

Identifying influential spreaders in complex networks

Maksim Kitsak,^{1,2} Lazaros K. Gallos,³ Shlomo Havlin,⁴ Fredrik Liljeros,⁵

Lev Muchnik,⁶ H. Eugene Stanley,¹ and Hernán A. Makse³

*¹Center for Polymer Studies and Physics Department,
Boston University, Boston, Massachusetts 02215, USA*

*²Cooperative Association for Internet Data Analysis (CAIDA),
University of California-San Diego, La Jolla, California 92093, USA*

*³Levich Institute and Physics Department,
City College of New York, New York, New York 10031, USA*

*⁴Minerva Center and Department of Physics,
Bar-Ilan University, Ramat Gan, Israel*

⁵Department of Sociology, Stockholm University, S-10691, Stockholm, Sweden

*⁶Information Operations and Management Sciences Department,
Stern School of Business, New York University, New York, New York 10012, USA*

(Dated: January 29, 2010— kghsm.tex)

Abstract

Networks portray a multitude of interactions through which people meet, ideas are spread, and infectious diseases propagate within a society. Identifying the most efficient “spreaders” in a network is an important step to optimize the use of available resources and ensure the more efficient spread of information. Here we show that, in contrast to common belief, the most influential spreaders in a social network do not correspond to the best connected people or to the most central people (high betweenness centrality). Instead, we find: *(i)* The most efficient spreaders are those located within the core of the network as identified by the k -shell decomposition analysis. *(ii)* When multiple spreaders are considered simultaneously, the distance between them becomes the crucial parameter that determines the extend of the spreading. Furthermore, we find that— in the case of infections that do not confer immunity on recovered individuals— the infection persists in the high k -shell layers of the network under conditions where hubs may not be able to preserve the infection. Our analysis provides a plausible route for an optimal design of efficient dissemination strategies.

Spreading is a ubiquitous process which describes many important activities in society [1–5]. The knowledge of the spreading pathways through the network of social interactions is crucial for developing efficient methods to either hinder spreading in the case of diseases, or accelerate spreading in the case of information dissemination. Indeed, people are connected according to the way they interact with each other in society and the large heterogeneity of the resulting network greatly determines the efficiency and speed of spreading. In the case of networks with a broad degree distribution (number of links per node) [6], it is believed that the most connected people (hubs) are the key players being responsible for the largest scale of the spreading process [6–8]. Furthermore, in the context of social network theory, the importance of a node for spreading is often associated with the betweenness centrality which is believed to determine who has more ‘interpersonal influence’ on others [9, 10].

Here we argue that the topology of the network organization plays an important role such that there are plausible circumstances under which the highly connected nodes or the highest betweenness nodes have little effect in the range of a given spreading process. For example, if a hub exists at the end of a branch at the periphery of a network, it will have a minimal impact in the spreading process through the core of the network, while a less connected person who is strategically placed in the core of the network will have a significant effect that leads to dissemination through a large fraction of the population. In order to identify the core and the periphery of the network we use the k -shell [11–13] decomposition of the network. Examining this quantity in a number of real networks allows us to identify the best individual spreaders in the network when the spreading originates in a single node. For the case of a spreading process originating in many nodes simultaneously we show that we can further improve the efficiency by considering spreading origins located at a determined distance from each other.

We study real-world complex networks that represent archetypical examples of social structures. We investigate (i) the friendship network between 5.5 million members of the *LiveJournal.com* community [14], (ii) the network of email contacts in the Computer Science Department of London’s Global University [15], (iii) the contact network of inpatients (CNI) collected from hospitals in Sweden [16], and (iv) the network of actors who have co-starred in movies labeled by imdb.com as adult [17] (see Section I for details).

To study the spreading process we apply the Susceptible-Infectious-Recovered (SIR) and Susceptible-Infectious-Susceptible (SIS) models [2, 3, 18] on the above networks. In the SIR

model, all nodes are initially in susceptible state (S) except for one node in the infectious state (I). At each time step, the I nodes attempt to infect their susceptible neighbors with probability β and then enter the recovered state (R) where they become immunized and cannot be infected again. The SIS model aims to describe spreading processes that do not confer immunity on recovered individuals: infected individuals return to the susceptible state with probability λ (here we use $\lambda = 0.8$) and can be reinfected at subsequent time steps, while they remain infectious with probability $1 - \lambda$. These models have been used to describe disease spreading as well as information and rumor spreading in social processes where the actor constantly needs to be reminded [19]. In our study we use relatively small β values for the infection probability, so that the infected percentage of the population remains small. In the case of large β values, where spreading reaches the entire population, the role of individual nodes is no longer important and spreading would cover all the network, independently of where it originated from.

Figures 1a-c illustrate the fact that the size of the population infected by a single node in the CNL network is not necessarily related to the degree of node: spreading may be very different even when it starts from hubs of similar degree as comparatively shown in Figs. 1a and b. Instead, we will show that the location of a node is more important to predict the spreading importance. The k -shell decomposition determines the location of a node in the network by assigning a k_S index to each node to distinguish nodes that are in the periphery of the network with low k_S index from those that belong to the innermost network core with high k_S , (see Fig. 1d and Section II for a definition) [11–13]. Indeed, in the examples of Fig. 1a and 1c, nodes with different degree but same k_S index produce similar spreading areas. This example suggests that the position of the node relative to the organization of the network determines its spreading influence more than a local property, such as the degree of the node, k .

Next, we quantify the influence of a given node i in an SIR spreading process by studying the average size of the population M_i infected in an epidemic originating at i . The extent of an epidemic that starts in a typical single node with (k_S, k) is estimated by the infected population averaged over all the origins with the same (k_S, k) values:

$$M(k_S, k) = \sum_{i \in \Upsilon(k_S, k)} \frac{M_i}{N(k_S, k)}, \quad (1)$$

where $\Upsilon(k_S, k)$ is the union of all $N(k_S, k)$ nodes with (k_S, k) values.

The analysis of $M(k_S, k)$ in the studied social networks reveals three general results (see Fig. 2): (a) For a fixed degree, there is a wide spread of $M(k_S, k)$ values. In particular, there are many hubs located in the periphery of the network (large k , low k_S) that are poor spreaders. (b) For a fixed k_S , $M(k_S, k)$ is approximately independent of the degree of the nodes. This result is revealed in the vertically layered structure of $M(k_S, k)$ suggesting that infected nodes located in the same k -shell produce similar epidemic outbreaks $M(k_S, k)$ independent of the value of k of the infection origin. (c) The most efficient spreaders are located in the inner-core of the network (large k_S region) fairly independently of their degree. These results indicate that the k -shell index of a node is better predictor of spreading influence. When an outbreak starts in a large k -shell node there exist many pathways through which a virus can infect the rest of the network. The existence of these pathways also implies that nodes located in high k -shell layers are infected earlier and they are more likely to be infected during an epidemic outbreak (see Section III). Similar results are obtained from the analysis of $M(k_S, C_B)$ in Fig. 2, where C_B is the betweenness centrality of a node in the network [9, 10]: the value of C_B is not a good predictor for spreading efficiency.

To quantify the above claim we calculate the “imprecision functions” $\epsilon_{k_S}(p)$, $\epsilon_k(p)$, and $\epsilon_{C_B}(p)$, for each of the three indicators k_S , k , or C_B , respectively. The function $\epsilon(p)$ quantifies the difference in the average spreading between the pN nodes ($0 < p < 1$) with highest k_S , k , or C_B from the average spreading of the pN most efficient spreaders (N is the number of nodes in the network, see the caption of Fig. 3a for a detailed definition of ϵ and Section IV). The k -shell strategy is found to be consistently more accurate than a method based on k in the studied p range (Fig. 3a). The C_B -based strategy gives poor results compared to the other two strategies, indicating that it is not an appropriate method for predicting the spreading efficiency of a node.

Our finding is not specific to the social networks shown in Fig. 2. In Section V we analyze the spreading efficiency in other networks not social in origin, like the Internet at the router level, with similar conclusions. The key insight of our finding is that in the studied networks a large number of hubs are often not located in the inner core, but instead are located in the peripheral low k -shell layers (Fig. 3b shows the location of the top hubs in the CNI; see also Section V). These hubs contribute poorly to spreading; their existence is a consequence of the rich topological structure of real networks. In contrast, in a randomized network obtained by randomly rewiring a real network preserving the node degree (see Section VI)

all the hubs are placed in the core of the network (see the red scatter plot in Fig. 3c) and they contribute equally largely to spreading. In such a randomized structure there is a monotonic increase $k_S \propto k$ (Fig. 3c) and the same information is contained in the k -shell as in the degree classification. Examples of real networks that are similar to random are the network of product space of economic goods [20] and the Internet at the AS level (which are analyzed in the Section V).

Our study highlights the importance of the relative location of a *single* spreading origin. Next, we address the question of the extent of an epidemic that starts in simultaneous *multiple* origins. Figure 3d shows the extent of SIR spreading in the CNI network when the outbreak simultaneously starts from the n nodes with the highest degree k or the highest k_S index. Surprisingly, even though the high k -shell nodes are the best single spreaders, the nodes with highest degree are more efficient than highest k_S in multiple spreading. This result can be attributed to the overlap of the areas infected by the multiple nodes: the high k_S nodes tend to be clustered close to each other, while hubs are more spread in the network. Clearly, the steps in the plot of high k_S nodes indicate that even if we infect all the nodes in a given shell, there will only be a minor influence to the infected area, while including just one node from a different shell can result in a significantly increased spreading. This suggests that a better spreading strategy using multiple n spreaders is to choose either the highest k or k_S nodes with the requirement that no two of the n spreaders are directly linked to each other. This scheme then provides the largest infected area of the network as shown in Fig. 3d.

