

Communicability across evolving networks

Article

Accepted Version

Physics Review E accepted

Grindrod, P., Parsons, M. C., Higham, D. J. and Estrada, E. (2011) Communicability across evolving networks. Physical Review E, 83 (4). 046120. ISSN 1539-3755 doi: 10.1103/PhysRevE.83.046120 Available at <https://reading-clone.eprints-hosting.org/19357/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

To link to this article DOI: <http://dx.doi.org/10.1103/PhysRevE.83.046120>

Publisher: American Physical Society

Publisher statement: Copyright American Physical Society 2011

All outputs in CentAUR are protected by Intellectual Property Rights law, including copyright law. Copyright and IPR is retained by the creators or other copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

www.reading.ac.uk/centaur

CentAUR

Central Archive at the University of Reading

Reading's research outputs online

Communicability Across Evolving Networks

Peter Grindrod* Desmond J. Higham[†] Mark C. Parsons[‡]
Ernesto Estrada[§]

Abstract

Many natural and technological applications generate time ordered sequences of networks, defined over a fixed set of nodes; for example time-stamped information about ‘who phoned who’ or ‘who came into contact with who’ arise naturally in studies of communication and the spread of disease. Concepts and algorithms for static networks do not immediately carry through to this dynamic setting. For example, suppose A and B interact in the morning, and then B and C interact in the afternoon. Information, or disease, may then pass from A to C, but not vice versa. This subtlety is lost if we simply summarize using the daily aggregate network given by the chain A-B-C. However, using a natural definition of a walk on an evolving network, we show that classic centrality measures from the static setting can be extended in a computationally convenient manner. In particular, communicability indices can be computed to summarize the ability of each node to broadcast and receive information. The computations involve basic operations in linear algebra, and the asymmetry caused by time’s arrow is captured naturally through the non-commutativity of matrix-matrix multiplication. Illustrative examples are given for both synthetic and real-world communication data sets. We also discuss the use of the new centrality measures for real-time monitoring and prediction.

PACS: 89.75.Fb, 89.75.Hc, 84.40.Ua

*Department of Mathematics, University of Reading, UK

[†]Department of Mathematics and Statistics, University of Strathclyde, UK

[‡]Department of Mathematics, University of Reading, UK

[§]Department of Mathematics and Statistics & Department of Physics, University of Strathclyde, UK

1 Introduction

At the heart of network science are the well established mathematical fields of deterministic and random graph theory, with concepts such as connectedness, pathlength, diameter, degree and clique playing key roles [1, 2]. The motivation for this work is that a new type of time-dependent network-based object is emerging from a range of digital technologies that requires a fundamentally different way of thinking.

In Figure 1 we show a simple example of an evolving network, where undirected connections between a fixed set of seven nodes is recorded over three days. If we regard the links as representing communication, for example, by telephone or email, then we see that A may pass a message to C through the links $A \leftrightarrow B$ and $B \leftrightarrow G$ on day 1 and then through the links $G \leftrightarrow E$ and $E \leftrightarrow C$ on day 2. However, there is no way for C to pass a message to A. Analogously, if the links represent physical proximity, then A may pass an infection to C but C cannot cause A to be infected. This asymmetry, which arises even though each individual network is symmetric, is caused by the arrow of time. It is clear that simply aggregating the individual networks would present a very misleading summary. This highlights a fundamental gap between the static and dynamic cases, and points out the need for a theory of evolving networks that

1. deals with the time ordering inherent in the edge lists when considering communication around the network,
2. respects the inherent asymmetry imposed by the arrow of time, even when each individual snapshot consists of an undirected network.

Many application areas give rise to connectivity patterns that change over time in this manner. As well as the traditional context of individual-to-individual contacts in epidemiology [3], the digital revolution is generating novel large scale examples, including

- networks of mobile users with a link denoting current “interaction”, i.e., either copresence in a location or logged contact through their mobile devices [4],
- networks of online social users (e.g., Facebook) interacting through messaging [4] or online chatting systems (e.g. MSN) [5],

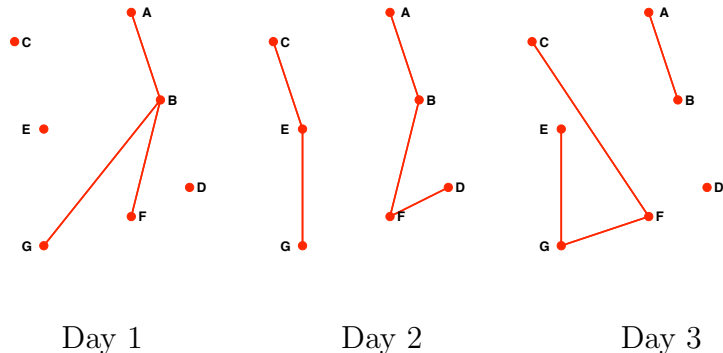


Figure 1: (Color online) Simple example of an evolving network.

