*Journal of Complex Networks* (2013) Page 1 of 22 doi:10.1093/comnet/cnt006

# Cascading behaviour in complex socio-technical networks

JAVIER BORGE-HOLTHOEFER\* AND RAQUEL A. BAÑOS

Instituto de Biocomputación y Física de Sistemas Complejos (BIFI), Universidad de Zaragoza, Mariano Esquillor s/n, 50018 Zaragoza, Spain \*Corresponding author: borge.holthoefer@gmail.com

Sandra González-Bailón

Oxford Internet Institute, University of Oxford, 1 St. Giles, Oxford OX1 3JS, UK

AND

YAMIR MORENO

Instituto de Biocomputación y Física de Sistemas Complejos (BIFI), Universidad de Zaragoza, Mariano Esquillor s/n, 50018 Zaragoza, Spain Departamento de Física Teórica, Universidad de Zaragoza, 50009 Zaragoza, Spain

Edited by: Ernesto Estrada

[Received on 21 March 2013; revised on 25 March 2013; accepted on 25 March 2013]

Most human interactions today take place with the mediation of information and communications technology. This is extending the boundaries of interdependence: the group of reference, ideas and behaviour to which people are exposed is larger and less restricted to old geographical and cultural boundaries; but it is also providing more and better data with which to build more informative models on the effects of social interactions, amongst them, the way in which contagion and cascades diffuse in social networks. Online data are not only helping us gain deeper insights into the structural complexity of social systems, they are also illuminating the consequences of that complexity, especially around collective and temporal dynamics. This paper offers an overview of the models and applications that have been developed in what is still a nascent area of research, as well as an outline of immediate lines of work that promise to open new vistas in our understanding of cascading behaviour in social networks.

Keywords: contagion; diffusion; social influence; computational social science; big data.

## 1. Introduction

By the end of 2012, Facebook had 1.06 billion monthly active users [1]. Over 60% of them were active on a daily basis. Twitter claims to have 200 million users [2] producing over 400 million tweets each day; and Google+ is the fastest-growing network ever with over 400 million subscribed users, 25% of them active [3]. According to the International Telecommunication Union, more than 6 billion mobile phones around the world are currently in use [4]. All these figures are indicative of the radical transformation that affects how we interact and communicate, but also how we confront the research of those communication patterns: we can now capitalize on massive amounts of data to advance theoretical approaches that, so far, had to rely on small datasets or analytical models lacking in external

<sup>©</sup> The Authors 2013. Published by Oxford University Press. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (http://creativecommons.org/licenses/by-nc/3.0/), which permits non-commercial reuse, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial reuse, please contact journals.permissions@oup.com

validity. In the middle of this transition, a new scientific paradigm [5] is brewing under the label of 'computational social science' (CSS) [6–8], a shorthand for the new avenues of research that the massive amounts of data generated by information and communications technology (ICTs) are opening up. This new paradigm brings together insights from physics, computer science and mathematics to revisit old theoretical questions at the core of social science research. Prominent amongst them are the structural properties of social systems, and the dynamical consequences of those structures, a question that has attracted the attention of social scientists for decades [9,10] and that the science of networks, powered by digital data, has contributed to advance to new exciting grounds [11]. This article aims to offer a perspective of some of the models and findings that are emerging at the intersection of those disciplines; in particular, it offers a survey review of the models and theories that have focused attention on one important dimension of social systems: cascading behaviour in contagion processes, and how the dynamics relate to the network topology in which they take place.

As a paradigm, CSS is uniquely placed to tackle that question, if only because it is based on the recognition that exploiting the new datasets now available is not the domain of a single discipline. The sheer volume of the data demands the joint efforts of approaches that can handle the logistical needs of storing and processing data in an efficient way, and also-most importantly-that can make sense of all that information. Often referred to as Big Data (where 'big' refers to storage size but also, and more specifically, to the higher spatial and temporal resolution of the data, and the granularity of observations), these new sources of information require an efficient approach to data manipulation and exploitation; but also new theoretical tools to model and scrutinize their inherent complexity-that is, the complexity of the social dynamics that the data track with unprecedented fidelity. This requires devising models that span the different levels of analysis that can now be analysed simultaneously the micro to macro link to which social scientists refer to illuminate the connection between individual actions and collective behaviour [12]; and also devising models that help identify the right time resolution now that longitudinal dynamics can be tracked down to the second. The success of this endeavour promises radical changes in how we think about innovation, economic growth, governance, health interventions and even political revolutions-all core issues in the social science agenda, and of great potential impact for their policy implications.

Before those changes are possible, however, we have to deal with many unanswered questions about the mechanisms of complexity, how they manifest in social systems and whether it is possible to harness those drivers to capitalize on the power of decentralized networks [13–15]. This requires breaking up the problem into its constituent parts, that is, into the different aspects of the connection linking individual actions and collective behaviour, as mediated by networks. One of the ways in which individual behaviour can lead to unanticipated collective dynamics is by means of social influence, or social contagion. This is also one of the most visible examples of how the confluence of different disciplinary backgrounds can help blaze new empirical trails. An increasing body of work, borrowing theories and models from a range of research traditions, considers the influence of network topology on the unfolding of cascades or chain reactions that start with an initial seed (or a set of them), placed randomly or in correlation with some network property. What follows aims to map those developments, with a particular focus on models and findings fed with the data that ICTs yield, and to give the coordinates of where we are and where we could be if this line of work is pursued further.

The article starts with an overview of the 'new' science of networks (where 'new' is used to mark the phase transition that Big Data caused in network research); the aim is to lay down the basic building blocks coming from graph theory on which research on cascading behaviour is based. Section 3 reviews the main theoretical approaches to the study of contagious behaviour, and the type of analytical models that were developed in the absence of better empirical data. These approaches attack diffusion dynamics from different fronts, and assuming different mechanisms; the most prominent are epidemic, threshold and rumour models. Many of these models, however, were developed under the limitations imposed by the need of analytical tractability; these limitations have now been levied by the availability of large datasets, which makes it possible to revisit many theoretical assumptions through the lens of better observational evidence. The insights thus gained are summarized in Section 4, which reviews some prominent recent studies on social influence and contagion in social networks. This empirical work has fed back on the development of theoretical models, mostly by means of generalizations drawn from the data on two fronts: cascade size distributions (with a focus on the frequency of large cascades, that is, chain reactions that percolate to reach system-wide proportions); and the topological underpinning of influence (that is, the structural roots of large cascades, or the position where cascades tend to start in the overall network structure). These findings have helped refine epidemic, threshold and rumour models in powerful ways. On the basis of these findings, future lines of research are outlined in the last section of the paper, which also considers some of the practical implications of this line of work.

#### 2. The building blocks of network science

The structure of networks has been studied using the language of graph theory, a branch of mathematics. There are many excellent reviews and books in the literature about the structure and dynamics of complex networks [16–19]. Here, we give an overview of the network features that are relevant for the work mentioned in subsequent sections.

A graph is a mathematical abstraction consisting of a set of N nodes or vertices, connected by a set of E edges or links. Nodes are usually depicted as labelled circles, and lines between them represent existing relationships. For example, a molecule can be thought of as a network where nodes are atoms and links are bonds between them. In the social realm, nodes tend to be people (but can also represent countries, or organizations), and connections map their interactions (for instance, communication, but also trade or collaboration in the case of countries and organizations). Networks can be classified according to several topological properties, but the simplest classification relies on the nature of the interactions. Networks can be directed or undirected, depending on whether the directionality of connections matters for the analysis and interpretation. Examples of directed networks include the World Wide Web, airline route maps, flow charts or even binary relations in mathematics. The most usual directed social network maps the structure of friendship: nominations can be reciprocated or asymmetrical. As far as the interaction strength is concerned, links may be weighted or unweighted: if two individuals talk frequently, their tie will be stronger (or heavier) than if they talk only occasionally; a weighted network records the information of this frequency of interaction.

Every graph may be represented in a matrix notation, through the so-called adjacency matrix, A, which is a  $N \times N$  matrix where the entries  $a_{ij} = w_{ij}$  indicate the existence of a link of weight  $w_{ij}$  from vertex *i* to *j*. Adjacency matrices standing for undirected networks are symmetric,  $a_{ij} = a_{ji}$ , whereas unweighted networks are represented by binary matrices,  $a_{ij} \in 0, 1$ .