Many contagious infections, including most sexually transmitted infections [21], do not confer full immunity after infection as assumed in the SIR model, and therefore are suitably described by the SIS epidemic model. In an SIS epidemic the number of infectious nodes eventually reaches a dynamic equilibrium “endemic” state at which exactly as many infectious individuals become susceptible as susceptible nodes become infectious [18]. The spreading efficiency of a given node i in SIS spreading is the persistence, $\rho_i(t)$, defined as the probability that node i is infected at time t [7]; in an endemic SIS state, $\rho_i(t \rightarrow \infty)$ becomes independent of t (see Section VII). The persistence $\rho_i(t \rightarrow \infty)$ has been shown to be higher in hubs which are reinfected frequently due to the large number of their neighbors [7, 22, 23]. However, we find that this result holds only in randomized network structures. In real network topologies we find that viruses persist mainly in high k_S layers instead.

For instance, in a random network the epidemic threshold for SIS is estimated by $\beta_c \equiv \lambda\langle k \rangle / \langle k^2 \rangle$ [22, 24] so that viruses persist when $\beta > \beta_c$. On the contrary, in the CNI network viruses persist even if $\beta < \beta_c = 1.7\%$ (Fig. 4a-c). In this case, viruses survive only within the high k -shell layers of the network, while virus persistence in lower k -shell layers is negligible. As a k -shell structure presupposes both assortative interaction and high clustering (at least locally in the higher shells), this result is in agreement with the observation that high clustering [25] and high positive assortativity [26] may decrease the epidemic threshold.

This picture is confirmed when we analyze the persistence function

$$\rho(k_S, k) \equiv \sum_{i \in \Upsilon(k_S, k)} \frac{\rho_i(t \rightarrow \infty)}{N(k_S, k)}, \quad (2)$$

as a function of (k_S, k) for different β above and below the random case threshold $\beta_c = 1.7\%$ (Fig. 4a and b, respectively). High k -shell layers in networks might be closely related to the concept core group in STI research [21]. The core groups are defined as distinct subgroups in the general population that are all characterized by high partner turnover rate and extensive intergroup interaction [21]. Similar to the core group, the dense sub-network formed by nodes in the innermost k -shells helps the virus to consistently survive locally in the inner-core area and infect other nodes adjacent to the area. These k -shells preserve the existence of a virus, in contrast to e.g. isolated hubs in the periphery. Note that virus can not survive in the degree-preserving randomized version of the CNI network due to the absence of high k -shells.

In conclusion, ranking the importance of individual nodes in spreading can influence the success of dissemination strategies. When spreading starts from a single node, the k_S value is enough for this ranking, while in the case of many simultaneous origins, spreading is greatly enhanced when we additionally “repel” spreaders. In the case of infections that do not confer immunity on recovered individuals, high k -shell layers form a reservoir where infection can survive locally even if its contagiousness is well below the epidemic threshold.

-
- [1] Caldarelli G., Vespignani A. (eds) *Large scale structure and dynamics of complex networks*. (World Scientific, Singapore, 2007).
- [2] Anderson, R. M., May, R. M., & Anderson, B. *Infectious Diseases of Humans: Dynamics and Control* (Oxford Science Publications, 1992).
- [3] Diekmann, O., & Heesterbeek, J. A. P. *Mathematical Epidemiology of Infectious Diseases: Model Building, Analysis and Interpretation* (Wiley Series in Mathematical & Computational Biology, New York, 2000).
- [4] Keeling, M. J., & Rohani, P. *Modeling Infectious Diseases in Humans and Animals* (Princeton University Press, 2008).
- [5] Rogers, E. M. *Diffusion of Innovation* (Free Press, New York, 4th ed, 1995).
- [6] Albert, R., Jeong, H., & Barabasi, A.-L. Error and attack tolerance of complex networks. *Nature* **406**, 378–482 (2000).
- [7] Pastor-Satorras, R. & Vespignani, A. Epidemic spreading in scale-free networks. *Phys. Rev. Lett.* **86**, 3200–3203 (2001).
- [8] Cohen, R., Erez, K., ben-Avraham, D. & Havlin, S. Breakdown of the Internet under intentional attack. *Phys. Rev. Lett.* **86**, 3682–3685 (2001).
- [9] Freeman, L. C. Centrality in social networks: Conceptual clarification. *Social Networks* **1**, 215–239 (1979).
- [10] Friedkin, N.E. Theoretical foundations for centrality measures. *Am. J. of Sociology* **96**, 1478–1504 (1991).
- [11] Bollobas, B. *Graph Theory and Combinatorics: Proceedings of the Cambridge Combinatorial Conference in honor of P. Erdős*, 35 (Academic, New York, 1984).
- [12] Seidman, S. B. Network structure and minimum degree. *Social Networks* **5**, 269–287 (1983).
- [13] Carmi, S., Havlin, S, Kirkpatrick, S., Shavitt, Y. & Shir, E. A model of Internet topology using k-shell decomposition. *Proc. Natl. Acad. Sci. USA* **104**, 11150–11154 (2007).
- [14] Live Journal, www.livejournal.com
- [15] Zhou, S. <http://www.cs.ucl.ac.uk/people/S.Zhou.html>
- [16] Liljeros, F., Giesecke, J. & Holme, P. The contact network of inpatients in a regional healthcare system. A longitudinal case study. *Mathematical Population Studies* **14**, 269–284 (2007).

- [17] *The Internet Movie Database* www.imdb.com
- [18] Hethcote, H. W. The mathematics of infectious diseases. *SIAM Rev.* **42**, 599–653 (2000).
- [19] Castellano, C., Fortunato, S. & Loretto V. Statistical Physics of Social Dynamics. *Rev. Mod. Phys.* **81** 591–646 (2009).
- [20] Hidalgo, C. A., Klinger, B., Barabasi, A-L., & Hausmann, R. The product space conditions the development of nations. *Science* **317**, 482–487 (2007).
- [21] Hethcote, H. & Rogers, J. A. *Gonorrhoea transmission dynamics and control* (New York, Springer-Verlag, 1984).
- [22] Pastor-Satorras, R. & Vespignani, A. Immunization of complex networks. *Phys Rev E* **65**, 036104 (2002).
- [23] Dezsó, Z. & Barabási, A.-L. Halting viruses in scale-free networks. *Phys. Rev. E* **65** 055103 (2002).
- [24] Cohen, R., Erez, K., ben-Avraham, D., & Havlin, S. Resilience of the Internet to random breakdowns. *Phys. Rev. Lett.* **85**, 4626–4630 (2000).
- [25] Newman, M. E. J. Properties of highly clustered networks. *Physical Review E* **68**, 026121 (2003).
- [26] Newman, M. E. J. Assortative mixing in networks. *Phys. Rev. Lett.* **89**, 20 (2002).
- [27] Large Network visualization tool, <http://xavier.informatics.indiana.edu/lanet-vi/>.
- [28] Alvarez-Hamelin, J. I., Dallásta, L., Barrat, A. & Vespignani, A. Large scale networks fingerprinting and visualization using the k-core decomposition. *Advances in Neural Information Processing Systems* **18**, 41–51 (2006).

Acknowledgements: We thank NSF-SES, ONR, Epiwork, and the Israel Science Foundation for support. FL is supported by Riksbankens Jubileumsfond. We thank L. Braunstein, J. Brujic, kc claffy, D. Krioukov, and C. Song for valuable discussions, and Shi Zhou for providing us with the email dataset.

FIG 1. When the hubs may not be good spreaders. **a-c** The extent of the efficiency of the spreading process cannot be accurately predicted based on a measure of the immediate neighborhood of the node, such as the degree k . For the contact network of inpatients (CNI), we compare infections originating from single nodes having the same degree $k = 96$ (nodes A and B) or the same index $k_S = 63$ (nodes A and C), with infection probability $\beta = 0.035$. In the corresponding plots, the colors indicate the probability that a node will be infected when spreading starts in the corresponding origin, as long as this probability is higher than 25%. In the first case, where origin A has $k_S = 63$, spreading reaches a much wider area more frequently, in contrast to origin B ($k_S = 26$), where the infection remains largely localized in the immediate neighborhood of B. Spreading is very similar between origins A and C, which have the same k_S value, although the degree of C is much smaller than A. The importance of the network organization is also highlighted when we randomly rewire the network (preserving the same degree for all nodes). In this case the standard picture is recovered: the extent of spreading coincides and both hubs contribute equally largely to spreading. **d**, A schematic representation of a network under the k -shell decomposition. Nodes are assigned to k -shell layers according to their remaining degree, which is obtained by successive pruning of nodes with degree smaller than the k_S value of the current layer. We start by removing all nodes with degree $k = 1$. After removing all the nodes with $k = 1$, some nodes may be left with one link, so we continue pruning the system iteratively until there is no node left with $k = 1$ in the network. The removed nodes, along with the corresponding links, form a k -shell with index $k_S = 1$. In a similar fashion, we iteratively remove the next k -shell, $k_S = 2$, and continue removing higher k -shells until all nodes are removed. As a result, each node is associated with a unique k_S index, and the network can be viewed as the union of all k -shells. The resulting classification of a node can be very different than when the degree k is used. For example, the two nodes of degree $k = 8$ (blue and yellow nodes) in this network are in different locations: one lies in the periphery, ($k_S = 1$) while the other hub is in the innermost core of the network, i.e. it has the largest k -shell index ($k_S = 3$).

FIG 2. The k -shell index predicts the outcome of spreading more reliably than the degree k or the betweenness centrality C_B . The networks used are (top to bottom): email contacts ($\beta = 8\%$), CNI network ($\beta = 4\%$), the actors network ($\beta = 1\%$), and the Livejournal.com friendship network ($\beta = 1.5\%$). **a, c, e, g** Average infected size of

the population $M(k_S, k)$ when spreading originates in nodes with (k_S, k) . **b, d, f, h** The infected size $M(k_S, C_B)$ when spreading originates in nodes of a given combination of k_S and C_B . In both cases, spreading is larger for nodes of higher k_S , while nodes of a given k or C_B value can result in either small or large spreading, depending on the value of k_S . (There is an exception at large k_S and small k of the livejournal database, which is due to artificial closed groups of virtual characters that connect with each other for the purpose of online gaming and do not correspond to regular users, as the rest of the database.)