- networks of travellers, vehicles or available routes defined over a dynamic transportation infrastructure [6, 7, 8],
- networks describing transient social interactions over cyberspace [9],
- networks describing individuals' attendance at regularly scheduled events over time [10],
- correlated neural activity in response to a functional task [11].

In this work, we show how centrality concepts that have proved useful for determining important nodes in static networks can be extended to this dynamic setting. Our approach is related to that of [12, 9, 13], in the sense that static graph concepts are directly generalized in a manner that respects the time dependency, but we take a walk counting viewpoint and focus on the type of centrality measures that are popular for social networks [14]. Earlier studies have also dealt with time-respecting paths. Berman [7] looked at a more restricted class of dynamic networks, where each edge has a single *start time* and a single *finish time*, and focused on issues such as the minimum number of nodes that must be deleted in order to disconnect all paths between a given pair of nodes. Further algorithmic and combinatorial issues

were considered in [15] for the case where each edge exists at a single instant of time. More recent work in [16] deals with the setting where each edge may exist at more than one time instance, and focuses on how quickly information or disease may spread across a contact network, whereas [17] looks at identifying edges that have the potential to pass information while it is fresh. We also note that dynamic networks are treated in [18], but the emphasis in that work is to discover communities, and a different approach is used, where extra links are added to represent the passage of time.

Let us emphasize at this stage that unlike in the well-studied ‘network growth’ context, where new nodes and accompanying edges are accumulated and only the final, aggregate network is of interest [19, 2], we are concerned here with a different time-dependent scenario where the population of nodes remains fixed from the outset, and the graph evolves through the appearance (birth) or the deletion (death) of edges.

2 Katz Centrality

To motivate our work we briefly discuss the case of a single, static network. Given a directed graph G defined over N nodes, we let A denote the corresponding N -by- N binary adjacency matrix, where a nonzero i, j entry records the presence of a link from node i to node j . We allow for $A_{ij} \neq A_{ji}$, so that the adjacency matrix may be unsymmetric.

Numerous measures have been proposed for quantifying important features of a network, many of them originating from the field of social network analysis [20]. We consider here the issue of assigning an importance ranking to each node, and focus on the idea of Katz [14], which has widely influenced the study of static social networks [21, 22, 23].

Although the definition can be derived by viewing a link as a ‘vote of confidence,’ in the style of Google’s PageRank algorithm [23, Chapter 7], we summarize here the original derivation of Katz. To quantify the propensity for node i to communicate, or interact, with node j , we may count how many walks of length $w = 1, 2, 3, \dots$ there are from i to j . We may then combine these counts into a single, cumulative total over all w . Allowing for the fact that shorter walks are generally more important (since, for example, the noise or cost of a transmission may increase with length), it makes sense to scale the counts according to the walk length. A particularly attractive choice is to scale walks of length w by a factor a^w , where a is a suitably

chosen scalar. A basic identity from graph theory shows that the k th power of the adjacency matrix has i, j element that counts the number of walks of length w from node i to node j . Introducing the identity matrix $I \in \mathbb{R}^{N \times N}$ for convenience, this leads us to the expansion $I + aA + a^2A^2 + a^3A^3 \dots$, which converges to the resolvent function $(I - aA)^{-1}$ when $a < 1/\rho(A)$. Here $\rho(\cdot)$ denotes the spectral radius; that is, the largest eigenvalue in modulus. Since $((I - aA)^{-1})_{ij}$ summarizes how well information can pass from node i to node j , the n th row sum

$$\sum_{k=1}^N ((I - aA)^{-1})_{nk} \quad (1)$$

is a centrality measure for node n . Following Newman [23] we refer to this as the *Katz centrality*.