The simplest and most extensively studied property of a node or vertex in a graph is the connectivity or degree of a node *i*,  $k_i$ , which counts the number of edges connecting node *i* to other nodes in the network. In directed graphs, the degree is often measured as indegree (the number of connections ending on a vertex) and outdegree (the number of connections starting at a vertex). If all vertices in a graph have the same degree,  $k_i = k$ , the graph is designated as a regular graph, where the degree probability distribution, P(k), is a Dirac delta function,  $P(k) = \delta(k_i - k)$ . Although connectivity is a local property of vertices, degree distributions often determine some important global characteristics of networks, and they help devise a classification according to the homogeneity of the degree distribution. Degree

distributions of homogeneous networks are characterized by a tail of P(k) that decays exponentially fast. The Erdös–Rényi model [20] is one of the most used methods for generating this type of graphs, where links are formed randomly according to a Poisson probability distribution, assuming the independence of dyads. On the other side, heterogeneous networks display tail distributions decaying as a power law,  $P(k) \sim k^{-\gamma}$ . The Barabási–Albert model [21] is a paradigmatic algorithm that uses a preferential attachment mechanism to generate scale-free random graphs; according to this mechanism, and under the assumption of network growth, well-connected nodes are a more likely target for new connections than nodes with lower connectivity.

The redundancy of connections amongst the neighbours of a vertex i is another important structural property of networks, and is usually described by the clustering coefficient,  $C_i = t_i/[k_i(k_i - 1)/2]$ . This coefficient is the quotient of the number of existing triangles attached to node *i*,  $t_i$ , out of the maximum possible number of such triangles,  $k_i(k_i - 1)/2$ , and measures the degree to which nodes in a graph tend to cluster together. The coefficient can also be calculated on the global level, as an average of the local coefficients or as the quotient of the total number of closed triads out of all possible triplets. This type of link redundancy is important because it provides social reinforcement for adoption [22].

Another important characteristic of a vertex is its betweenness centrality, a measure based on the concept of shortest path,  $l_{ij}$ . For every pair of nodes *i* and *j*, the shortest path is the minimum, in terms of the number of hops, of the possible paths starting at node *i* and ending at *j*. It just corresponds to the number of edges comprising the path, and allows another categorization of networks in two groups: connected graphs, where there is at least one shortest path between all pairs of vertices ('strongly connected' if directionality is taken into account); and non-connected graphs, which are made up of broken subgraphs. The study of the giant connected component in a network and the size distribution of finite connected subgraphs offers another dimension in the description of network topology. The betweenness centrality of a vertex relies on this concept of global connectivity and is defined as follows: let s(i,j) > 0 be the number of shortest paths between vertices *i* and *j*, and let s(i, v, j) be the number of shortest paths between other nodes that run through vertex *v*, normalized by the total number of shortest paths. This metric provides information about the importance of a node in terms of the relative distance to the rest of the network, and therefore in terms of how central it is in the flow allowed by the network.

The *k*-core offers another local property that relies on global network structure. This metric gauges the existence of cohesive subgroups of nodes in a network. The network can be seen as a set of successively enclosed substructures or *k*-cores, comprising vertices having at least degree *k*. This partition of the whole graph assigns an integer number to every node in the network obtained by a recursive pruning of the vertices. One starts with isolated nodes, which are assigned a  $k_c = 0$ . Then, vertices with k = 1 are removed along with their links, and assigned  $k_c = 1$ . If any of the remaining nodes is left with *k* connections, it is also removed and contained in the  $k_c = 1$  core. The process continues with  $k_c = 2, 3, \ldots$  until every node has been assigned to a  $k_c$  shell. This measure of centrality goes beyond degree because it takes into account the centrality of a vertex neighbours to define the centrality of that vertex.

A more sophisticated version of the degree centrality is the so-called eigenvector centrality [23]. Defining the vector of centralities  $x = (x_1, x_2, ...)$ , we can rewrite this equation in matrix form as  $\lambda x = Ax$  and hence we see that x is an eigenvector of the adjacency matrix with eigenvalue  $\lambda$ . Assuming that we wish the centralities to be non-negative,  $\lambda$  must be the largest eigenvalue of the adjacency matrix and x the corresponding eigenvector. Eigenvector centrality assigns relative scores to all nodes in the network based on the principle that connections to high-scoring nodes contribute more to the score of the node in question than equal connections to low-scoring nodes; interestingly, it turns out to be a revealing measure in many situations. For example, a variant of eigenvector centrality (PageRank, [24])

is employed by the well-known Web search engine Google to rank Web pages and has been usefully applied in other contexts.

### 3. Models and analytical solutions

The study of diffusion dynamics has a long tradition in the social sciences [25]. Most of these studies, however, are based on aggregated data of adoption rates, which are compatible with a number of individual-level mechanisms, including learning, externalities, contagion or influence [26]. The study of contagion dynamics often makes explicit the effects of network structure on adoption rates: the assumption is that information or behaviour diffuses in a population because adopters are exposed to previous adopters via their networks, which delineate the boundaries of their group of reference.

Unlike social influence, which can also derive from exposure to a common source of information like mass media, contagion assumes that influence dynamics are channelled locally, through the paths that networks open. However, network data are often lacking from the earlier diffusion studies—with a few exceptions [27,28]—which forces the identification of contagion effects using proxies like geographical distance or bursts of activity [29,30]. ICTs are generating the kind of network data that were missing before; however, the nature of social networks (and in particular the composition of the group of reference to which individuals are exposed) is likely to have changed compared with how social networks operated before, especially as measured using surveys or census data. Prior to the irruption of the Internet, networks were more local and narrower; for this reason, previous empirical analyses of diffusion offer an inappropriate benchmark for comparison with online contagion: it is too contingent on the data and contextual circumstances analysed. Models that were built to overcome the lack of network data, and explore the generic principles that govern network dynamics, offer a more appropriate point of comparison.

In the absence of appropriate data, simulation and analytical models filled the empirical gap. These models were developed under the influence of three streams of research: threshold models of collective behaviour [31–33], epidemiological studies [34–36] and rumour dynamics [37–39]. All these approaches to contagion revolve around a common mechanism: an agent (in the 'inactive', 'susceptible' or 'ignorant' state) that decides whether to adopt a given behaviour as a function of the neighbouring agents who have already adopted (those 'active', or in the 'infectious' or 'spreader' class). While in epidemic- and rumour-like dynamics the decisions to adopt are taken independently with probability p for each successive contact (these are 'independent interaction models' [40]), in threshold models the decision depends on a critical proportion of previous adoptions: an actor will only join the adoption curve if she registers that the critical proportion is satisfied.

#### 3.1 Threshold models

The first attempt to interweave cascading phenomena and complex networks [32] built on previous work on the diffusion effects of interdependent decision-making [31,41,42]. In this article, Watts provides an analytic approach to discern the conditions under which global cascades may occur in structured sparse topologies. Using percolation methods, the model explores how network topology and individual thresholds interact in the spreading of behaviour. First, a network with an arbitrary degree distribution  $p_k$  is chosen from an ensemble of graphs. Each node is then assigned a fixed threshold  $\phi$  drawn from a distribution  $0 \le f(\phi) \le 1$  and, with the exception of a small initial seeding set, each agent is marked as *inactive*. An agent *i* updates her state calculating the fraction of active neighbours  $a_i/k_i$ : if  $a_i/k_i > \phi_i$ she activates. The simulation evolves following this logic until an equilibrium is reached, i.e. no more updates occur. Given this set-up, the *cascade condition* is derived from the growth of the initial fraction of *innovators*  $\rho_0$ . The simulations show that large cascades can only occur if the subnetwork of early adopters percolates, if the average vulnerable cluster size  $\langle n \rangle$  diverges. Using a generating function approach, this condition is met at

$$G_0'' = \sum_k k(k-1)\rho_k p_k = \langle k \rangle, \qquad (3.1)$$

where 
$$\rho_k = \begin{cases}
1, & k = 0, \\
F(1/k), & k > 0.
\end{cases}$$
(3.2)

For  $G_0'' < \langle k \rangle$  all vulnerable clusters are small, and the seed cannot grow beyond isolated groups of early adopters; on the contrary, a small seed set may unleash—with finite probability—global cascades when  $G_0'' > \langle k \rangle$ . Accordingly, simulations show that cascades are strongly constrained by the network's connectivity: low  $\langle n \rangle$  allows for system-wide cascades (power-law distribution, Fig. 1, top panel, in red) because the bulk of nodes are vulnerable; but rich local connectivities yield large sets of locally stable



FIG. 1. Top panel: Cumulative distributions of cascade sizes  $\rho$  for a Poisson undirected random graph of  $N = 10^5$  nodes and a single early adopter, with mean degrees at the lower ( $\langle k \rangle = 1.05$ ) and upper ( $\langle k \rangle = 6.14$ ) critical points. Cascades at the lower critical point are power-law distributed. Bottom left panel: Average cascade size  $\rho$  (colour-coded) as a function of the constant threshold  $\phi$ , and the average degree  $\langle k \rangle$ , for a seed fraction of  $\rho_0 = 0.01$ . Bottom right panel: Values of  $\rho$  at  $\phi = 0.18$  and different values of seed fractions. Numerical simulations have been averaged over 100 randomizations (each realization consisting of a randomly generated network and a set of  $\rho_0$  randomly selected early adopters).

nodes, hampering the adoption of the new behaviour (exponential decay, Fig. 1, top panel, in green). Left-lower panel in Fig. 1 illustrates perfectly the match between the analytical cascade condition and numerical simulations, with a well-defined region where large cascades are possible.