FIG 3. k -shell structure of the contact network of inpatients. **a**, The imprecision function $\epsilon(p)$ tests the merit of using k -shell, k and C_B to identify the most efficient spreaders in the CNI networks. For $\beta = 4\%$ and a given fraction of the system p we first identify the Np most efficient spreaders as measured by M_i (which we designate by Υ_{eff}). Similarly, we identify the Np individuals with the highest k -shell index (Υ_{k_S}). We define the imprecision of k -shell identification as $\epsilon_{k_S}(p) \equiv 1 - M_{k_S}/M_{\text{eff}}$, where M_{k_S} and M_{eff} are the average infected percentages averaged over the Υ_{k_S} and Υ_{eff} groups of nodes respectively. ϵ_k and ϵ_{C_B} are defined similar to ϵ_{k_S} . Even though both k -shell and k identification strategies yield comparable results for $p = 2\%$, the k -shell strategy is consistently more accurate for $2\% < p < 10\%$ with ϵ_{k_S} approximately twice lower than ϵ_k . The C_B identification of the most efficient spreaders is the least accurate, with ϵ_{C_B} exceeding 40%. **b**, We visualize the CNI network as a set of concentric circles of nodes representing inpatients, each circle corresponding to a particular k -shell layer. k_S indices of the k -shell layers increase as one moves from the periphery to the center of the network [27, 28]. Node size is proportional to the logarithm of the degree of the node. We highlight the 25 inpatients with the largest degree values. Note that inpatients with high k values are not concentrated at the “center” of the network but instead are scattered throughout different k -shell layers. We highlight the position of the three nodes A, B, and C, of the origins that were used in the example of Fig. 1. **c**, Scatter plot of the node degree k as a function of k_S for all the nodes in the CNI network (black symbols) and the degree-preserving randomized version of the same network (red symbols). Note that there are many inpatients with large k and low k_S values in the original network while in the randomized email network all the hubs are located in the inner core of the network. We also show the position of the three origins used in Fig. 1. **d**, When spreading starts from multiple origins, the set of nodes with highest degree (blue continuous line) can spread significantly more than the set of highest- k_S nodes (red continuous line),

because in the latter case most of these nodes are connected to each other. If we only consider in this set nodes that are not directly linked, then both the sets of highest k or k_S nodes yield a similar result (dashed lines), where spreading is significantly enhanced. Results are shown for $\beta = 3\%$ in the CNI.

FIG 4. SIS spreading in the Contact Network of Inpatients. 20% of the individuals are initially infected. **a, b**, Virus persistence $\rho(k_S, k)$ as a function of k and k_S values of inpatients in the CNI network for, $\beta = 2\%$, and $\beta = 4\%$, respectively. The infection survives mainly in nodes with large k_S values. **c**, We form four groups of nodes of the CNI network based on their k -shell values. For all values of β , virus persistence is consistently higher in the inner k -shell layers. Note that the infection survives locally in high k -shells even when $\beta < \beta_c$.

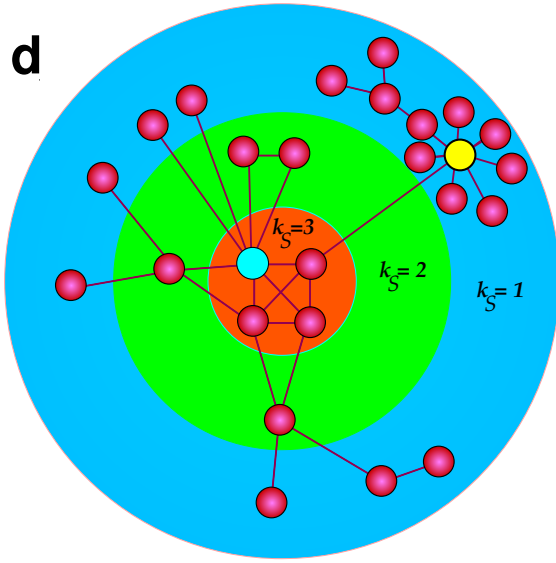
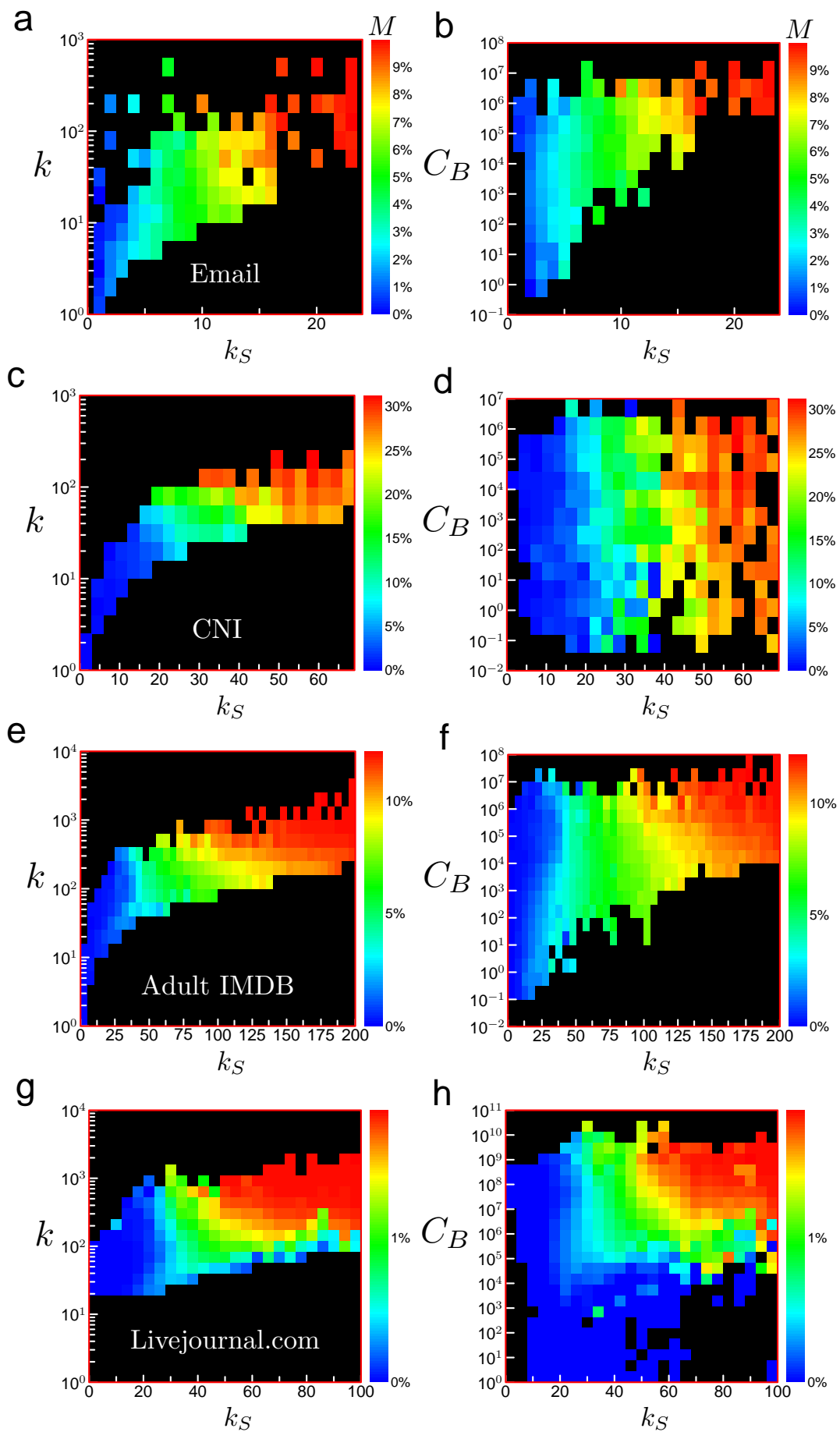


FIG. 1:



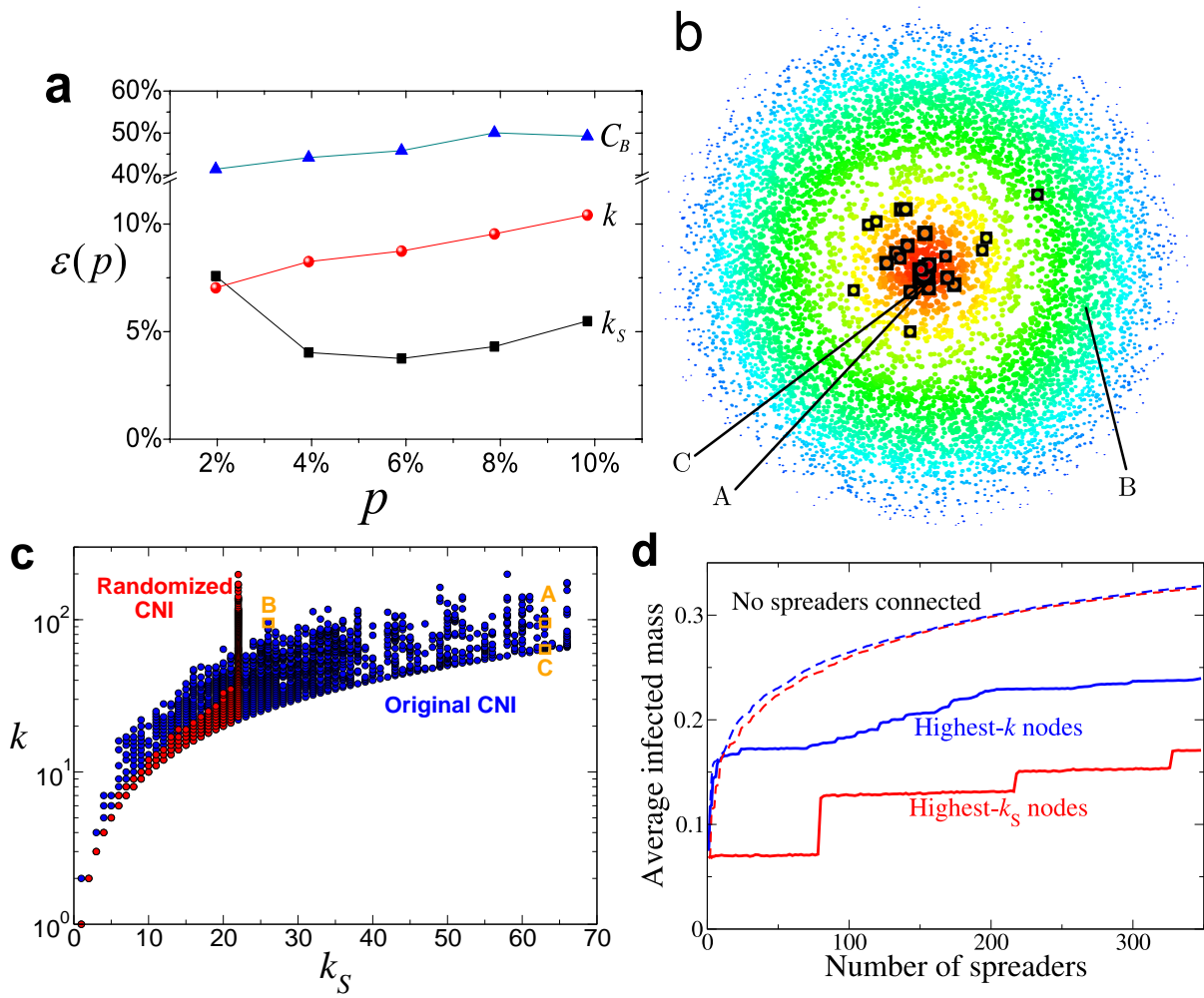


FIG. 3:

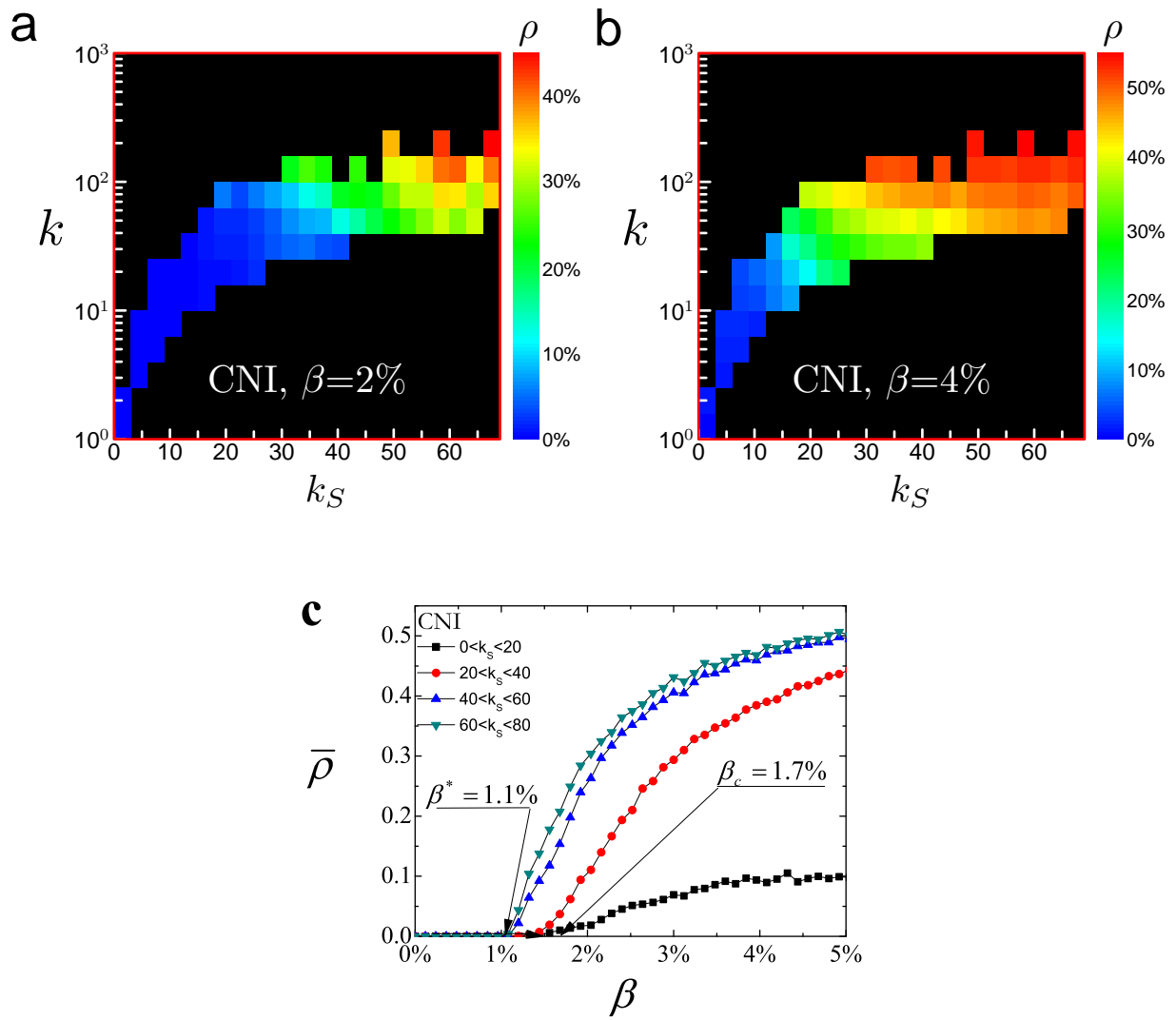


FIG. 4:

Appendix

I. DATASETS

In this study we have mainly focused on social networks, but our results can be extended to networks from practically any discipline. The datasets that were used in the paper are the following:

a) *Contact Network of Inpatients*. We use records from Swedish hospitals [16] and establish a link between two inpatients if they have both been hospitalized in the same quarters. We restrict the recording period to one week. There are 8622 inpatients in the largest component, with an average degree of around 35.1.

b) *IMDB actors in adult films*. We have created a network of connections between actors who have co-starred in films, whose genre has been labeled by the Internet Movie Database [17] as ‘adult’. This network is a largely isolated sub-set of the original actor collaboration network. Additionally, all these films have been produced during the last few decades, rendering the network more focused in time. The largest component comprises 47719 actors/actresses in 39397 films. The average degree of the network is 46.0.

c) *Email Contact Network*. The network of email contacts is based on email messages sent and received at the Computer Sciences Department of London’s Global University. The data have been collected in the time window between December 2006 and May 2007. Nodes in the network represent email accounts. We connect two email accounts with an undirected link in the case where at least two emails have been exchanged between the accounts (at least one email in each direction). There are 12701 nodes with an average degree of 3.2.

d) *LiveJournal.com*. The network of friends in the LiveJournal community, as recorded in a 2008 snapshot. We only consider reciprocal links, i.e. when two members are in each other’s list of friends. There are 3453394 nodes in the largest component, and the average degree is 12.4.

e) *Cond-mat collaboration network*. This is the network of collaborations between scientists that have posted reprints in the ‘cond-mat’ e-print archive, between 1995 and 2005. The nodes of the network represent the authors, who are connected if they have co-authored at least one paper. The cond-mat collaboration dataset consists of 17628 authors with average degree 6.0

Network Name	N	N_E	$\langle k \rangle$	$\langle k^2 \rangle$	β_c	β	$k_{S_{max}}$
Contact Network of Inpatients	8622	151649	35.1	1633	1.7%	4%	66
Actor Network	44719	1028537	46.0	17483	0.21%	1%	199
Email Contacts	12701	20417	3.2	351.1	0.73%	8%	23
Live Journal	3453394	21378154	12.38	892.45	1.1%	1.5%	100
Cond-mat Collaboration Network	17628	52884	7.0	109.4	5.1%	10%	22
RL Internet	493312	808844	3.3	71.9	4.6%	6%	36
AS Internet	20556	62920	6.1	2111.2	0.23%	n/a	41
Product Space	765	40164	104.8	16931	0.0050%	n/a	100

TABLE I: Properties of the real-world networks studied in this work. Here N is the number of nodes, N_E is the number of edges, $\langle k \rangle$ is the average degree in the network, $\langle k^2 \rangle$ is the average squared degree in the network, β_c is the epidemic threshold ($\beta_c \approx \lambda \langle k \rangle / \langle k^2 \rangle$), $\lambda = 0.8$ in SIS simulations, β is the value we used in SIR simulations and $k_{S_{max}}$ is the highest k -shell index of the network. We consider only the largest connected cluster of the network if the original network is disconnected.

f) *The Internet at the router level (RL)*. The nodes of the RL Internet network are the Internet routers. Two routers are connected if there exists a physical connection between them. Data have been gathered from the DIMES project [13]. The largest connected component of the analyzed dataset contains 493312 routers with an average degree of 3.3.

g) *The Internet at the autonomous system level (AS)*. The nodes are autonomous systems which are connected if there exists a physical connection between them. An autonomous system is a collection of connected IP routing prefixes under the control of one or more network operators that presents a common, clearly defined routing policy to the Internet. Data have been gathered by the DIMES project [13]. The largest connected component of the AS Internet consists of 20556 autonomous systems with average degree 6.1.

h) *Product space of economic goods*. This is the network of proximity between products according to Ref. [20]. We use a proximity threshold 0.3, and we recover similar results for different thresholds, as well.

We outline some of the basic properties for these networks in Table I.

II. THE k -SHELL DECOMPOSITION METHOD

In order to classify the nodes into k -shells we employ the k -shell decomposition algorithm. First, we remove all nodes with degree $k=1$. After this first stage of pruning there may appear new nodes with $k=1$. We keep on pruning these nodes, as well, until all nodes with degree $k=1$ are removed. The removed nodes along with the links connecting them form the $k_S = 1$ k -shell. Next, we repeat the pruning process in a similar way for the nodes of degree $k=2$ to extract the $k_S = 2$ k -shell and subsequently for higher values of k until all nodes are removed. As a result, the network can be viewed as a set of adjacent k -shells (see Fig. 5).

