0.2013 | 0.0009 | 0.4882 | 0.0033 |
| 0.25 | 2.0 | 0.1904 | 0.0011 | 0.4659 | 0.0032 |
| 0.25 | 3.0 | 0.1865 | 0.0011 | 0.4513 | 0.0034 |
| 0.5 | 1.1 | 0.2588 | 0.0007 | 0.5469 | 0.0029 |
| 0.5 | 2.0 | 0.2521 | 0.0008 | 0.5327 | 0.0030 |
| 0.5 | 3.0 | 0.2496 | 0.0006 | 0.5249 | 0.0031 |
| 0.75 | 1.1 | 0.3269 | 0.0007 | 0.6060 | 0.0031 |
| 0.75 | 2.0 | 0.3232 | 0.0007 | 0.5991 | 0.0029 |
| 0.75 | 3.0 | 0.3217 | 0.0008 | 0.5951 | 0.0029 |

Table S1: Mean and standard deviation of Spearman's correlations between node strengths of naive and reconstructed networks.

## Comparing cascade metric distributions

In the Information cascade structure section of the main text, we calculate the depth, breadth, and structural virality of cascades and compare the distributions of these metrics across different reconstruction approaches. To determine whether these distributions differ significantly, we perform Kolmogorov-Smirnov two-sample tests for all possible comparisons. The results for depth, maximum breadth, and structural virality are presented in Tables S2, S3, and S4, respectively. To account for multiple comparisons, we apply a Bonferroni correction across 45 comparisons, treating each metric as its own family of comparisons. Out of 135 possible comparisons, 122 (90%) are significant ($P < 0.05$).

Table S2: Kolmogorov-Smirnoff statistics for comparing depth distributions. Rows containing "TID" represent comparisons to distributions based on the Time-Inferred Diffusion method. All values are rounded to two decimal points.