From this seminal work a good number of developments emerged. Centola et al. [33] devised a threshold model with a hard-wired  $\phi_i$ ,  $\forall i$ . Though a simplification, this set-up allows the authors to explore the existence of *critical points*  $\phi_c$ , such that, given a network and a certain  $\phi > \phi_c$  value, large cascades are likely to happen. The authors derive the value  $\phi_c$  for a wide range of topologies, from regular lattices to scale-free networks. Gleeson & Cahalane [43] extend the cascading condition (Equation (3.1)) to second order, which provides an even more accurate matching between analytical approximation and simulation, and allows quantifying the impact of the size of the initial seed  $\rho_0$  on the probability of obtaining large cascading events. As the initial seed set grows larger, system-wide cascades are possible for a wider range of  $\langle k \rangle$ ; see the right-lower panel in Fig. 1. In another related work, Gleeson [44] examines the cascade condition for correlated and modular random networks. The author develops a general framework which analytically reduces the threshold model to a site (node) and bond (link) percolation process—already hinted in [32]. Most interestingly, numerical simulations of the model on modular networks reveal that large cascades either choke-because communities act as topological traps for their growth—or the rate of activation occurs in a sequence of cascade fronts, which signals the existence of structural bottlenecks. This prediction matches—at least qualitatively—empirical observations, as the next section reviews. Finally, in an effort to approach models to real-world systems, Hacket et al. [45] offer analytical and numerical results for a class of clustered networks. Indeed, most analytical results derive from the class of random networks defined by the so-called *configura*tion model, which renders tree-like (non-clustered) local structures. Their conclusions suggest that, for certain  $\langle k \rangle$  regimes, clustering will decrease (3 >  $\langle k \rangle$  > 29) or increase (3 <  $\langle k \rangle$  < 29) the probability of obtaining large cascades.

## 3.2 Epidemic and rumour models

The mathematics of epidemic spreading were originally developed, unsurprisingly, in the fields of Medicine and Biology [34]. Their application to information cascades has been rather indirect, through physicists and computer scientists who found in epidemic spreading a fecund metaphor of information propagation. This approach assumes that information travels through social networks as viral infections and that personal interactions open the diffusion routes [38,39].

According to these models, contagion dynamics evolve following a simple scheme: at each time step, infected individuals propagate the contagion to susceptible neighbours with probability  $\lambda$ . Additionally, infected individuals can recover at a rate  $\mu$  (as in the susceptible–infected–recovered, or SIR, models); or they can revert to the susceptible state with probability  $\mu$  (as in the susceptible–infected–susceptible, or SIS, models) [36]. These transitions can be expressed as differential equations under a simple form, which yield valuable insights within the framework of complex networks. For instance, Pastor-Satorras & Vespignani [46] analytically established, for the SIS model, that the critical point (or epidemic threshold) in uncorrelated scale-free networks is given by  $\lambda_c = \langle k \rangle / \langle k^2 \rangle$ , leading to  $\lambda_c \to 0$  as  $N \to \infty$  when  $2 < \gamma \leq 3$ . Taking this as a starting point, Leskovec *et al.* [47] exploit epidemic processes to replicate real cascade size distributions in the blogosphere. Tuning the infection probability, they can reproduce the seven most frequent cascades as well as match cascade size distributions.

A different approach to contagion goes deeper into the mechanisms that allow epidemic dynamics to unfold. Unlike what happens with viral epidemics, social contagion relies on the effects of social influence, which is at the core of sociological research [48,49] and has inspired the recent distinction

between simple and complex contagion: information, like viruses, can propagate with a single exposure, but the spread of behaviour often requires multiple exposure from multiples sources [50]; evidence of this type of contagion has been found in a number of online settings [22,51].

A question that has naturally followed from the study of these contagious dynamics is where seeds are in the network topology—that is, if the leaders of the process have a specific network position. Marketing experts seek to find those actors to engineer the spreading of product adoption [52] much in the same way as epidemiologists try to identify the spreaders of a disease, but for the opposite reason [53].

Across all these areas of application, disease spreading (and especially the SIR model) has become a rather usual benchmark to identify the network features—mainly degree and centrality measures that perform better when it comes to spotting outstanding spreaders, i.e. the nodes in the network that trigger larger cascades. This is the case of the work by Kitsak *et al.* [53], for instance, which explores whether the degree of a node k or its k-core can help predict the spreading capabilities of a certain node. They modelled the underlying dynamics using the SIR and SIS frameworks, because of the wide range of real-world phenomena they can be mapped onto. The author's findings indicate that centrality, rather than connectivity, is the key topological feature to understand the spreading power of a node. In addition to the empirical validity of such a claim—further elaborated in the following section—this work triggered a number of efforts to determine which, among centrality descriptors, performed better at spotting influential spreaders. In this vein, Klemm *et al.* [15], for instance, propose a *dynamical influence* (DI) measure which capitalizes on eigenvector centrality—as opposed to a static, purely topological approach. To demonstrate that DI outperforms k-core centrality, they used a variety of benchmarks including the SIR scheme and the voter model in opinion dynamics [54].

Finally, rumour dynamics have also been modelled in parallel to more general models of social influence. These models sprung directly from the canonical SIR, renaming the susceptible, infected and recovered classes (SIR) to ignorant (who has not heard the rumour yet), spreader (who knows the rumour and is 'infecting' ignorants) and stifler (who also knows the rumour, but has decided not to spread it further). Although rumour models are often regarded as a simple mapping of its epidemic counterpart, a number of differences set them apart. First, SIR is an attempt to model a real process, whereas researchers on rumour dynamics-which typically seek to maximize influence for the sake of technological and commercial applications-are free to design the rules of epidemic infection in order to reach the desired result. This affects mainly the transition from spreader to stifler, which can be implemented under different plausible forms. Secondly, rumour models can be applied to social systems the connectivity of which can be changed: for instance, in peer-to-peer file-sharing systems, the connectivity distribution of the nodes can be changed in order to maximize the performance of the protocols, as informed by the models [55]. Thirdly, the dynamics are also different: the transition to the class of 'recovered' in SIR happens spontaneously (at a certain rate), while classical rumour spreading allows the transition to 'stifle' (at a certain rate) only after a 'spreader' interacts with either another spreader or a stifler, i.e. spreaders learn that the rumour has lost its 'news value' when they encounter neighbours already informed. For all these reasons the outcomes of the rumour model may present significant differences when compared with simulations of SIR models. This is the case in Borge-Holthoefer & Moreno [56], who, motivated by the aforementioned theoretical predictions [53] and some empirical findings related to social movements and political mobilization [57], attempted to identify super-influencers in real networks on top of which rumour dynamics were performed. Surprisingly, the subtle differences between SIR and rumour dynamics suffice to flatten the reported 'k-core effect', i.e. cascade size and k-core appear to be uncorrelated; see the left panel in Fig. 2. Contrary to the expectations, hubs and high k-core nodes act as firewalls—they turn stiflers early in the dynamics, Fig. 2 (right)—preventing the diffusion of the rumour to large fractions of the underlying structure. A similar result was obtained



FIG. 2. Left: Average stifler density  $\rho$  for a rumour process triggered at nodes with k-core  $k_c$  on an e-mail network. Different conditions were tested, and yet no correlation is observed between  $\rho$  and the initiator's centrality, i.e. absence of influential spreaders in rumour dynamics. Right: central nodes (those in  $k_c^{max}$ , blue dots) acting as firewalls of the rumour spreading: these nodes are among the first to become stiflers (time-to-stifler t(s) is low), thus acting as topological barriers for the dynamics. For these central nodes, the time it takes them to turn into stiflers is even lower for those with the highest degree (normalized in the *x*-axis by the  $k_{max}$  within a *k*-shell). The contrast is clear if compared with lower cores (red circles) or, in general, to the rest of the network (gray dots). Adapted from [56].

previously [58], under the paradigm of threshold models. These consistent results, obtained under a number of different modelling assumptions, hint at the existence of a class of spreaders (the 'hidden influentials', topologically unexceptional) which hold the key to trigger most system-wide cascades. Two simple variations on classical rumour models (one of them coupling nodes with complex activity patterns [59], the other adding a new transition from ignorant to stifler) recover the observed positive connectivity-cascade reach correlation [60].