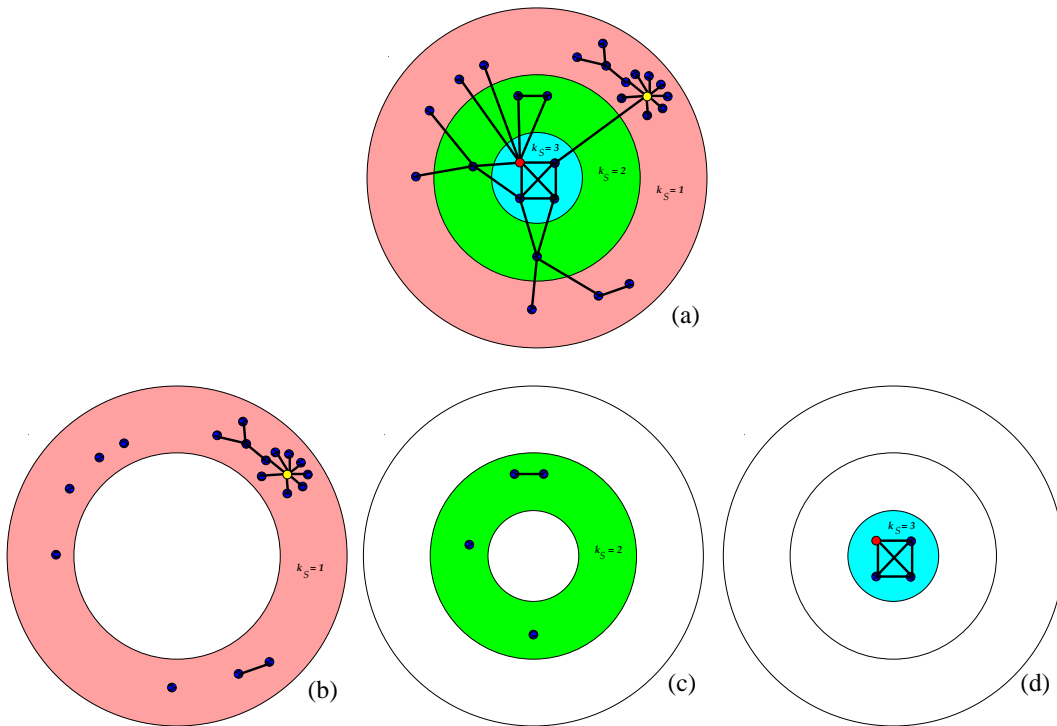


FIG. 5: **The illustration of the k -shell extraction method.** **a**, A schematic network is represented as a set of 3 successively enclosed k -shells labeled accordingly. **b**, Nodes with edges forming $k_S = 1$ shell of the network. **c**, Nodes with edges forming $k_S = 2$ shell of the network. **d**, Nodes with edges forming $k_S = 3$ shell of the network.

The k -shell decomposition method assigns a unique k_S value to each node, that corresponds to the index of the k -shell this node belongs to. The k_S index provides a different type of information on a node than that provided by the degree k . By definition, a given k -shell layer with index k_S can be occupied with nodes of degree $k \geq k_S$. In the case of

random model networks, such as the configurational model, there is a strong correlation between k and the k_S index of a node and, therefore, both quantities provide the same type of information. Thus, the low-degree nodes are generally in the periphery, and the high-degree nodes are generally in the innermost k -shells. In real networks, however, this relation is often not true. In real networks hubs may have very different k_S values and can be located both in the periphery (yellow node in Fig. 5) or in the core (red node in Fig. 5) of the network.

III. PROBABILITY AND TIME OF INFECTION

We have demonstrated that the location of a node, as described through the k_S index, is important for the extent of spreading M_i when this node is the spreading origin. Here, we show that nodes with high k_S are more probable to be infected during an epidemic outbreak and are infected earlier than nodes with low k_S , when spreading starts at a random node. We introduce the quantity E_i , as the probability that a node i is going to be infected during an epidemic outbreak originating at a random location, and T_i , as the average time before node i is infected during the same process.

As shown in Figs. 6a-d all three quantities that characterize the role of a node in an epidemics process, M_i , E_i and T_i are strongly correlated. The nodes that are infected by a given node i form a cluster of size \overline{M}_i , and they are statistically the nodes that can reach i when they act as origins themselves. Thus, the probability E_i to reach this node in general is directly proportional to the size M_i , as shown in the plots. The average time T_i to reach a node is inversely proportional to its spreading efficiency M_i , which emphasizes the fact that these nodes are easily reachable from different network locations. In conclusion, the nodes with the largest k_S values consistently a) are infecting larger parts of the network, b) are infected more frequently, and c) are infected earlier, than nodes with smaller k_S values.

IV. THE IMPRECISION FUNCTIONS

We quantify the spreading efficiency of an individual origin i through the infected number of nodes M_i . In order to compare the different methods, we rank all network nodes according to their spreading efficiency, independently of their other properties, and we consider a

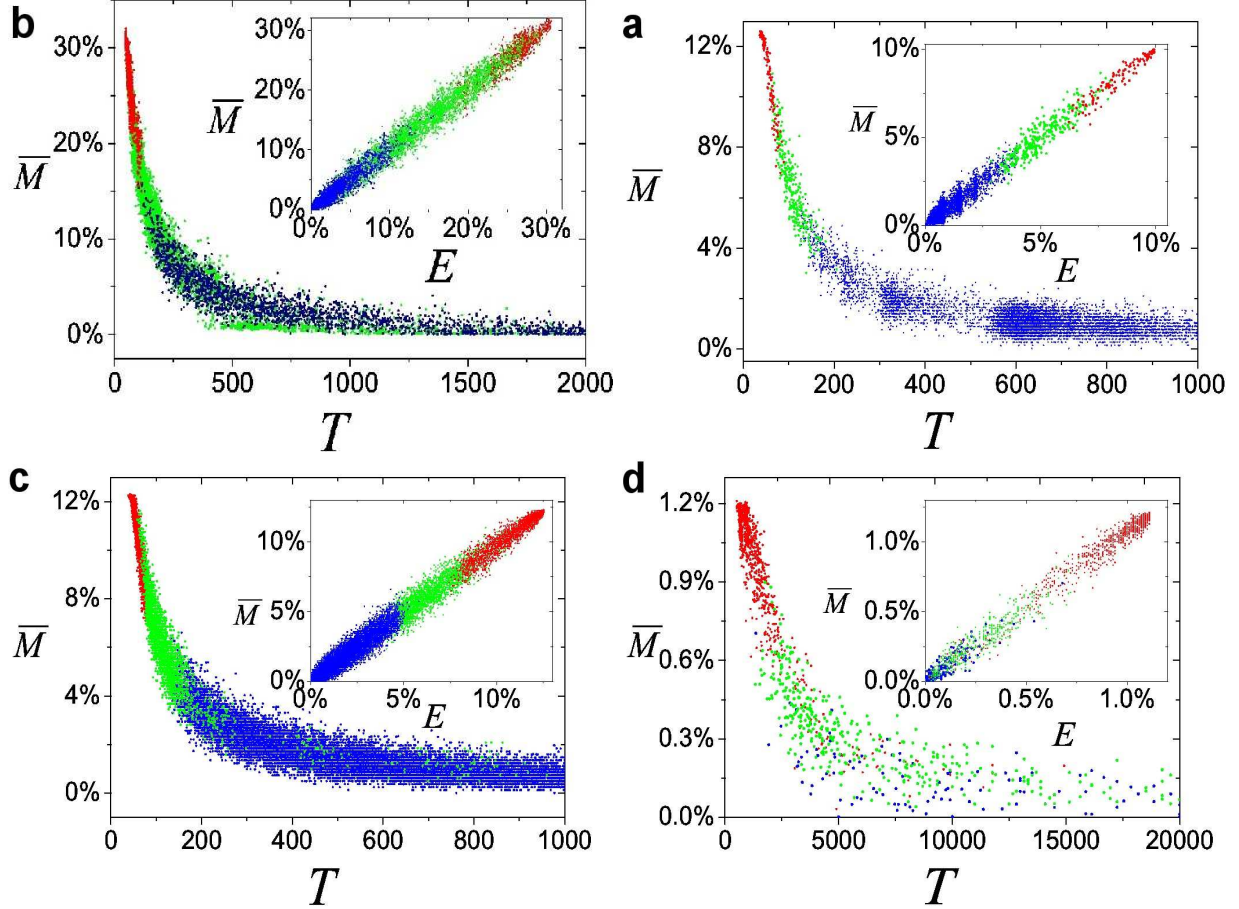


FIG. 6: Cross-plots of M_i as a function of T_i , and M_i as a function of E_i (inset) for a) hospital inpatients, b) email, c) actor network and d) RL Internet. Every point denotes the corresponding quantities for a given node, and the color denotes the k -shell layer of this node. The k_S values are aggregated and highlighted with red (large k_S regime), green (intermediate k_S regime) and blue (low k_S values) colors, respectively. A high level of correlation between M_i and E_i indicates that the most efficient spreaders (as measured by M_i) are the most likely to be infected during an epidemic outbreak originating at random inpatient in the network. On the other hand, the anti-correlation between M_i and T_i indicates that the most efficient spreaders are typically infected earlier than other nodes during an epidemic outbreak.

fraction p of the most efficient spreaders ($p \in [0, 1]$). We designate this set by $\Upsilon_{eff}(p)$. Similarly, we define $\Upsilon_{k_S}(p)$ as the set of individuals with highest k -shell values. In order to assess the merit of using k -shell decomposition to identify the most efficient SIR spreaders one needs to compare the two sets $\Upsilon_{eff}(p)$ and $\Upsilon_{k_S}(p)$. In order to consider individual M_i

values, we calculate the average $M_{eff}(p)$ and $M_{k_S}(p)$ values for the sets $\Upsilon_{eff}(p)$ and $\Upsilon_{k_S}(p)$ respectively: $M_{k_S}(p) \equiv \sum_{i \in \Upsilon_{k_S}(p)} M_i / Np$ and $M_{eff}(p) \equiv \sum_{i \in \Upsilon_{eff}(p)} M_i / Np$, where Np is the number of nodes that we consider in the comparison. By definition, $M_{eff}(p) \geq M_{k_S}(p)$, and the equality is only reached if $\Upsilon_{eff}(p) = \Upsilon_{k_S}(p)$. We assess the imprecision of k -shell identification by calculating the ratio between $M_{eff}(p)$ and $M_{k_S}(p)$:

$$\epsilon_{k_S}(p) \equiv 1 - \frac{M_{k_S}(p)}{M_{eff}(p)}. \quad (3)$$

Similarly, we can define $\epsilon_k(p)$ and $\epsilon_{C_B}(p)$:

$$\epsilon_k(p) \equiv 1 - \frac{M_k(p)}{M_{eff}(p)}, \quad \epsilon_{C_B}(p) \equiv 1 - \frac{M_{C_B}(p)}{M_{eff}(p)}. \quad (4)$$

A value for ϵ close to 0 denotes a very efficient process, since the nodes that are chosen are practically those that contribute most to epidemics. In all cases, the k_S method yields a spreading that is closer to the optimum than either the degree or the betweenness centrality. Additionally, this behavior is independent on the fraction of spreaders p that we consider in each case.