| # | $\gamma_1$ | $\alpha_1$ | $\gamma_2$ | $\alpha_2$ | statistic | $P$ | $P$ adj.[†] | Sig. |
|---|---|---|---|---|---|---|---|---|
| 1 | 0.25 | 1.10 | 0.25 | 2.00 | 0.02 | 0.00 | 0.00 | *** |
| 2 | 0.25 | 1.10 | 0.25 | 3.00 | 0.04 | 0.00 | 0.00 | *** |
| 3 | 0.25 | 1.10 | 0.50 | 1.10 | 0.05 | 0.00 | 0.00 | *** |
| 4 | 0.25 | 1.10 | 0.50 | 2.00 | 0.04 | 0.00 | 0.00 | *** |
| 5 | 0.25 | 1.10 | 0.50 | 3.00 | 0.03 | 0.00 | 0.00 | *** |
| 6 | 0.25 | 1.10 | 0.75 | 1.10 | 0.12 | 0.00 | 0.00 | *** |
| 7 | 0.25 | 1.10 | 0.75 | 2.00 | 0.11 | 0.00 | 0.00 | *** |
| 8 | 0.25 | 1.10 | 0.75 | 3.00 | 0.10 | 0.00 | 0.00 | *** |
| 9 | 0.25 | 1.10 | TID | TID | 0.35 | 0.00 | 0.00 | *** |
| 10 | 0.25 | 2.00 | 0.25 | 3.00 | 0.02 | 0.00 | 0.00 | ** |
| 11 | 0.25 | 2.00 | 0.50 | 1.10 | 0.08 | 0.00 | 0.00 | *** |
| 12 | 0.25 | 2.00 | 0.50 | 2.00 | 0.06 | 0.00 | 0.00 | *** |
| 13 | 0.25 | 2.00 | 0.50 | 3.00 | 0.05 | 0.00 | 0.00 | *** |
| 14 | 0.25 | 2.00 | 0.75 | 1.10 | 0.14 | 0.00 | 0.00 | *** |
| 15 | 0.25 | 2.00 | 0.75 | 2.00 | 0.13 | 0.00 | 0.00 | *** |
| 16 | 0.25 | 2.00 | 0.75 | 3.00 | 0.12 | 0.00 | 0.00 | *** |
| 17 | 0.25 | 2.00 | TID | TID | 0.36 | 0.00 | 0.00 | *** |
| 18 | 0.25 | 3.00 | 0.50 | 1.10 | 0.09 | 0.00 | 0.00 | *** |
| 19 | 0.25 | 3.00 | 0.50 | 2.00 | 0.07 | 0.00 | 0.00 | *** |
| 20 | 0.25 | 3.00 | 0.50 | 3.00 | 0.06 | 0.00 | 0.00 | *** |
| 21 | 0.25 | 3.00 | 0.75 | 1.10 | 0.15 | 0.00 | 0.00 | *** |
| 22 | 0.25 | 3.00 | 0.75 | 2.00 | 0.14 | 0.00 | 0.00 | *** |
| 23 | 0.25 | 3.00 | 0.75 | 3.00 | 0.14 | 0.00 | 0.00 | *** |
| 24 | 0.25 | 3.00 | TID | TID | 0.37 | 0.00 | 0.00 | *** |
| 25 | 0.50 | 1.10 | 0.50 | 2.00 | 0.02 | 0.00 | 0.00 | *** |
| 26 | 0.50 | 1.10 | 0.50 | 3.00 | 0.03 | 0.00 | 0.00 | *** |
| 27 | 0.50 | 1.10 | 0.75 | 1.10 | 0.07 | 0.00 | 0.00 | *** |
| 28 | 0.50 | 1.10 | 0.75 | 2.00 | 0.06 | 0.00 | 0.00 | *** |
| 29 | 0.50 | 1.10 | 0.75 | 3.00 | 0.05 | 0.00 | 0.00 | *** |
| 30 | 0.50 | 1.10 | TID | TID | 0.31 | 0.00 | 0.00 | *** |
| 31 | 0.50 | 2.00 | 0.50 | 3.00 | 0.01 | 0.00 | 0.18 | |
| 32 | 0.50 | 2.00 | 0.75 | 1.10 | 0.08 | 0.00 | 0.00 | *** |
| 33 | 0.50 | 2.00 | 0.75 | 2.00 | 0.07 | 0.00 | 0.00 | *** |
| 34 | 0.50 | 2.00 | 0.75 | 3.00 | 0.07 | 0.00 | 0.00 | *** |
| 35 | 0.50 | 2.00 | TID | TID | 0.32 | 0.00 | 0.00 | *** |
| 36 | 0.50 | 3.00 | 0.75 | 1.10 | 0.09 | 0.00 | 0.00 | *** |
| 37 | 0.50 | 3.00 | 0.75 | 2.00 | 0.09 | 0.00 | 0.00 | *** |
| 38 | 0.50 | 3.00 | 0.75 | 3.00 | 0.08 | 0.00 | 0.00 | *** |
| 39 | 0.50 | 3.00 | TID | TID | 0.33 | 0.00 | 0.00 | *** |
| 40 | 0.75 | 1.10 | 0.75 | 2.00 | 0.01 | 0.03 | 1.00 | |
| 41 | 0.75 | 1.10 | 0.75 | 3.00 | 0.02 | 0.00 | 0.00 | ** |
| 42 | 0.75 | 1.10 | TID | TID | 0.31 | 0.00 | 0.00 | *** |
| 43 | 0.75 | 2.00 | 0.75 | 3.00 | 0.01 | 0.19 | 1.00 | |
| 44 | 0.75 | 2.00 | TID | TID | 0.31 | 0.00 | 0.00 | *** |
| 45 | 0.75 | 3.00 | TID | TID | 0.31 | 0.00 | 0.00 | *** |

Significance codes: *** $P < 0.001$, ** $P < 0.01$, * $P < 0.05$
† Using Bonferroni's method with 45 comparisons

Table S3: Kolmogorov-Smirnoff statistics for comparing maximum breadth distributions. Rows containing "TID" represent comparisons to distributions based on the Time-Inferred Diffusion method. All values are rounded to two decimal points.