For the sake of exposition, we have left out many other modelling approaches which branch off the main ones outlined here, for instance [61–65], among many others. It is important to note that contradicting simulation results do not cancel out at the theoretical level. Incompatible outcomes (i.e. that influencers exist, or that they do not) simply highlight the fact that all models recover real phenomena only partially. As is often the case, the incorporation of empirical data adds important caveats to analytical conclusions and facilitates feedback that helps refine and improve the theoretical models [60]. The following section expands more on the contribution that empirical analyses can make to the study of contagious behaviour.

#### 4. Validation: findings and theoretical developments

There is a fast-growing literature that is now revisiting the theoretical models discussed in the previous section through the lens of the massive datasets generated by e-mail communication, weblogs and social networking sites (SNSs). Other online forms of communications—like telephone calls [66–68], chatrooms or discussion forums—will not be considered in this review, although they have also provided interesting insights. Online data contain information of the relationship between users (the structural dimension of social systems), but also of the dynamics of their interactions, both on the temporal and spatial levels.

Although the properties of online networks often differ drastically from what is known about offline, face-to-face networks, they can often be used as a good proxy to those social networks [69–71].

The extent to which the study of cascades in online networks is applicable to cascades in offline networks is an empirical question, which depends on which online network is being analysed (i.e. the map of informal interactions drawn from email communication in a large organization is a good representation of offline interactions amongst the members of that organization; see [72]). With that caveat, this section aims to identify, on the one hand, the features that can characterize the structural and temporal dimension of cascades in networks; and, on the other, the network statistics that are most useful to predict the likelihood that a cascade will grow viral. The main goal is to highlight consistent findings obtained from different approaches and methods, and to point out where the theoretical predictions discussed in the previous section match (or not) the observed empirical trends.

Delivering these aims presents at least two problems. First, the affordances of online technologies differ from platform to platform: communication in the blogosphere does not follow the same rules as in SNSs. Even within the same platform, some differences might arise (see [73,74]); for instance, information about someone's activity in an online network can be public by default, accessible to selected friends, or to the wider set of friends of friends (who may be strangers to the focal user). These differences have an impact on behaviour, and on the collective dynamics that such behaviour can trigger. As a consequence, the standard terminology of 'influence', 'virality', 'early adopter', etc. hides in fact a significant variance in how cascades are operationalized across platforms. And this diversity puts some constraints on the comparability and generalizability of results. Secondly, there are a number of technical issues (such as the sampling that application programming interfaces [APIs] impose to data collection) that might hamper the validity of some conclusions. For instance, we ignore whether the data retrieved through publicly available APIs is a random sample of all generated activity, or how significant the bias can be [75]. Also, we do not know yet how dynamical classes [76] and the plurality of collective attention patterns [77] relate to observed activity in SNSs-even in the ideal scenario of data from the same SNS and an agreed conceptualization of a cascade. As an illustration, Fig. 3 shows how communication around different topics have evolved for over a year. Clearly most topics present bursty activity ('15 m', 'Elections', 'Reform', 'Strike', 'Sinde') due to events in the real world (black arrows in the Figure indicate such exogenous factors); whereas 'crisis' presents a chatter-like pattern [64], with users continuously discussing at moderate levels and lacking outstanding spikes. Most likely, the mechanisms governing activity during bursty or chatter-like activity are different, but this hypothesis has not been tackled so far (to the best of our knowledge).

#### 4.1 Cascade definitions

4.1.1 *Content-based cascades.* Most empirical work on cascades revolve around 'content chains': the basic criterion to include a node *i* in a diffusion tree starting at *j* is to guarantee that (i) *i* and *j* became friends at  $t_1$  (the notion of 'friend' changes across online platforms and must be understood broadly here); (ii) *i* received a piece of information from *j*, who had previously sent it out, at time  $t_2$ ; and finally (iii) the node *i* sends out the same piece of information at time  $t_3$ . Note that no strict time restriction exists besides the fact that  $t_1 < t_2 < t_3$ , the emphasis being placed on whether what flows is the *same* content. This is the case for e-mail chain letters [78], URL forwarding [74,79] and re-tweeting [73], fanning in Facebook pages [80] or picture spread in Flickr [81].

4.1.2 *Time-constrained activity cascades.* Online platforms allow users to share contents, but also to spread behaviour. When a user likes a Facebook page, she is sending a signal about the content of the page but she is also setting a behavioural precedent, and makes the 'liking' activity more prevalent amongst her neighbours in the network, even if they end up liking completely different pages. So



FIG. 3. Activity time series for six political and economic topics in Twitter (a topic is in this case all those tweets containing at least one hashtag from an arbitrary, predefined closed list). All but one topic ('crisis') exhibit spiky behaviour due to key dates in which exogenous, real-world events triggered activity—demonstrations, election day, etc. These dates are highlighted with black arrows. The topic 'crisis', instead, shows a rather constant—chatter-like [64]—pattern, because it is a rather daily topic—since 2008 at least. How these different trends affect contagion mechanisms—and therefore research approaches to cascading events—is not clear.

focusing exclusively on content to define cascades excludes other interesting diffusion events that also take place in online networks. They include, for instance, the *conversational* [82] or *collaborative* dimensions of SNSs, which can connect groups of people in critical situations [83]. The ability to address other users (like the @mention feature in Twitter, for instance) accentuates these alternative features [73,82], and observed patterns of link reciprocity [84] hint at the use of some SNSs as instant messaging systems, in which non-identical pieces of information around a topic may be circulating (typically over short time spans) in many-to-many interactions, along direct or indirect information pathways [85].

The definition of a time-constrained cascade is useful to measure this type of diffusion, less focused on content and more on behaviour: it uncovers how—and how often—users get involved in sequential message exchange, for which the strict repetition of the same content is not necessary (possibly not even frequent). As in the content-based definition of cascades, time-constrained cascades also assume that conditions (i) to (iii) above are met, except that, for *i* to be included in an avalanche started at *j*, the piece of information being transmitted does not need to be the same, and  $t_3 - t_2 \leq \Delta \tau$ , where  $\tau$  is an arbitrary (typically up to one day) time lapse. In this way, two aspects of critical phenomena (bursty behaviour and avalanches) meet, through the concept of time-resolute cascades.

It is worth noting that content-based and time-constrained cascades do not differ much in their modelling, except in the way they stipulate strong and weak conditions: the former strongly accentuates the strict-content copy condition, with loose temporal constraints (though these exist); whereas the latter lays a tight temporal condition, relaxing the content constraint (though, again, content still matters). This has been the approach in [57,86,87], where content similarity was guaranteed by the limitation of activity to a closed list of hashtags (which referred to specific topics) on Twitter.

4.1.3 Other remarks. Some theoretical approaches (see previous section) have addressed the importance of considering the size of the seed set (typically expressed as  $\rho_0$ , the initial density of cascade originators). When analysing the empirics of cascades, different strategies have been employed regarding seeds. In all cases, the seed is defined as an independent originator of cascades (none in the personal network of a seed has shown any previous activity); but often cascades originate in different regions of the topology (such as two non-connected users liking the same Facebook page, or two unrelated bloggers posting links to the same source). In those cases researchers allow for a multiple-source scheme, in which cascades can merge [80] (or collide [47]). Finally, all these definitions do not control for exogenous factors, and the impact they can have in seeding the network in parallel to, or reinforcing, cascade activation. For instance, the decision to like a Facebook page might result from exposure to friends doing so before, or might in fact be a consequence of an offline association with the person or organization publishing that page. Several strategies can be applied to diminish the effects of this noise (for instance, [80] can be taken as an hybrid cascade definition because they impose both content and temporal conditions, with  $\tau = 24$  h). However, it is fair to note that these definitions of cascades (informational or behavioural) most probably overestimate cause and effect in the endogenous emergence of avalanches.