We emphasize that this centrality measure is based on the combinatorics of *walks*, which allow nodes and edges to be reused during a traversal, rather than *paths* or *shortest paths*. A practical advantage of the walk counting approach is that the combinatorics can be conveniently described and implemented in terms of basic operations in linear algebra. Two further justifications are that (a) information does not necessarily flow along paths or geodesics [21, 22, 24] and (b) walk counting is more tolerant of errors (missing and spurious edges) than path counting. We also note that a related walk based measure of centrality was proposed for this static network case in [25], and the idea has been shown to lead to very powerful tools that are useful across a range of application areas [26, 27, 28, 29].

In the case of a directed network, the Katz centrality (1) quantifies the ability of node n to send out information along the directed links. By an entirely analogous argument, we could use a column sum instead of a row sum (or, equivalently, replace the adjacency matrix by its transpose) to give a measure that quantifies the ability of node n to acquire information:

$$\sum_{k=1}^N ((I - aA)^{-1})_{kn} . \quad (2)$$

3 Dynamic Centralities

We now return to our main theme of dynamic networks. To formalize our ideas, given a set of N nodes we consider an ordered sequence $\{G^{[k]}\}$ for

$k = 0, 1, 2, \dots, M$, where each $G^{[k]}$ is an unweighted graph defined over those nodes. We think of a corresponding ordered sequence of time points $t_0 \leq t_1 \leq \dots \leq t_M$, so that $G^{[k]}$ records the state of the network at time t_k . Each graph may then be represented by its adjacency matrix, $A^{[k]}$.

To address the question of how well information can be passed between pairs of nodes, we generalize the static graph concept of a walk as follows.

Definition 1 *A dynamic walk of length w from node i_1 to node i_{w+1} consists of a sequence of edges $i_1 \rightarrow i_2, i_2 \rightarrow i_3, \dots, i_w \rightarrow i_{w+1}$ and a non-decreasing sequence of times $t_{r_1} \leq t_{r_2} \leq \dots \leq t_{r_w}$ such that $A_{i_m, i_{m+1}}^{[r_m]} \neq 0$. We also define the lifetime of this walk to be $t_{r_w} - t_{r_1}$.*

We note that an analogous definition of a *dynamic path* can be made by insisting that no node is visited more than once—that concept was developed recently in [4], and, as described in section 1, time-respecting paths for different temporal models have appeared in earlier work [7, 16, 15, 17].

We emphasize that the sequence of times $t_{r_1}, t_{r_2}, \dots, t_{r_w}$ in Definition 1 must be nondecreasing, in order to respect the arrow of time, but

- repeated times are allowed: for example, if $r_1 < r_2 = r_3 < r_4$ then precisely two edges are followed at time t_{r_2} ,
- times are not required to be consecutive: for example, if $r_2 > r_1 + 1$ then the networks corresponding to times in between t_{r_1} and t_{r_2} have not been used during the walk.

Of course, depending on the application area, it may be reasonable to alter these features; forcing at most one edge per time level and/or forcing time levels to be consecutive. The ideas presented here could be adjusted accordingly.

Our key observation, which generalizes the static walk-counting identity mentioned in section 2, is that the matrix product $A^{[r_1]} A^{[r_2]} \dots A^{[r_w]}$ has i, j element that counts the number of dynamic walks of length w from node i to node j on which the m th step of the walk takes place at time t_{r_m} .

Now, in this new dynamic setting, we may apply the arguments that were used to derive the Katz centrality measure (1). We wish to quantify the propensity for node i to communicate, or interact, with node j . For each length $w = 1, 2, \dots$, we may count the number of dynamic walks from i to j ,

downweighting walks of length w by a factor a^w . In the matrix multiplication framework, this leads to the task of summing all products of the form

$$a^w A^{[r_1]} A^{[r_2]} \dots A^{[r_w]}, \quad \text{where } r_1 \leq r_2 \leq \dots \leq r_w. \quad (3)$$

These arguments motivate the matrix product

$$\mathcal{Q} := (I - aA^{[0]})^{-1} (I - aA^{[1]})^{-1} \dots (I - aA^{[M]})^{-1}. \quad (4)$$

The use of the identity matrices in (4) is crucial in our target case of large, sparse networks—it allows a message to ‘wait’ at a node until a suitable connection appears at a later time.