V. SIR SPREADING EFFICIENCY

In the main text we present results for $M(k_S, k)$ for the email network, the CNI, the actor network and the Livejournal network. Here, we present additional results of the k -shell analysis of the scientific collaboration network and the Internet at the Router Level (RL). Figure 8 shows the results for $M(k_S, k)$ and $M(k_S, C_B)$. The conclusion on the spreading importance of high k_S nodes is exactly the same as for the social networks in the main text.

We also highlight the location of the 25 largest hubs in the k -shell structure of the studied networks. Fig. 9 shows the results for the collaboration, actor, email, RL Internet, AS Internet, and Product space networks. High-degree nodes in most of the studied networks are scattered at different k -shell layers: the high- k nodes appear both in the periphery (starting as low as $k_S = 1$) and in the network center (large k_S value). In certain cases, such as in the actors network, the largest hubs are located in the highest k_S layers. The relation of k_S and k in the AS Internet and the product space is strongly monotonic, and there are very few nodes where k_S is large or small compared to the degree k . This is a typical

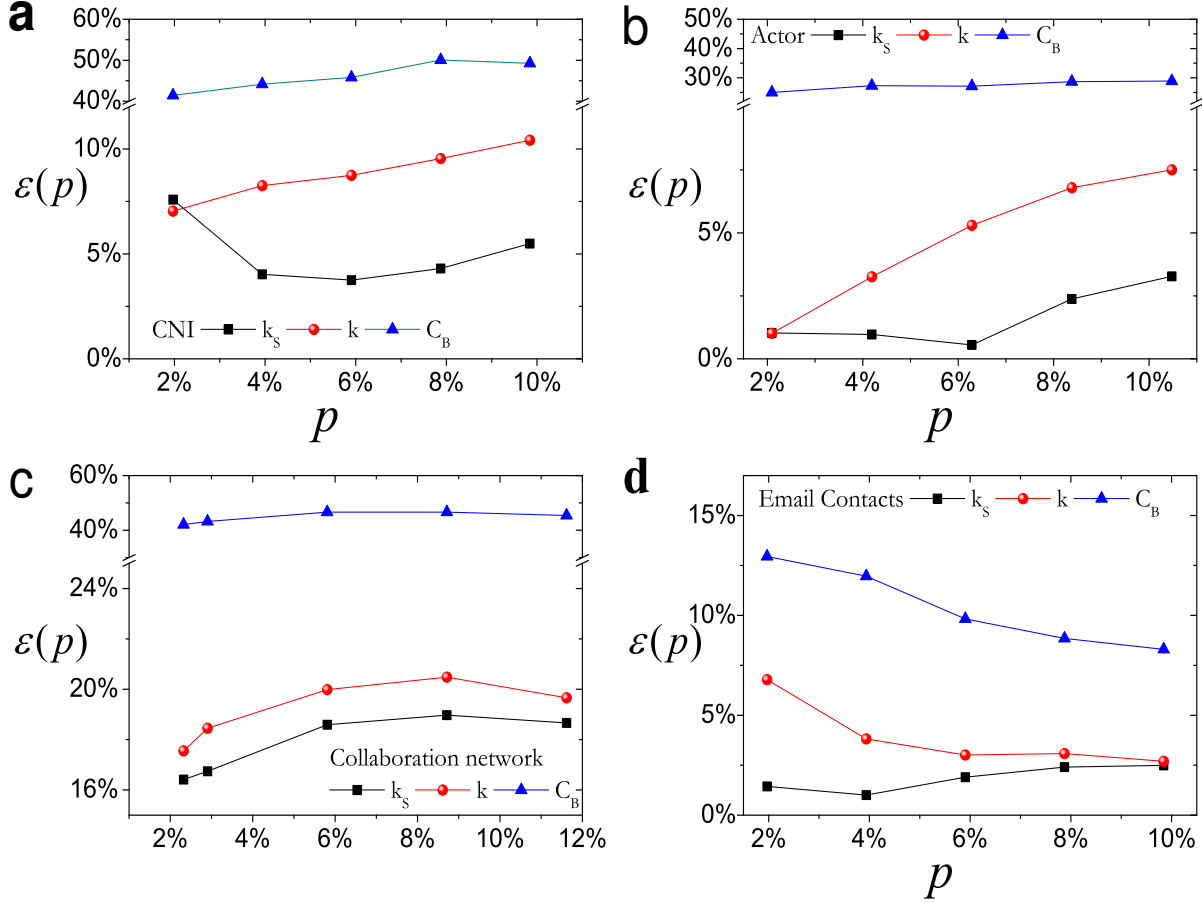


FIG. 7: The imprecision functions $\epsilon(p)$ test the merit of using k -shell, k and C_B to identify the most efficient spreaders in the CNI, actor, collaboration, and email contact networks. The k -shell based identification method yields consistently lower imprecision compared to the k and C_B based methods.

behavior for random networks, and the structure of these two networks is significantly close to their randomized counterparts. In these cases, choosing a node based on its degree or its k -shell index does not make a difference, since they practically lead to the same nodes.

The complex organization of the nodes in the k -shells is highlighted when we randomly rewire the links in the networks, yet preserving the nodes degree. This rewiring ‘restores’ all the hubs to the innermost k -shell of the system and imposes a strict hierarchy of nodes in terms of both k and k_S . The bottom row of plots in Fig. 9 shows the scatter-plots of degree k as a function of k -shell index k_S for every node in the network. In all cases, a monotonic relation of k vs k_S is followed in the ‘rewired’ networks (red symbols), where now all the hubs appear in the highest k -shell) as opposed to the weak correlation between k and k_S in

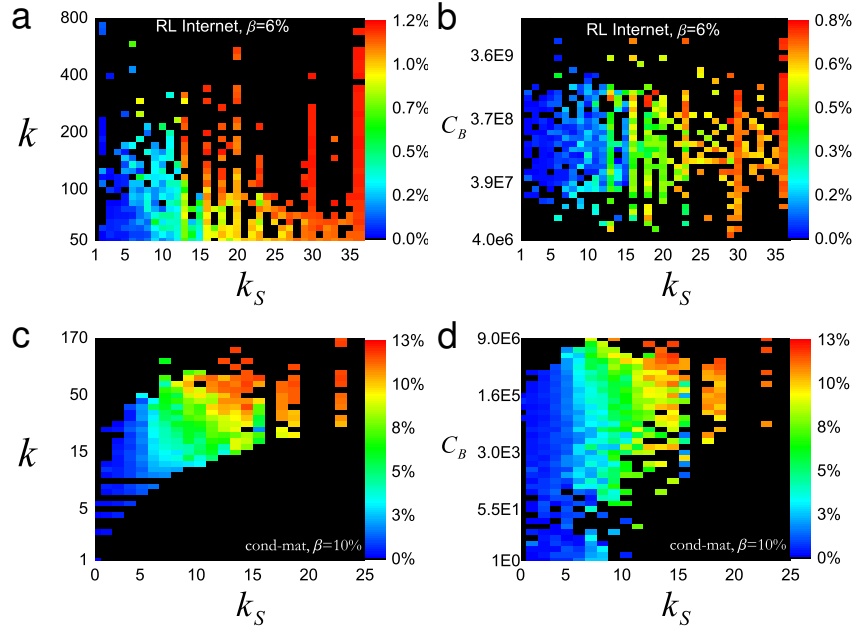


FIG. 8: The shell index k_S predicts the outcome of spreading more reliably than the degree k or the betweenness centrality C_B . The networks that were analyzed are: (a, b) the RL Internet and (c, d) the collaboration network. a and c, The average infected size $M(k_S, k)$ as a function of (k_S, k) values of the infection origin nodes. b and d, The average infected size $M(k_S, C_B)$ as a function of (k_S, C_B) values of the infection origin nodes.

the original networks (shown in black).