## 4.2 Characterizing cascades

4.2.1 *Global structure of cascades.* Figure 4 shows the nodes contained in 15 cascades. The examples have been chosen for the sake of visualization (larger as well as smaller cascades do exist). From the point of view of a user–user network, a cascade can be represented as a connected tree-like sub-graph, where the inclusion of a node is driven by activity dynamics. An obvious first question is what these sub-graphs look like. Almost all works addressing this question coincide in the report that most cascades have the shape of ultra-shallow, typically star-graphs [47,57,73,74,79,80,86–88]. The immediate conclusion is that most events do not spread at all, and large-scale cascades are uncommon in the dynamics of social networks. The exception to these robust findings is Liben-Nowell & Kleinberg [78], who report on narrow, deep propagation trees in e-mail letter-forwarding. This may result from the specificity of e-mail communication or, most probably, from fundamental differences in the methodology employed: while cascading behaviour is typically analysed using all initiated cascades, Liben-Nowell and Kleinberg restrict their analysis to the chain letter of a widely circulated petition known to have spread widely. The general trend is that cascade size distributions are stretched, typically under the form of a power law (see the top panel in Fig. 5) [47,57,74] (to cite just a few), but also of lognormal distributions [79].

4.2.2 Temporal and topological penetration. The histogram of most frequent cascade sub-graphs and the size distributions already suggest how far and how long the diffusion of information or behaviour can typically travel through a network: one direct implication is that most initiated cascades die quickly and convey information to near by locations. To get a better idea of how widely and how quickly information propagates, however, additional measurements are needed. Two of these measures are topological penetration,  $\Delta r$ , which can be defined as the shortest path between the seed of the cascade and the farthest node involved in the cascade; and temporal penetration, which can be understood as the lifetime  $\Delta t$  of a cascade. See Fig. 6 for an illustration of both quantities.

The review of the literature suggests that conclusions regarding topological penetration converge, and can thus be taken as robust. Typical social networks—like many other complex networked systems—exhibit a diameter  $D \ll N$  [18]. As such, as soon as cascades grow even to short distances



FIG. 4. Different structures of real cascades occurred on Twitter. Node sizes are proportional to the node degree and links represent the follower relationship between users.

 $(\geq 4$  hops away from the seed), the number of activated nodes escalates very fast; see [73,81,86]. Regarding time, a well-established observation is that interest (for a certain topic, hashtag, etc.) decays very fast –or, conversely, reactions occur mostly soon after the information appeared [22,47]. It is presumably in this narrow time window that large cascades happen, although this fact has been largely overlooked in the literature. On the other hand, for several events there is not a clear pattern of adoption over time [77,80], so temporal penetration may present a rich distribution with a few cascades—the largest ones—lasting for months [86].

#### 4.3 Influence: super-spreaders or hidden influentials?

As Section 3 showed, the concept of influence has been widely discussed in theoretical models, without succeeding to agree on a way to quantify it. The question remains whether there is a set of privileged (presumably topologically salient) nodes that are in a better position to trigger large cascades. A first remarkable finding is that influence should not be simply mapped to connectivity [89]; authority is gained through specialization and concerted efforts to limit communication activity to a single topic. Sun *et al.* [80], Bakshy *et al.* [74] and Kwak *et al.* [73] put to test the 'million follower fallacy' measuring a number of descriptors, possible candidates to grasp influence on a network: number of followers ( $k_{out}$ ), number of friends ( $k_{in}$ ), Pagerank, activity rates (mentions, retweets), and other related measures, and—more or less explicitly—agree on the fact that there is not a clear-cut measure for influence.



FIG. 5. Top: distribution of cascade sizes ( $n_c$ ) suggests that only a few cascades percolate to affect most users, and that the vast majority die in the early stages of diffusion. This result is robust across SNSs and for different  $\Delta \tau$  –in the case of time-constrained cascades. It is also robust to different activity regimes–low, non-spiky period or bursty ones. Middle and bottom: we can observe a positive correlation between normalized cascade sizes  $n_c/N$  and network connectivity—middle—and centrality (measured by the classification of nodes in *k*-cores)—bottom—respectively, suggesting that well-connected users suffice to release global chains of information diffusion. Adapted from [87].

It seems then that large connectivity—being a hub—might be a sufficient, though not necessary, condition for a cascade to occur. Indeed, high connectivity sometimes guarantees the occurrence of large-scale cascades, for instance in the shape of a shallow, wide tree (for example, as soon as a celebrity shows some activity acting as an initiator). A positive correlation between degree (and k-core) and final cascade size confirms this [57,87]; see the middle and lower panels in Fig. 5. Nevertheless, as predicted in [56], it is possible to observe a counterintuitive 'hub-firewall' effect by which cascades may die out when they encounter a hub [86]. A simple Twitter follower vs. friend scatter plot (Fig. 7) provides some keys to this dual behaviour: news media and celebrities' accounts have a disproportionate number of followers relative to those who they follow, such that they behave like sinks (for other's information) and successful sources (when they generate information) due to the striking concentration of attention on these outstanding users [90,91].

All this evidence suggests that, if super-spreaders do not exist, another class of users might be feeding system-wide cascades, users that could go under the label of 'hidden influentials'. Some theoretical



FIG. 6. Structure of a real cascade of duration  $\Delta t = 6\tau$  and maximum topological penetration of  $\Delta r = 4$ . The initial seed (white node) emitted a message at time  $t_0$  that was spread over a subgraph of the network and reached 966 different users. Colours indicate the instant when nodes first listened to the message (left), or their distance (shortest path length) to the initial seed (right).



FIG. 7. Twitter follower vs. friend scatter plot, for two topic-constrained datasets. The observed trends suggest that the huge amount of attention some users receive—those with many more followers than friends—causes these nodes to display a dual behaviour, i.e. as effective spreaders, but also as firewalls.

models coincide in this aspect [58], as mentioned before, and different empirical studies [74,86,92] provide consistent evidence that such a category of users play a crucial role in diffusion dynamics. More specifically, in [86], the authors establish how these nodes, who do not occupy key topological positions ( $10^2 < k_{out} < 10^3$ , in a network with  $k_{max} \approx 35000$ ), cause a large multiplicative effect that results on a high spreading efficiency.

Summing up, theoretical models devise two possibilities: either there is a small subset of special individuals who, given their centrality in the network, can influence a disproportionate number of others; or influence accumulates through the smaller networks of a critical mass of less central people who, on the aggregate, will generate large cascades. What the data reveal is that both views are compatible.

### 4.4 Topological barriers: community structure in social networks

Community structure is a typical feature of social networks, which has also been observed in the online context—the blogosphere as well as in SNSs [93–95]. Originally developed by social analysts [96], there are many available formalizations of the idea of communities and methods to identify the subgroups of individuals within a network [97]. The interest in modular structure lies on the idea that topologically dense clusters impose restrictions to dynamical processes, which has been proved correct in a wide range of phenomena including information transfer [98]. This is also the case in the threshold model [44], as pointed out earlier in this work.

The question remains whether this is the case in actual information cascades. Baños *et al.* [86] first apply a modularity optimization algorithm [97,99,100] to detect the modular structure of the follower-friend Twitter graph. From the resulting partition, a two-level analysis is performed. At the module level, they measure how often a cascade spills over the community where it was triggered. Interestingly, small- to medium-sized cascades (compared with the system size N) mainly stay within their original community, which hints at the fact that influence occurs within specialized topics [74,89]—assuming that people with similar interests tend to gather [101,102]. The strong tendency, for a large fraction of cascades (the smaller ones), to stay within modular boundaries confirms that topological bottlenecks play an important role to hinder large-scale events. At the individual level, two main trends are observed: first, local leaders—nodes with larger-than-average intra-modular connectivity—have a higher probability to trigger large cascades; and secondly, connector nodes—those who link users in other communities besides their own—also have better chances to spread information widely. Note that connector nodes may or may not exhibit large connectivity, which—again, and from a very different level of analysis—strongly suggests that influence may be found in nodes which are not outstanding when classified with typical descriptors.