Overall, as required, the matrix \mathcal{Q} records the sum of all terms of the form (3). We may therefore use \mathcal{Q}_{ij} as our summary of how well information can be passed from node i to node j . The n th row and column sums

$$C_n^{\text{broadcast}} := \sum_{k=1}^N \mathcal{Q}_{nk} \quad \text{and} \quad C_n^{\text{receive}} := \sum_{k=1}^N \mathcal{Q}_{kn} \quad (5)$$

are centrality measures that quantify how effectively node n can *broadcast* and *receive* messages, respectively¹.

Because we are interested in the relative values of the centrality measures across all nodes, rather than their absolute sizes, we are free to multiply \mathcal{Q} by any positive scalar. We can use this freedom to avoid under or overflow in the computations. A normalized version, say $\hat{\mathcal{Q}}$, could be computed as the result $\hat{\mathcal{Q}}^{[M]}$ of an iteration such as

$$\hat{\mathcal{Q}}^{[k]} = \frac{\hat{\mathcal{Q}}^{[k-1]} (I - aA^{[k]})^{-1}}{\|\hat{\mathcal{Q}}^{[k-1]} (I - aA^{[k]})^{-1}\|}, \quad k = 0, 1, 2, \dots, M,$$

with $\hat{\mathcal{Q}}^{[-1]} = I$, where $\|\cdot\|$ denotes any convenient matrix norm. In our computations we use the Euclidean norm.

These new centralities are a direct generalization of Katz centralities (1) and (2) to the case of more than one time point.

Two features of this new approach are immediately apparent.

¹An alternative is to specify $k \neq n$ in the summations, so that closed walks are not included. However, such closed walks can play an important role as indicators of centrality [30].

- The basic computational tasks are linear system solves, which are convenient and efficient for large, sparse networks.
- The inherent asymmetry caused by the dynamics is captured directly through the non-commutativity of matrix multiplication.

To help understand the role of the downweighting parameter, a , we first note that for a fixed collection of network data, in the limit $a \rightarrow 0$ the centrality measures reduce to multiples of the aggregate out and in degrees, shifted by unity;

$$\lim_{a \rightarrow 0^+} \frac{C_n^{\text{broadcast}} - 1}{a} = \sum_{k=1}^N \left(\sum_{p=0}^M A^{[p]} \right)_{nk},$$

$$\lim_{a \rightarrow 0^+} \frac{C_n^{\text{receive}} - 1}{a} = \sum_{k=1}^N \left(\sum_{p=0}^M A^{[p]} \right)_{kn}.$$

At the other extreme, to guarantee that each resolvent $(I - aA^{[s]})^{-1}$ in (4) exists, we require that $a < 1/\rho(A^{[s]})$, for all s . Furthermore, choosing a close to $1/\max_s \rho(A^{[s]})$ will cause the corresponding time t_s to dominate the overall communicability matrix \mathcal{Q} . In practice, a suitable choice of a would be sufficiently below $1/\max_s \rho(A^{[s]})$ that the results are not sensitive to small changes in a and sufficiently above zero that they do not collapse to the shifted aggregate out and in degrees.

Let us consider how the asymmetry of \mathcal{Q} arises, and hence that between the centrality measures, in the case of an evolving undirected graph. Here, all the adjacency matrices are symmetric. For any N -by- N matrix, B , we define $\mathcal{S}(B) := \frac{1}{2}(B + B^T)$ and $\mathcal{AS}(B) := \frac{1}{2}(B - B^T)$ to be the projections of B onto the space of symmetric matrices and the orthogonal space of anti-symmetric matrices, respectively. The anti-symmetric part of \mathcal{Q} governs the differences between the column and row sums of \mathcal{Q} , since

$$2\mathcal{AS}(\mathcal{Q})\mathbf{1} = \mathbf{C}^{\text{broadcast}} - \mathbf{C}^{\text{receive}},$$

where $\mathbf{1} = (1, 1, \dots, 1)^T$, and $\mathbf{C}^{\text{broadcast}} = (C_1^{\text{broadcast}}, \dots, C_N^{\text{broadcast}})^T$ are N -vectors, with $\mathbf{C}^{\text{receive}}$ defined analogously.