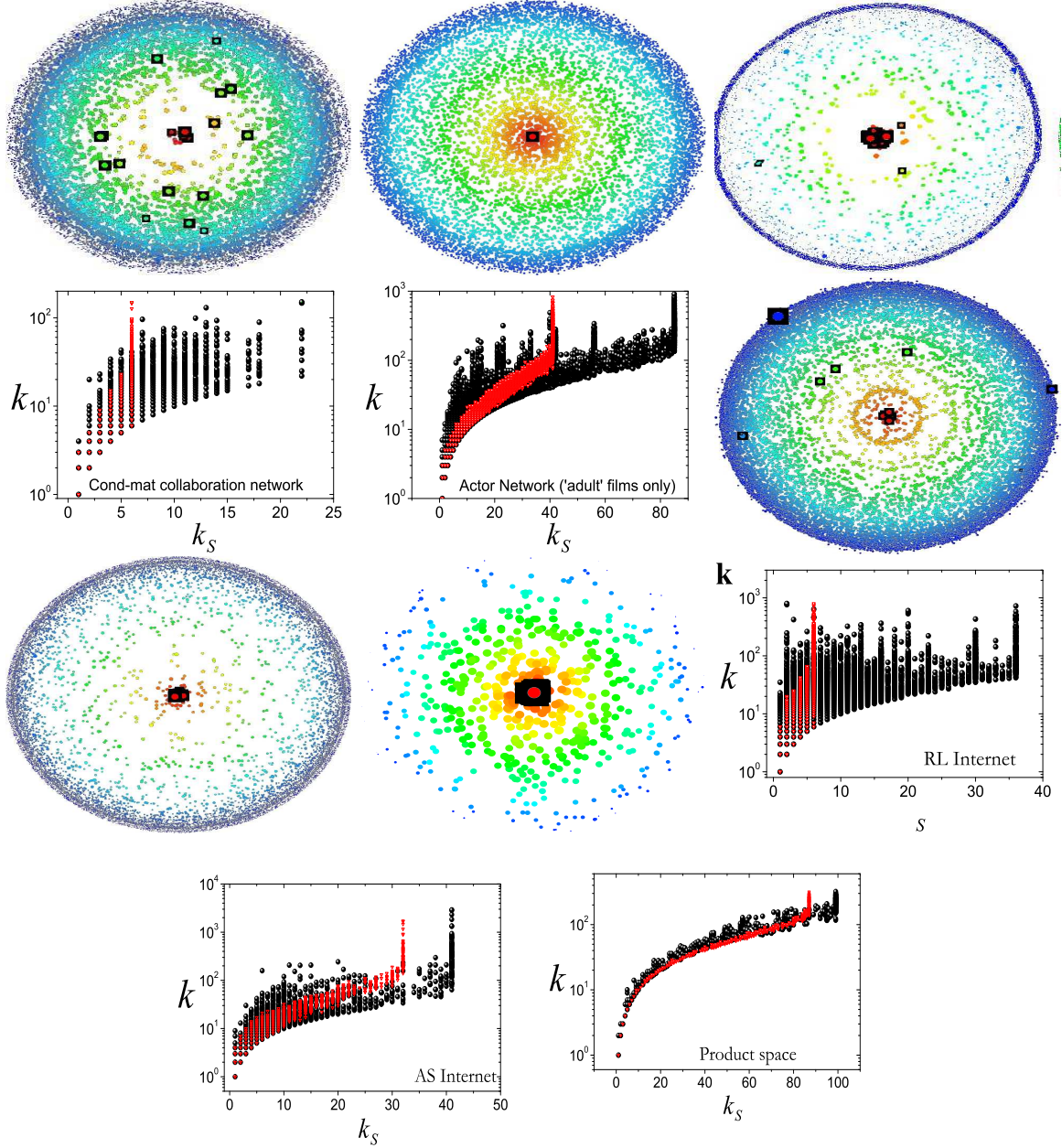


FIG. 9: *k*-shell structure of the analyzed networks. (**Top row**): Visualization of the *k*-shell structure. We represent networks as sets of concentric circles of nodes, each one corresponding to the particular *k*-shell, with low k_S values in the periphery and large k_S values towards the center of the network. The size of each visualized node is proportional to the logarithm of its degree value. We highlight the 25 highest degree nodes with black squares. Many of the hubs are found in outer layers. (**Bottom row**): Scatter plots of node degree k as a function of its *k*-shell index k_S for the original networks (black symbols) and the degree-preserving randomized version of the networks (red symbols). The networks correspond to: the cond-mat collaboration network, the actor network, the email contact network, the RL Internet, the AS Internet, and the Product Space network.

VI. REWIRING HIGHLIGHTS THE IMPORTANCE OF k -SHELL

In Figs. 1a and 1b of the main text we show that the extent of infection can be remarkably different, although we start from two origins with similar degree. The importance of the structure in the dynamics of spreading can be highlighted if we randomly rewire the network. During this process the original degrees of all nodes are preserved, but random neighbors are chosen for each node, destroying thus any correlations and any patterns in the local connectivity. We denote by $P(M|i)$ the probability that a percentage M of the total population will be infected if a disease originates on node i . In Figs. 1a,b of the main text and in Fig.10a we show that two nodes #1 and #2 with similar degree may yield markedly different distributions $P(M|1)$ and $P(M|2)$. After rewiring, these distributions become practically indistinguishable (see Fig. 10b).

VII. VIRUS PERSISTENCE IN SIS

Many infectious diseases, including most sexually transmitted infections, do not confer immunity after infection, so that they cannot be described via the SIR model. These cases are better simulated through the SIS epidemic model [18]. The dynamics of SIS epidemics is different, since the number of infected nodes eventually reaches a dynamic equilibrium “endemic” state at which exactly as many infectious individuals become susceptible as susceptible nodes become infected [18]. The quantity characterizing the role of nodes in SIS spreading is the persistence, $\rho_i(t)$, defined as the probability that node i is infected at time t [7]. In an endemic SIS state, which is reached asymptotically, ρ_i becomes independent of t . The persistence ρ has been shown to be higher in hubs which are reinfected frequently due to the large number of their neighbors [7, 22, 23]. To uncover the role of k -shell layers in SIS spreading we use the joint persistence function

$$\rho(k_S, k) \equiv \sum_{i \in \Upsilon(k_S, k)} \frac{\rho_i}{N(k_S, k)}. \quad (5)$$

Here we present results for the virus persistence in the Actor, Collaboration, Email and RL Internet Networks. Similar to Fig. 4, we depict $\rho(k_S, k)$ in both supercritical ($\beta > \beta_c$) and subcritical ($\beta^* < \beta < \beta_c$) regimes. In the supercritical regime, $\rho(k, k_S)$ increases with both k and k_S , with maximum values corresponding to hubs in the innermost k -shell layers

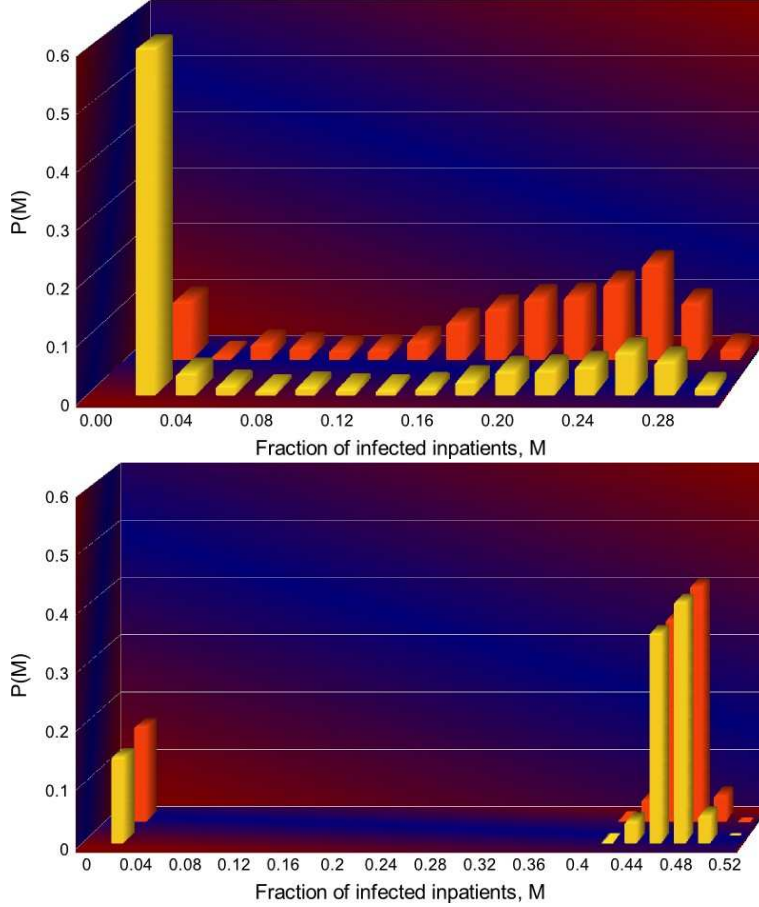


FIG. 10: **Why the hubs may not be good spreaders.** The probability distribution $P(M|i)$ of the infected percentage for the contact network of inpatients, when the epidemic starts at two of the origin hubs in Fig. 1 $i = A, B$ with the same degree ($k = 96$), but different k_S values ($k_S = 63$ and $k_S = 26$, respectively). In each histogram, we use 1000 random realizations of the simulation, starting an SIR epidemic from the same given origin i . Despite the fact that the two origins of the epidemic spreading have the same degree, the two histograms present a radically different character. In one case (red histogram), the hub infects up to 30% of the population, while most of the spreading attempts from the other hub (yellow histogram) practically cannot propagate the infection at all. The importance of the organization of the network is highlighted when we randomly rewire the network (preserving the same degree for all nodes). In this case both distributions $P(M|A)$ and $P(M|B)$ coincide and both hubs contribute equally to spreading. Notice also that spreading in the rewired network extends over a much larger size of the population.

(see Fig. 12). As depicted in Fig. 12, in the subcritical regime, viruses persist only in the

highest k -shell layers, while the probability of finding an infected node in low k -shell layers is negligible.

In order to determine in the above networks the actual epidemic threshold β^* we study the behavior of SIS spreading for the entire range of β values: $0 < \beta < \beta_c$. In order to highlight the role of k -shell layers in spreading, we organize several groups of nodes based on the k -shell layers of each network. Every such group comprises approximately 100 randomly chosen nodes with the corresponding k -shell indices. In order to achieve similar average degree in each of the groups, we pick nodes with uniform probability based on their degree. As shown in Fig. 11, virus persistence is consistently higher in the inner k -shell layers for all values of β . Moreover, we find substantially lower epidemic thresholds $\beta^* < \beta_c$ in all considered networks except for the Email Contact network.