#### 5. Discussion and future work

Online networks are core to many of the daily activities in which we are involved: some, like gossiping through SNSs, are more mundane than others—for instance, using those networks to access political news that would otherwise be censored by a repressive State. But whatever their use, the one thing that online networks help create is a better view of the connectedness of our actions—of the things we read and do—and the explosive consequences that such interdependence is capable of generating. Complex systems are all about the unpredictable consequences that small changes may generate on a global level; but when those systems are social, and are formed by actors capable of building their own representations of the networks they inhabit, complexity gets another twist and extends into a

whole new level of feedback reactions. The study of social influence, diffusion and contagion (terms that have all been used somehow interchangeably in the course of this article, but that hide nuanced differences in the mechanisms involved) tries to assess how interactions shape decisions and behaviour. Exposure to information or previous behaviour shifts perceptions and attitudes, and propels people to behave in a way that differs from what they would have done in isolation. Cascades are one way of approximating the dynamics of this interdependence. The models designed for their study aim to disentangle the network mechanisms that brought them into existence, and delineate their impact and consequences for the system in which they take place. This review has laid down the basic theoretical tools developed in recent years to attain those aims, and it has assessed those tools in view of the empirical work that online data are facilitating.

Although many of those findings (like the distribution of cascade size) have been replicated across methodologies and datasets-and are therefore robust and consistent-there are still many unknowns that encourage further developments in this area of work. Here, we will outline three. First, more work is required to illuminate the spill-over effects that online contagion dynamics have in the offline world. Epidemics in online games like World of Warcraft can be tragic for the players involved, but nobody would question that flesh-and-bone epidemics are more consequential. Luckily, there is no mechanism that can transfer the spread of a virus from one world to the other, but in many other areas of human behaviour, what happens online has a direct impact on offline actions. Those working in marketing are obviously interested in translating online buzz into higher sales, and some researchers have seriously considered the ways in which online networks might be capitalized for that purpose [103]. Other researchers have actually pushed the boundaries of what is possible with current data by linking the influence dynamics in an online network (Facebook) around self-declared voting behaviour with actual voting records, and finding positive spill-over effects [69]. Identifying these effects is not an easy task, if only because of legitimate privacy concerns; but it is crucial if this line of work is to make an impact not only on our understanding of social systems, but on the way in which we devise interventions, hopefully to promote the public good, like increased civic participation.

A second area of work that requires further developments has to do with disentangling the joint effects of local versus global information in adoption rates. Networks, as this review has explained, channel social influence by exposing individuals to the behaviour of their contacts or friends. To the extent that every actor inhabits a different local context, influence will flow differently in different parts of the network. However, these streams of information often coexist, and interact with, the effects of common exposure to a single global source. This might take the form of mass media, exogenous to the network or—depending on the affordances of the platform—some metric that summarizes global activity, like trending topics in Twitter. The interaction between local and global influence in shaping adoption rates has been considered before [104] but not in the context of complex networks. Related to this, current efforts to model multiplex and time-varying networks might feed into this goal of taking into account the several layers through which influence spreads [105].

A third area that can benefit our understanding of contagion behaviour refers to its microfoundations, that is, to the psychological or cognitive triggers that make people want to join a cascade. There are exciting developments in experimental psychology [106] and neuroscience [107] that aim to pin down the mechanisms of information propagation. What this approach suggests is that emotions or sentiment play a significant role in predicting the behaviour that allows content to go viral or at least sets its preconditions. This, in turn, points to another fascinating and related area of work in machine learning and NLP that aims to quantify the subjectivity of human communication, with a special focus on social media [108]. The metrics that come out of these classifiers can be used to explain why some content might generate larger cascades: if emotions are triggers of behaviour, and messages offer stimuli 18 of 22

capable of arousing certain emotions, this creates a connection to the type of mechanisms that cognitive scientists are exploring at the brain and behavioural level.

Digital technologies, and their increasing prevalence in every dimension of social life, promise to yield the data that can help advance research on those three fronts—and thus enhance our understanding of contagious behaviour. This is important not only for the sake of scientific satisfaction but also for the implications that such knowledge can have on improving governance and public good interventions. ICTs have encouraged the emergence of a new form of organization that defies the hierarchical nature of bureaucracies by harnessing the power of interdependent decision-making. Examples include crowd-sourcing creative projects; platforms that help improve local governance; and prize-backed competitions that decentralize policy-making. These initiatives rely on the mobilizing power of networks, and the chain reactions that influence and contagion can produce. The potential of decentralized networks as a mechanism for decision-making can transform the way in which governments work and citizens self-organize—but a better understanding of the complex mechanisms that govern those networks is first necessary. This review has given an overview of how much has been learned so far, and outlined where to go from here.

### Funding

R.A.B. was supported by the FPI program of the Government of Aragón, Spain. This work has been partially supported by MINECO through Grants FIS2011-25167 and FIS2012-35719; Comunidad de Aragón (Spain) through a grant to the group FENOL; and by the EC FET-Proactive Projects PLEX-MATH (grant 317614) and MULTIPLEX (grant 317532).

#### References

- 1. http://investor.fb.com/releasedetail.cfm?ReleaseID=736911.
- 2. https://business.twitter.com/whos-twitter.
- 3. GUNDOTRA, V. (2012) https://plus.google.com/+VicGundotra/posts/2YWhK1K3FA5.
- **4.** (2012) Measuring the Information Society. *Technical Report*. International Telecommunication Union. http://www.itu.int/pub/D-IND-ICTOI-2012.
- 5. KUHN, T. S. (1962) The Structure of Scientific Revolutions. Chicago: University of Chicago press.
- LAZER, D., PENTLAND, A. S., ADAMIC, L., ARAL, S., BARABÁSI, A. L., BREWER, D., CHRISTAKIS, N., CONTRACTOR, N., FOWLER, J., GUTMANN, M., JEBARA, T., KING, G., MACY, M., ROY, D. & VAN ALSTYNE, M. (2009) Life in the network: the coming age of computational social science. *Science*, 323, 721.
- 7. CONTE, R., GILBERT, N., CIOFFI-REVILLA, C., DEFFUANT, G., KERTESZ, J., LORETO, V., MOAT, S., NADAL, J.-P., SANCHEZ, A., NOWAK, A., FLACHE, A., SAN MIGUEL, M., & HELBING, D. (2012) Manifesto of computational social science. *Eur. Phys. J. Special Topics*, **214**, 325–346.
- 8. GILES, J. (2012) Computational social science: making the links. *Nature*, 488, 448–450.
- 9. MILGRAM, S. (1977) The Individual in a Social World: Essays and Experiments. London: Pinter & Martin.
- 10. SCHELLING, T. C. (1978) Micromotives and Macrobehavior. London: Norton.
- 11. WATTS, D. J. (2004) The 'new' science of networks. Annu. Rev. Sociol., 30, 243-270.
- **12.** COLEMAN, J. S. (1990) *Foundations of Social Theory*. Cambridge, MA: Belknap Press of Harvard University Press.
- 13. LIU, Y.-Y., SLOTINE, J.-J. & BARABÁSI, A.-L. (2011) Controllability of complex networks. *Nature*, 473, 167–173.
- **14.** VESPIGNANI, A. (2011) Modelling dynamical processes in complex socio-technical systems. *Nat. Phys.*, **8**, 32–39.
- 15. KLEMM, K., SERRANO, M., EGUILUZ, V. & MIGUEL, M. (2012) A measure of individual role in collective dynamics: spreading at criticality. *Sci. Rep.*, **2**, 292.