Working with the non-normalized version of \mathcal{Q} in (4), we have

$$\mathcal{Q} = I + a \sum_{p=0}^M A^{[p]} + a^2 \sum_{p=0}^M \sum_{p'=p}^M A^{[p]} A^{[p']} + O(a^3).$$

It follows that

$$\mathcal{S}(\mathcal{Q}) = I + a \sum_{p=0}^M A^{[p]} + O(a^2), \quad (6)$$

and

$$2\mathcal{AS}(\mathcal{Q}) = a^2 \sum_{p=0}^M \sum_{p'=p+1}^M [A^{[p]}, A^{[p']}] + O(a^3), \quad (7)$$

where $[A, B] := AB - BA$ denotes the commutator of matrices A and B . Since each separate graph is undirected, this shows that the leading anti-symmetric terms arise only from interactions over distinct pair of time steps.

We emphasize that this work focuses on the case of a fixed set of data. In some applications, $A^{[k]}$ will itself be an aggregate of activity over a time window; for example, in sections 4.2 and 4.3 we consider daily telephone and email communication. If we reduce the time window down to hours, minutes, seconds, \dots , then, before we can think about the asymptotic limit, we must first be clear about the nature of the data.

If we consider email communication, and assume that messages are passed instantaneously, then in the limit of time resolution the communicability matrix \mathcal{Q} would not change—we reach a point where no further walks can be created through refinement. In the linear algebra setting of (4) only empty networks are added to the sequence and no new non-identity factors will arise. In this ultra-high-frequency regime, if the edges are undirected it would be natural to replace (4) with the “at most one link per time window” version

$$\mathcal{Q} := (I + aA^{[0]}) (I + aA^{[1]}) \cdots (I + aA^{[M]}),$$

in order to eliminate trivial closed walks such as $i \mapsto j \mapsto i$ that occur over a single window. However, the idea that the finest time level gives the most accurate picture must be treated with caution—in the email context the order in which messages are read or acted upon does not necessarily reflect the order of arrival.

By contrast, suppose that dynamic edges represent telephone communications with prescribed start and finish times. As we refine the time window in this setting beyond the point where new information is added, we simply start to collect contiguous repeats of a finite number of adjacency matrices. In this limit, suppose we allow the downweighting parameter a to scale inversely with the length of the time window. Then letting A denote the adjacency matrix representing the connectivity pattern that remains constant

over some period (because no calls are started or finished over this time), we will be computing as a factor in (4) a quantity of the form

$$\lim_{K \rightarrow \infty} \left(I - \frac{a}{K} A \right)^{-K} = \exp(aA).$$

Measures based on the scaled matrix exponential have been studied extensively, see, for example, [31], and we see that, in this asymptotic limit, the dynamic communicability matrix \mathcal{Q} is formed as a product of such factors.

We also mention that the definition (4) extends readily to the case where a varies with the time point; this might be natural, for example, if non-uniform time windows are used or if some external property, such as the cost of making a phone call or the likelihood of a batch of spam email, varies over time.

4 Computational Tests

In this section we describe some illustrative computations with the new dynamic centralities. For convenience, we let amax denote the upper limit $1/\max_s \rho(A^{[s]})$ for the downweighting parameter a . In order to compare results with those for the Katz centralities we let A^{star} denote the binarized version of the aggregated adjacency matrix, so $(A^*)_{ij} = 1$ if $(A^{[k]})_{ij} = 1$ for any k , and $(A^*)_{ij} = 0$ otherwise. We then let amax^* denote the upper limit $1/\rho(A^*)$.

4.1 Synthetic Data

Figure 2 shows a proof-of-principle test of the ideas behind (4). Here we used $N = 1001$ nodes and simulated networks at 31 time points; that is, a month of daily data. At each time point, for nodes 1 to 1000 we constructed, independently, a classical, undirected, Erdos/Renyi random graph—each was chosen uniformly from the collection of all graphs with 1000 nodes and 1000 edges. Then at each time the final node 1001 was connected to the two nodes with largest degree. In this way, node number 1001 is distinguished only by the time-sensitive ‘quality’ of its links—at each time t_k it has a degree that matches that of the average node, and it will never be among the highest degree nodes at any time; so any static or aggregative measure is likely to fail to identify this node as being special. The upper pictures in the figure

scatter the (normalized) broadcast and receive centralities (5) for each node, with node 1001 identified by a circle. In this case, $\text{amax} = 0.26$, and we show results for $a = 0.2$ (left) and $a = 0.1$ (right). We see that the new measures correctly identify the fact that this node can communicate well, despite never enjoying a high degree. Because each network is undirected, A^* is symmetric, and the Katz centralities (1) and (2) are equivalent. Here, $\text{amax}^* = 0.016$ and in the lower pictures we give a histogram for the Katz centrality in the case $a = 0.01$ (left) and $a = 0.005$ (right). The centrality for node 1001 is marked with a vertical dashed line, and we see that its advantageous connectivity at each time point has been lost in the aggregation process.