| # | $\gamma_1$ | $\alpha_1$ | $\gamma_2$ | $\alpha_2$ | statistic | $P$ | $P$ adj.[†] | Sig. |
|---|---|---|---|---|---|---|---|---|
| 1 | 0.25 | 1.10 | 0.25 | 2.00 | 0.02 | 0.00 | 0.00 | *** |
| 2 | 0.25 | 1.10 | 0.25 | 3.00 | 0.04 | 0.00 | 0.00 | *** |
| 3 | 0.25 | 1.10 | 0.50 | 1.10 | 0.05 | 0.00 | 0.00 | *** |
| 4 | 0.25 | 1.10 | 0.50 | 2.00 | 0.04 | 0.00 | 0.00 | *** |
| 5 | 0.25 | 1.10 | 0.50 | 3.00 | 0.04 | 0.00 | 0.00 | *** |
| 6 | 0.25 | 1.10 | 0.75 | 1.10 | 0.10 | 0.00 | 0.00 | *** |
| 7 | 0.25 | 1.10 | 0.75 | 2.00 | 0.09 | 0.00 | 0.00 | *** |
| 8 | 0.25 | 1.10 | 0.75 | 3.00 | 0.09 | 0.00 | 0.00 | *** |
| 9 | 0.25 | 1.10 | TID | TID | 0.20 | 0.00 | 0.00 | *** |
| 10 | 0.25 | 2.00 | 0.25 | 3.00 | 0.02 | 0.00 | 0.00 | *** |
| 11 | 0.25 | 2.00 | 0.50 | 1.10 | 0.07 | 0.00 | 0.00 | *** |
| 12 | 0.25 | 2.00 | 0.50 | 2.00 | 0.06 | 0.00 | 0.00 | *** |
| 13 | 0.25 | 2.00 | 0.50 | 3.00 | 0.06 | 0.00 | 0.00 | *** |
| 14 | 0.25 | 2.00 | 0.75 | 1.10 | 0.12 | 0.00 | 0.00 | *** |
| 15 | 0.25 | 2.00 | 0.75 | 2.00 | 0.12 | 0.00 | 0.00 | *** |
| 16 | 0.25 | 2.00 | 0.75 | 3.00 | 0.11 | 0.00 | 0.00 | *** |
| 17 | 0.25 | 2.00 | TID | TID | 0.22 | 0.00 | 0.00 | *** |
| 18 | 0.25 | 3.00 | 0.50 | 1.10 | 0.09 | 0.00 | 0.00 | *** |
| 19 | 0.25 | 3.00 | 0.50 | 2.00 | 0.08 | 0.00 | 0.00 | *** |
| 20 | 0.25 | 3.00 | 0.50 | 3.00 | 0.07 | 0.00 | 0.00 | *** |
| 21 | 0.25 | 3.00 | 0.75 | 1.10 | 0.13 | 0.00 | 0.00 | *** |
| 22 | 0.25 | 3.00 | 0.75 | 2.00 | 0.13 | 0.00 | 0.00 | *** |
| 23 | 0.25 | 3.00 | 0.75 | 3.00 | 0.13 | 0.00 | 0.00 | *** |
| 24 | 0.25 | 3.00 | TID | TID | 0.23 | 0.00 | 0.00 | *** |
| 25 | 0.50 | 1.10 | 0.50 | 2.00 | 0.01 | 0.00 | 0.02 | * |
| 26 | 0.50 | 1.10 | 0.50 | 3.00 | 0.02 | 0.00 | 0.00 | *** |
| 27 | 0.50 | 1.10 | 0.75 | 1.10 | 0.05 | 0.00 | 0.00 | *** |
| 28 | 0.50 | 1.10 | 0.75 | 2.00 | 0.05 | 0.00 | 0.00 | *** |
| 29 | 0.50 | 1.10 | 0.75 | 3.00 | 0.04 | 0.00 | 0.00 | *** |
| 30 | 0.50 | 1.10 | TID | TID | 0.16 | 0.00 | 0.00 | *** |
| 31 | 0.50 | 2.00 | 0.50 | 3.00 | 0.01 | 0.03 | 1.00 | |
| 32 | 0.50 | 2.00 | 0.75 | 1.10 | 0.06 | 0.00 | 0.00 | *** |
| 33 | 0.50 | 2.00 | 0.75 | 2.00 | 0.05 | 0.00 | 0.00 | *** |
| 34 | 0.50 | 2.00 | 0.75 | 3.00 | 0.05 | 0.00 | 0.00 | *** |
| 35 | 0.50 | 2.00 | TID | TID | 0.17 | 0.00 | 0.00 | *** |
| 36 | 0.50 | 3.00 | 0.75 | 1.10 | 0.06 | 0.00 | 0.00 | *** |
| 37 | 0.50 | 3.00 | 0.75 | 2.00 | 0.06 | 0.00 | 0.00 | *** |
| 38 | 0.50 | 3.00 | 0.75 | 3.00 | 0.06 | 0.00 | 0.00 | *** |
| 39 | 0.50 | 3.00 | TID | TID | 0.18 | 0.00 | 0.00 | *** |
| 40 | 0.75 | 1.10 | 0.75 | 2.00 | 0.01 | 0.14 | 1.00 | |
| 41 | 0.75 | 1.10 | 0.75 | 3.00 | 0.01 | 0.00 | 0.08 | |
| 42 | 0.75 | 1.10 | TID | TID | 0.15 | 0.00 | 0.00 | *** |
| 43 | 0.75 | 2.00 | 0.75 | 3.00 | 0.01 | 0.47 | 1.00 | |
| 44 | 0.75 | 2.00 | TID | TID | 0.15 | 0.00 | 0.00 | *** |
| 45 | 0.75 | 3.00 | TID | TID | 0.15 | 0.00 | 0.00 | *** |