- **16.** NEWMAN, M., BARABÁSI, A. & WATTS, D. (2006) *The Structure and Dynamics of Networks*. Princeton, NJ: Princeton University Press.
- 17. ESTRADA, E. (2011) The Structure of Complex Networks. Oxford: Oxford University Press.
- 18. BOCCALETTI, S., LATORA, V., MORENO, Y., CHAVEZ, M. & HWANG, D. (2006) Complex networks: structure and dynamics. *Phys. Rep.*, 424, 175–308.
- 19. NEWMAN, M. (2010) Networks: An Introduction. Oxford: Oxford University Press.
- 20. ERDÖS, P. & RÉNYI, A. (1959) On random graphs. Publ. Math. (Debrecen), 6, 290–297.
- 21. BARABÁSI, A. & ALBERT, R. (1999) Emergence of scaling in random networks. Science, 286, 509.
- 22. CENTOLA, D. (2010) The spread of behavior in an online social network experiment. *Science*, 329, 1194–1197.
- 23. BONACICH, P. (1972) Factoring and weighting approaches to status scores and clique identification. J. Math. Sociol., 2, 113–120.
- 24. PAGE, L., BRIN, S., MOTWANI, R. & WINOGRAD, T. (1998) The PageRank citation ranking: Bringing order to the web. *Technical Report*. Standford, CA: Stanford Digital Library Technologies Project.
- 25. ROGERS, E. M. (2003) Diffusion of Innovations. New York, NY: Free Press.
- 26. DIMAGGIO, P. & GARIP, F. (2012) Network effects and social inequality. Annu. Rev. Sociol., 38, 93–118.
- 27. COLEMAN, J., KATZ, E. & MENZEL, H. (1957) The diffusion of an innovation among physicians. *Sociometry*, 20, 253–270.
- 28. VALENTE, T. W. (2012) Network interventions. Science, 337, 49–53.
- **29.** BIGGS, M. (2005) Strikes as forest fires: Chicago and Paris in the late nineteenth century1. *Am. J. Sociol.*, **110**, 1684–1714.
- HEDSTRÖM, P. (1994) Contagious collectivities: on the spatial diffusion of Swedish trade unions, 1890–1940. Am. J. Sociol., 99, 1157–1179.
- 31. GRANOVETTER, M. (1978) Threshold models of collective behavior. Am. J. Sociol., 83, 1420–1443.
- **32.** WATTS, D. (2002) A simple model of global cascades on random networks. *Proc. Natl Acad. Sci.*, **99**, 5766–5771.
- 33. CENTOLA, D., EGUÍLUZ, V. & MACY, M. (2007) Cascade dynamics of complex propagation. *Phys. A: Stat. Mech. Appl.*, 374, 449–456.
- **34.** BAILEY, N. T. (1975) *The Mathematical Theory of Infectious Diseases and its Applications*, 2nd edn. High Wycombe: Charles Griffin & Company Ltd.
- 35. MURRAY, J. (1993) Mathematical Biology. Berlin: Springer.
- 36. HETHCOTE, H. W. (2000) The mathematics of infectious diseases. SIAM Rev., 42, 599-653.
- **37.** RAPOPORT, A. (1953) Spread of information through a population with socio-structural bias I. Assumption of transitivity. *Bull. Math. Biophys.*, **15**, 523–533.
- GOFFMAN, W. & NEWILL, V. A. (1964) Generalization of epidemic theory. An application to the transmission of ideas. *Nature*, 204, 225–228.
- 39. DALEY, D. J. & KENDALL, D. G. (1964) Epidemics and rumours. Nature, 204, 1118.
- 40. DODDS, P. S. & WATTS, D. J. (2004) Universal behavior in a generalized model of contagion. *Phys. Rev. Lett.*, 92, 218701.
- **41.** SCHELLING, T. C. (1973) Hockey helmets, concealed weapons, and daylight saving: a study of binary choices with externalities. *J. Conflict Resolution*, **17**, 381–428.
- **42.** VALENTE, T. W. (1995) *Network Models of the Diffusion of Innovations*. Quantitative Methods in Communication Series. Cresskill, NJ: Hampton Press.
- **43.** GLEESON, J. & CAHALANE, D. (2007) Seed size strongly affects cascades on random networks. *Phys. Rev. E*, **75**, 056103.
- 44. GLEESON, J. (2008) Cascades on correlated and modular random networks. Phys. Rev. E, 77, 046117.
- **45.** HACKETT, A., MELNIK, S. & GLEESON, J. (2011) Cascades on a class of clustered random networks. *Phys. Rev. E*, **83**, 056107.
- 46. PASTOR-SATORRAS, R. & VESPIGNANI, A. (2001) Epidemic spreading in scale-free networks. *Phys. Rev. Lett.*, 86, 3200–3203.

- 47. LESKOVEC, J., MCGLOHON, M., FALOUTSOS CH, G. & HURST, M. (2007) Cascading behavior in large blog graphs. *Proceedings of the 7th SIAM International Conference on Data Mining (SDM)*, Philadelphia, USA, pp. 29406–29413.
- **48.** COLEMAN, J. S., KATZ, E. & MENZEL, H. (1966) *Medical Innovation: A Diffusion Study*. Indianapolis: Bobbs-Merrill.
- **49.** KATZ, E. & LAZARSFELD, P. (1955) *Personal Influence: The Part Played by People in the Flow of Communications.* Glencoe, IL: Free Press.
- 50. CENTOLA, D. & MACY, M. (2007) Complex contagions and the weakness of long ties. Am. J. Sociol., 113, 702–734.
- ROMERO, D. M., MEEDER, B. & KLEINBERG, J. (2011) Differences in the mechanics of information diffusion across topics: idioms, political hashtags, and complex contagion on twitter. *Proceedings of the 20th International Conference on World Wide Web*. New York: ACM, pp. 695–704.
- **52.** ARAL, S. & WALKER, D. (2011) Creating social contagion through viral product design: a randomized trial of peer influence in networks. *Manag. Sci.*, **57.9**, 1623–1639.
- 53. KITSAK, M., GALLOS, L., HAVLIN, S., LILJEROS, F., MUCHNIK, L., STANLEY, H. & MAKSE, H. (2010) Identification of influential spreaders in complex networks. *Nat. Phys.*, **6**, 888–893.
- 54. CASTELLANO, C., FORTUNATO, S. & LORETO, V. (2009) Statistical physics of social dynamics. *Rev. Modern Phys.*, **81**, 591.
- 55. MORENO, Y., NEKOVEE, M. & PACHECO, A. (2004) Dynamics of rumor spreading in complex networks. *Phys. Rev. E*, **69**, 066130.
- 56. BORGE-HOLTHOEFER, J. & MORENO, Y. (2012) Absence of influential spreaders in rumor dynamics. *Phys. Rev. E*, **85**, 026116.
- 57. GONZÁLEZ-BAILÓN, S., BORGE-HOLTHOEFER, J., RIVERO, A. & MORENO, Y. (2011) The dynamics of protest recruitment through an online network. *Sci. Rep.*, **1**, 197.
- 58. WATTS, D. J. & DODDS, P. S. (2007) Influentials, networks, and public opinion formation. J. Consumer Res., 34, 441.
- 59. BARABÁSI, A. (2005) The origin of bursts and heavy tails in human dynamics. Nature, 435, 1251–1251.
- **60.** BORGE-HOLTHOEFER, J., MELONI, S., GONÇALVES, B. & MORENO, Y. (2012) Emergence of influential spreaders in modified rumor models. *J. Stat. Phys.*, **148**, 1–11.
- **61.** GOLDENBERG, J., LIBAI, B. & MULLER, E. (2001) Talk of the network: a complex systems look at the underlying process of word-of-mouth. *Mark. Lett.*, **12**, 211–223.
- 62. GUARDIOLA, X., DIAZ-GUILERA, A., PEREZ, C. J., ARENAS, A. & LLAS, M. (2002) Modeling diffusion of innovations in a social network. *Phys. Rev. E*, 66, 026121.
- **63.** KEMPE, D., KLEINBERG, J. & TARDOS, É. (2003) Maximizing the spread of influence through a social network. *Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York: ACM, pp. 137–146.
- **64.** GRUHL, D., GUHA, R., LIBEN-NOWELL, D. & TOMKINS, A. (2004) Information diffusion through blogspace. *Proceedings of the 13th International Conference on World Wide Web.* New York: ACM, pp. 491–501.
- 65. FOWLER, J. (2005) Turnout in a Small World. Philadelphia: Temple University Press, pp. 269–287.
- 66. ONNELA, J.-P., SARAMÄKI, J., HYVÖNEN, J., SZABÓ, G., DE MENEZES, M. A., KASKI, K., BARABÁSI, A.-L. & KERTÉSZ, J. (2007) Analysis of a large-scale weighted network of one-to-one human communication. *New J. Phys.*, 9, 179.
- 67. WANG, P., GONZÁLEZ, M. C., HIDALGO, C. A. & BARABÁSI, A.-L. (2009) Understanding the spreading patterns of mobile phone viruses. *Science*, **324**, 1071–1076.
- **68.** MIRITELLO, G., MORO, E. & LARA, R. (2011) Dynamical strength of social ties in information spreading. *Phys. Rev. E*, **83**, 045102.
- **69.** BOND, R. M., FARISS, C. J., JONES, J. J., KRAMER, A. D. I., MARLOW, C., SETTLE, J. E. & FOWLER, J. H. (2012) A 61-million-person experiment in social influence and political mobilization. *Nature*, **489**, 295–8.
- **70.** Kossinets, G. & Watts, D. J. (2006) Empirical analysis of an evolving social network. *Science*, **311**, 88–90.