The results of Figs. 11 and 12 suggest that the observed persistence of a virus is due to the dense sub-network formed by nodes in the innermost k -shell, which helps the virus to consistently survive locally in this area. Indeed, the innermost k -shell layers can be regarded as a small subgraph exclusively consisting of hubs. By definition, all nodes in this innermost k -shell will have degrees $k \geq k_{S_{max}}$. Therefore, as a simple approximation, one can regard the innermost core of a network as a regular graph consisting of nodes with the same degree $k = k_{S_{max}}$.

The mean-field solution of the SIS spreading in a regular graph can be found, for instance in Ref. [22]. We reproduce this solution below for the sake of convenience.

The master equation describing the time evolution at a mean-field level of the average density of infected individuals $\rho(t)$:

$$\frac{d\rho(t)}{dt} = -\rho(t) + \beta k \rho(t)(1 - \rho(t)), \quad (6)$$

where k is the degree of all nodes in the regular graph. The first term on the right hand side of Eq. (6) accounts for infected nodes becoming healthy. The second term on the right hand side of Eq. (6) accounts for healthy nodes becoming infected: a randomly chosen node is healthy with probability $1 - \rho(t)$, this healthy node can be infected by either of its k neighbor nodes with total probability of $\beta k \rho(t)$. The stationary endemic state is reached when $\frac{d\rho(t)}{dt} = 0$ which leads to

$$\rho = 1 - \frac{1}{\beta k}, \quad (7)$$

indicating the existence of a nonzero epidemic threshold of $\beta = 1/k$. The innermost core of a network consisting only of nodes with degrees $k \geq k_{S_{max}}$ will have epidemic threshold

$$\beta^* \leq 1/k_{S_{max}}. \quad (8)$$

The above inequality holds for all considered networks. Moreover, this inequality becomes an equality for CNI and collaboration networks where nearly all nodes in the innermost cores have degree $k \approx k_{S_{max}}$.

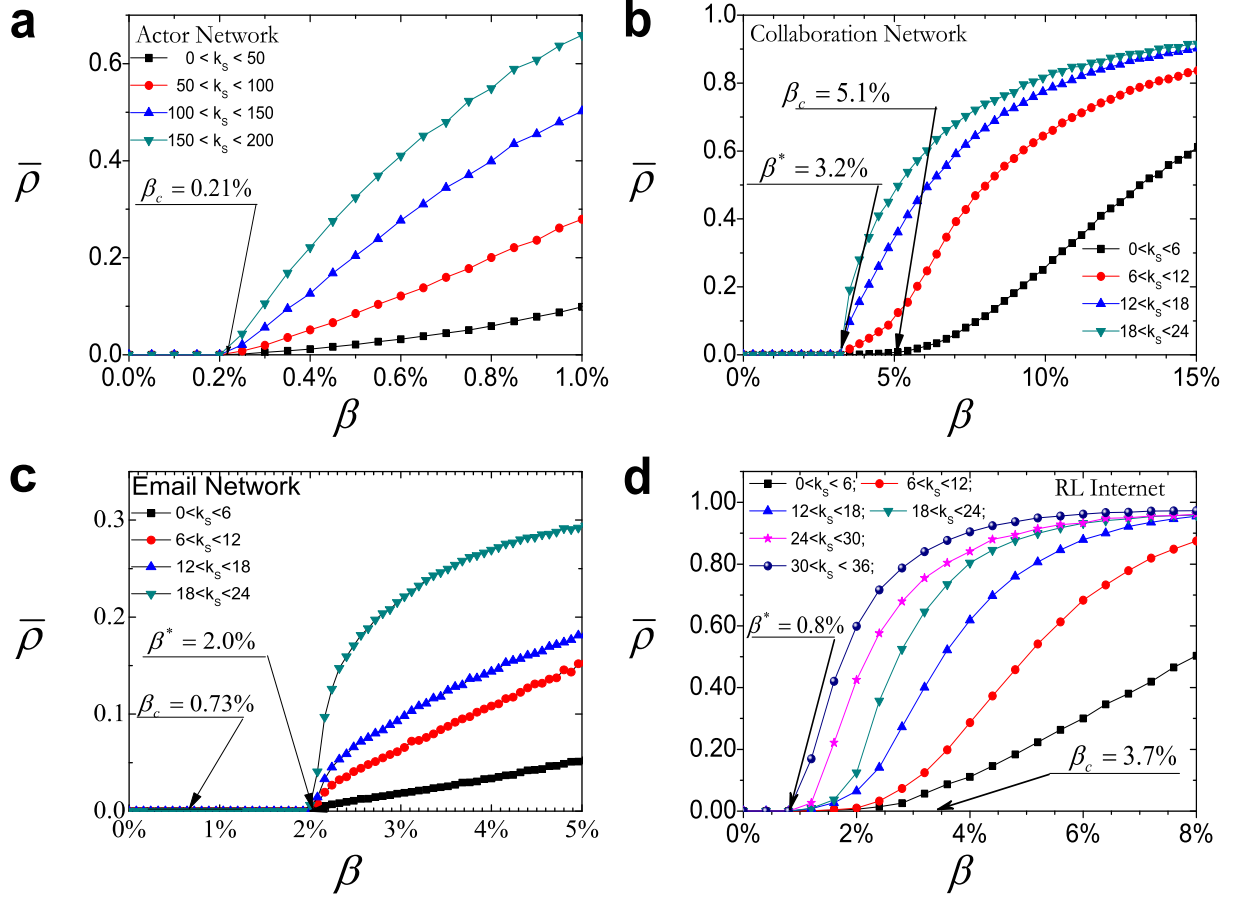


FIG. 11: **How average SIS persistence in different k -shell layers depends on virus contagiousness.** For every network we randomly sample several groups of nodes based on k -shell index. We plot the measure average virus persistence $\bar{\rho}$ for every group of nodes as a function of β for Email, Actor, Collaboration and RL Internet networks. Virus persistence is higher for nodes located in higher k -shell layers. In all considered networks besides the Email network viruses may persist even in the subcritical regime ($\beta^* < \beta < \beta_c$) owing to highly connected nodes in the innermost k -shell layers. The network of email contacts, on the contrary has fewer k -shell layers than its random counterpart [see Fig. 9]. As a result, its actual threshold β^* is larger than β_c .

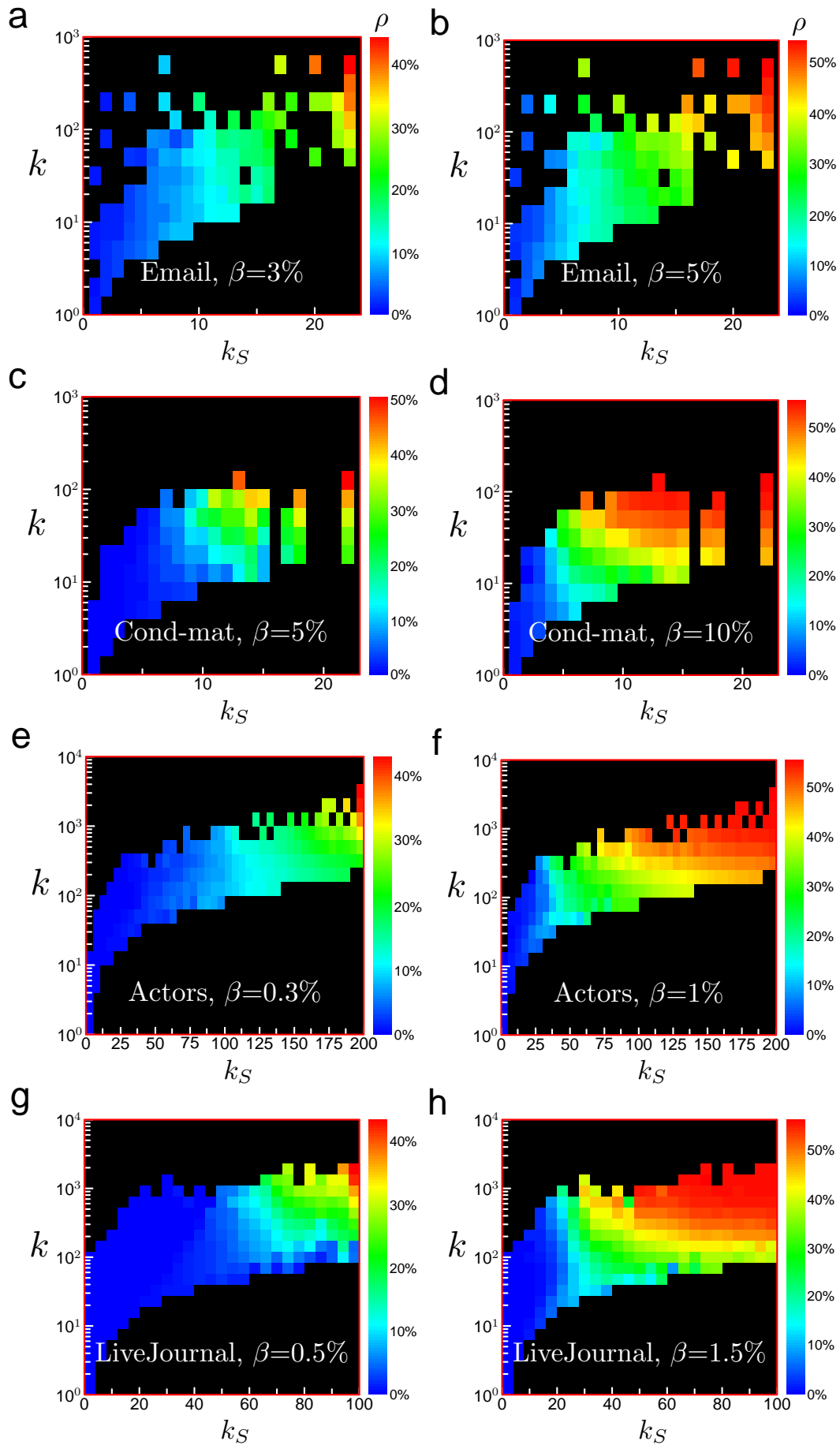


FIG. 12: SIS maps