Significance codes: *** $P < 0.001$, ** $P < 0.01$, * $P < 0.05$
† Using Bonferroni's method with 45 comparisons

Table S4: Kolmogorov-Smirnoff statistics for comparing structural virality distributions. Rows containing "TID" represent comparisons to distributions based on the Time-Inferred Diffusion method. All values are rounded to two decimal points.

| # | $\gamma_1$ | $\alpha_1$ | $\gamma_2$ | $\alpha_2$ | statistic | $P$ | $P$ adj.[†] | Sig. |
|---|---|---|---|---|---|---|---|---|
| 1 | 0.25 | 1.10 | 0.25 | 2.00 | 0.01 | 0.00 | 0.01 | * |
| 2 | 0.25 | 1.10 | 0.25 | 3.00 | 0.02 | 0.00 | 0.00 | *** |
| 3 | 0.25 | 1.10 | 0.50 | 1.10 | 0.05 | 0.00 | 0.00 | *** |
| 4 | 0.25 | 1.10 | 0.50 | 2.00 | 0.04 | 0.00 | 0.00 | *** |
| 5 | 0.25 | 1.10 | 0.50 | 3.00 | 0.04 | 0.00 | 0.00 | *** |
| 6 | 0.25 | 1.10 | 0.75 | 1.10 | 0.09 | 0.00 | 0.00 | *** |
| 7 | 0.25 | 1.10 | 0.75 | 2.00 | 0.08 | 0.00 | 0.00 | *** |
| 8 | 0.25 | 1.10 | 0.75 | 3.00 | 0.08 | 0.00 | 0.00 | *** |
| 9 | 0.25 | 1.10 | TID | TID | 0.18 | 0.00 | 0.00 | *** |
| 10 | 0.25 | 2.00 | 0.25 | 3.00 | 0.01 | 0.02 | 0.89 | |
| 11 | 0.25 | 2.00 | 0.50 | 1.10 | 0.06 | 0.00 | 0.00 | *** |
| 12 | 0.25 | 2.00 | 0.50 | 2.00 | 0.05 | 0.00 | 0.00 | *** |
| 13 | 0.25 | 2.00 | 0.50 | 3.00 | 0.05 | 0.00 | 0.00 | *** |
| 14 | 0.25 | 2.00 | 0.75 | 1.10 | 0.10 | 0.00 | 0.00 | *** |
| 15 | 0.25 | 2.00 | 0.75 | 2.00 | 0.09 | 0.00 | 0.00 | *** |
| 16 | 0.25 | 2.00 | 0.75 | 3.00 | 0.09 | 0.00 | 0.00 | *** |
| 17 | 0.25 | 2.00 | TID | TID | 0.19 | 0.00 | 0.00 | *** |
| 18 | 0.25 | 3.00 | 0.50 | 1.10 | 0.07 | 0.00 | 0.00 | *** |
| 19 | 0.25 | 3.00 | 0.50 | 2.00 | 0.06 | 0.00 | 0.00 | *** |