- 71. Kossinets, G. & Watts, D. J. (2009) Origins of homophily in an evolving social network1. *Am. J. Sociol.*, 115, 405–450.
- 72. ARAL, S. & VAN ALSTYNE, M. (2011) The diversity-bandwidth trade-off. Am. J. Sociol., 117, 90–171.
- **73.** KWAK, H., LEE, C., PARK, H. & MOON, S. (2010) What is Twitter, a social network or a news media? *Proceedings of the 19th International Conference on World Wide Web*. New York: ACM, pp. 591–600.
- 74. BAKSHY, E., HOFMAN, J., MASON, W. & WATTS, D. (2011) Everyone's an influencer: quantifying influence on twitter. *Proceedings of the 4th ACM International Conference on Web Search and Data Mining*. New York: ACM, pp. 65–74.
- **75.** GONZÁLEZ-BAILÓN, S., WANG, N., RIVERO, A., BORGE-HOLTHOEFER, J. & MORENO, Y. (2012) Assessing the bias in communication networks sampled from twitter. arXiv:1212.1684.
- **76.** CRANE, R. & SORNETTE, D. (2008) Robust dynamic classes revealed by measuring the response function of a social system. *Proc. Natl Acad. Sci.*, **105**, 15649–15653.
- LEHMANN, J., GONÇALVES, B., RAMASCO, J. J. & CATTUTO, C. (2012) Dynamical classes of collective attention in twitter. *Proceedings of the 21st International Conference on World Wide Web*. New York: ACM, pp. 251–260.
- **78.** LIBEN-NOWELL, D. & KLEINBERG, J. (2008) Tracing information flow on a global scale using Internet chain-letter data. *Proc. Natl Acad. Sci.*, **105**, 4633–4638.
- **79.** LERMAN, K., GHOSH, R. & SURACHAWALA, T. (2012) Social contagion: an empirical study of information spread on Digg and Twitter follower graphs. arXiv:1202.3162.
- SUN, E., ROSENN, I., MARLOW, C. & LENTO, T. (2009) Gesundheit! modeling contagion through facebook news feed. Proc. 3rd Internat. Conf. Weblogs Soc. Media (ICWSM), 22–30.
- **81.** CHA, M., MISLOVE, A. & GUMMADI, K. (2009) A measurement-driven analysis of information propagation in the Flickr social network. *Proceedings of the 18th International Conference on World Wide Web*. New York: ACM, pp. 721–730.
- HONEY, C. & HERRING, S. C. (2009) Beyond microblogging: conversation and collaboration via twitter. 42nd Hawaii International Conference on System Sciences, HICSS'09. Washington, DC: IEEE, pp. 1–10.
- 83. MUNGIU-PIPPIDI, A. & MUNTEANU, I. (2009) Moldova's 'Twitter Revolution'. J. Democracy, 20, 136–142.
- **84.** GARLASCHELLI, D. & LOFFREDO, M. I. (2004) Patterns of link reciprocity in directed networks. *Phys. Rev. Lett.*, **93**, 268701.
- **85.** KOSSINETS, G., KLEINBERG, J. & WATTS, D. (2008) The structure of information pathways in a social communication network. *Proceeding of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York: ACM, pp. 435–443.
- **86.** BAÑOS, R., BORGE-HOLTHOEFER, J. & MORENO, Y. (2013) The role of hidden influentials in the diffusion of online information cascades. arXiv:1303.4629.
- 87. BORGE-HOLTHOEFER, J., RIVERO, A. & MORENO, Y. (2012) Locating privileged spreaders on an online social network. *Phys. Rev. E*, 85, 066123.
- 88. WANG, D., WEN, Z., TONG, H., LIN, C.-Y., SONG, C. & BARABÁSI, A. L. (2011) Information spreading in context. *Proceedings of the 20th International Conference on World Wide Web*. New York: ACM, pp. 735–744.
- **89.** CHA, M., HADDADI, H., BENEVENUTO, F. & GUMMADI, K. (2010) Measuring user influence in twitter: The million follower fallacy. *Proceedings of the 4th International Conference on Weblogs and Social Media (ICWSM)*, Palo Alto, CA, USA, pp. 10–17.
- 90. BORGE-HOLTHOEFER, J., RIVERO, A., GARCÍA, I., CAUHÉ, E., FERRER, A., FERRER, D., FRANCOS, D., IÑIGUEZ, D., PÉREZ, M., RUIZ, G., SANZ, F., SERRANO, F., VIÑAS, C., TARANCÓN, A. & MORENO Y. (2011) Structural and dynamical patterns on online social networks: the Spanish May 15th movement as a case study. *PloS One*, 6, e23883.
- 91. WU, S., HOFMAN, J., MASON, W. & WATTS, D. (2011) Who says what to whom on twitter. Proceedings of the 20th International Conference on World Wide Web. New York, NY, USA: ACM, pp. 705–714.

- 92. GONZÁLEZ-BAILÓN, S., BORGE-HOLTHOEFER, J. & MORENO, Y. (2013) Broadcasters and Hidden Influentials in Online Protest Diffusion. Am. Behav. Sci. doi:10.1177/0002764213479371.
- **93.** ADAMIC, L. & GLANCE, N. (2005) The political blogosphere and the 2004 US election: divided they blog. *Proceedings of the 3rd International Workshop on Link Discovery*. New York: ACM, pp. 36–43.
- 94. CONOVER, M., RATKIEWICZ, J., FRANCISCO, M., GONÇALVES, B., FLAMMINI, A. & MENCZER, F. (2011) Political polarization on twitter. *Proceedings of the 5th International Conference on Weblogs and Social Media (ICWSM)*, Palo Alto, CA, USA, pp. 89–96.
- **95.** GRABOWICZ, P. A., RAMASCO, J. J., MORO, E., PUJOL, J. M. & EGUILUZ, V. M. (2012) Social features of online networks: the strength of intermediary ties in online social media. *PloS One*, **7**, e29358.
- **96.** WASSERMAN, S. & FAUST, K. (1994) *Social Network Analysis: Methods and Applications*. Cambridge: Cambridge University Press.
- 97. FORTUNATO, S. (2010) Community detection in graphs. Phys. Rep., 486, 75-174.
- **98.** DANON, L., ARENAS, A. & DÍAZ-GUILERA, A. (2008) Impact of community structure on information transfer. *Phys. Rev. E*, **77**, 36103.
- **99.** NEWMAN, M. E. J. & GIRVAN, M. (2004) Finding and evaluating community structure in networks. *Phys. Rev. E*, **69**.
- DANON, L., DUCH, J., ARENAS, A. & DÍAZ-GUILERA, A. (2005) Comparing community structure identification. J. Stat. Mech.: Theory Exp., 2005, P09008.
- 101. MCPHERSON, M., SMITH-LOVIN, L. & COOK, J. M. (2001) Birds of a feather: homophily in social networks. *Annu. Rev. Sociol.*, 27, 415–444.
- **102.** CENTOLA, D., GONZALEZ-AVELLA, J. C., EGUILUZ, V. M. & SAN MIGUEL, M. (2007) Homophily, cultural drift, and the co-evolution of cultural groups. *J. Conflict Resolution*, **51**, 905–929.
- **103.** ARAL, S. & WALKER, D. (2011) Forget viral marketing—make the product itself viral. *Harv. Bus. Rev.*, **89**, 34–35.
- 104. ONNELA, J.-P. & REED-TSOCHAS, F. (2010) Spontaneous emergence of social influence in online systems. *Proc. Natl Acad. Sci.*, 107, 18375–18380.
- 105. GÓMEZ, S., DÍAZ-GUILERA, A., GÓMEZ-GARDEÑES, J., PÉREZ-VICENTE, C., MORENO, Y. & ARENAS, A. (2013) Diffusion dynamics on multiplex networks. *Phys. Rev. Lett.*, 110, 028701.
- 106. BERGER, J. & MILKMAN, K. (2012) What makes online content viral? J. Mark. Res., 49, 192–205.
- 107. FALK, E. B., O'DONNELL, M. B. & LIEBERMAN, M. D. (2012) Getting the word out: neural correlates of enthusiastic message propagation. *Front. Hum. Neurosci.*, **6**. doi:10.3389/fnhum.2012.00313.
- **108.** PALTOGLOU, G. & THELWALL, M. (2011) Twitter, MySpace, Digg: unsupervised sentiment analysis in social media. *ACM Trans. Intell. Syst. Technol.*, **3**, 66.