| 20 | 0.25 | 3.00 | 0.50 | 3.00 | 0.06 | 0.00 | 0.00 | *** |
| 21 | 0.25 | 3.00 | 0.75 | 1.10 | 0.10 | 0.00 | 0.00 | *** |
| 22 | 0.25 | 3.00 | 0.75 | 2.00 | 0.10 | 0.00 | 0.00 | *** |
| 23 | 0.25 | 3.00 | 0.75 | 3.00 | 0.10 | 0.00 | 0.00 | *** |
| 24 | 0.25 | 3.00 | TID | TID | 0.20 | 0.00 | 0.00 | *** |
| 25 | 0.50 | 1.10 | 0.50 | 2.00 | 0.01 | 0.01 | 0.54 | |
| 26 | 0.50 | 1.10 | 0.50 | 3.00 | 0.02 | 0.00 | 0.00 | *** |
| 27 | 0.50 | 1.10 | 0.75 | 1.10 | 0.06 | 0.00 | 0.00 | *** |
| 28 | 0.50 | 1.10 | 0.75 | 2.00 | 0.05 | 0.00 | 0.00 | *** |
| 29 | 0.50 | 1.10 | 0.75 | 3.00 | 0.05 | 0.00 | 0.00 | *** |
| 30 | 0.50 | 1.10 | TID | TID | 0.16 | 0.00 | 0.00 | *** |
| 31 | 0.50 | 2.00 | 0.50 | 3.00 | 0.01 | 0.19 | 1.00 | |
| 32 | 0.50 | 2.00 | 0.75 | 1.10 | 0.06 | 0.00 | 0.00 | *** |
| 33 | 0.50 | 2.00 | 0.75 | 2.00 | 0.06 | 0.00 | 0.00 | *** |
| 34 | 0.50 | 2.00 | 0.75 | 3.00 | 0.06 | 0.00 | 0.00 | *** |
| 35 | 0.50 | 2.00 | TID | TID | 0.17 | 0.00 | 0.00 | *** |
| 36 | 0.50 | 3.00 | 0.75 | 1.10 | 0.07 | 0.00 | 0.00 | *** |
| 37 | 0.50 | 3.00 | 0.75 | 2.00 | 0.07 | 0.00 | 0.00 | *** |
| 38 | 0.50 | 3.00 | 0.75 | 3.00 | 0.06 | 0.00 | 0.00 | *** |
| 39 | 0.50 | 3.00 | TID | TID | 0.17 | 0.00 | 0.00 | *** |
| 40 | 0.75 | 1.10 | 0.75 | 2.00 | 0.01 | 0.64 | 1.00 | |
| 41 | 0.75 | 1.10 | 0.75 | 3.00 | 0.01 | 0.16 | 1.00 | |
| 42 | 0.75 | 1.10 | TID | TID | 0.16 | 0.00 | 0.00 | *** |
| 43 | 0.75 | 2.00 | 0.75 | 3.00 | 0.00 | 0.95 | 1.00 | |
| 44 | 0.75 | 2.00 | TID | TID | 0.16 | 0.00 | 0.00 | *** |
| 45 | 0.75 | 3.00 | TID | TID | 0.16 | 0.00 | 0.00 | *** |

Significance codes: *** $P < 0.001$, ** $P < 0.01$, * $P < 0.05$

† Using Bonferroni's method with 45 comparisons