4.2 Telecommunication Data

We now consider telecommunication data from [32]. We have daily “who phoned who” information between 106 individuals based at M.I.T. over 365 days, with starting date 20th July 2004. Because phone conversations are bi-directional, we have symmetrized the data, so $A_{ij}^{[k]} = 1$ if individuals i and j had at least one interaction on day k . Figure 3 shows a summary of the adjacency matrices aggregated into 28 day intervals (day 365 omitted). We notice a decrease in activity outside the traditional academic teaching periods.

The upper left picture in Figure 4 shows the daily edge count. For this data $\text{amax} = 0.12$, and the figure shows centrality results for $a = 0.1$. In the upper right picture we scatter plot on a log-log scale the broadcast and receive centralities (5). Here, and in all other scatter plots, the correlation coefficient “corr” and the Kendall τ index “ τ ” for a pair of raw (not log transformed) centralities are quoted to two decimal places in the figure caption. We see that even though the individual adjacency matrices are symmetric, there is no strong correlation between the two centralities. The lower pictures scatter plot the broadcast and receive centralities for each node against the total degree; that is, the sum of the node’s degrees over all days. This makes it clear that the new centralities are not simply repeating the degree information. The figure captions also quantify the overlap between the sets of nodes ranked among the top twenty. In Figure 4, the top twenty nodes ordered from twentieth place to first place in terms of the three measures are

broadcast : 27, 32, 38, 44, 47, 7, 45, 6, 2, 4, 10, 3, 30, 49, 26, 1, 46, 8, 5, 102
receive : 94, 58, 76, 95, 15, 20, 12, 89, 93, 30, 19, 49, 6, 35, 39, 52, 42, 8, 13, 53,

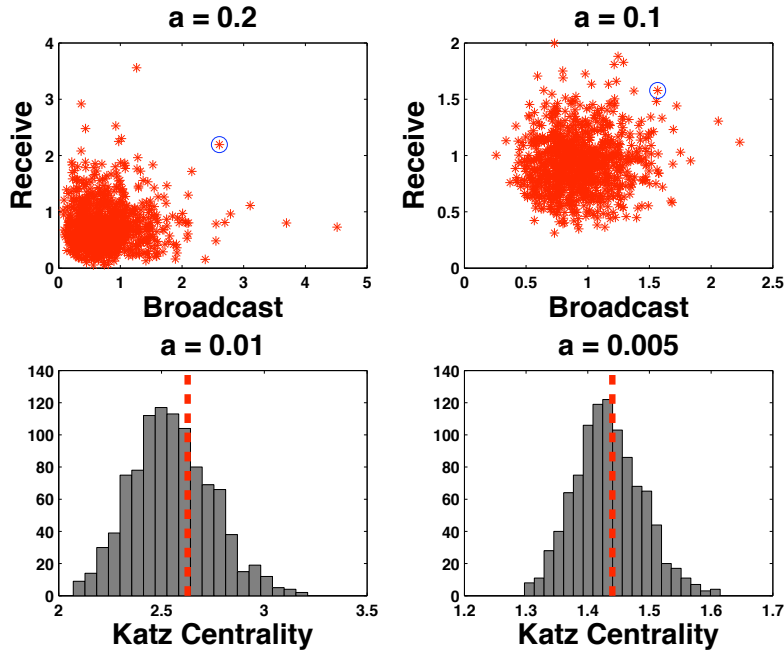


Figure 2: (Color online) Comparison of centralities synthetically generated network of 1001 nodes, where node number 1001 is designed to have average activity, as measured by the aggregate degree, but enjoys high quality connections at each time point. Upper pictures scatter the broadcast and receive centralities (5) with node 1001 circled. Lower pictures show histograms of Katz centrality for the binarized aggregate network, with centrality for node 1001 marked with a vertical dashed line.

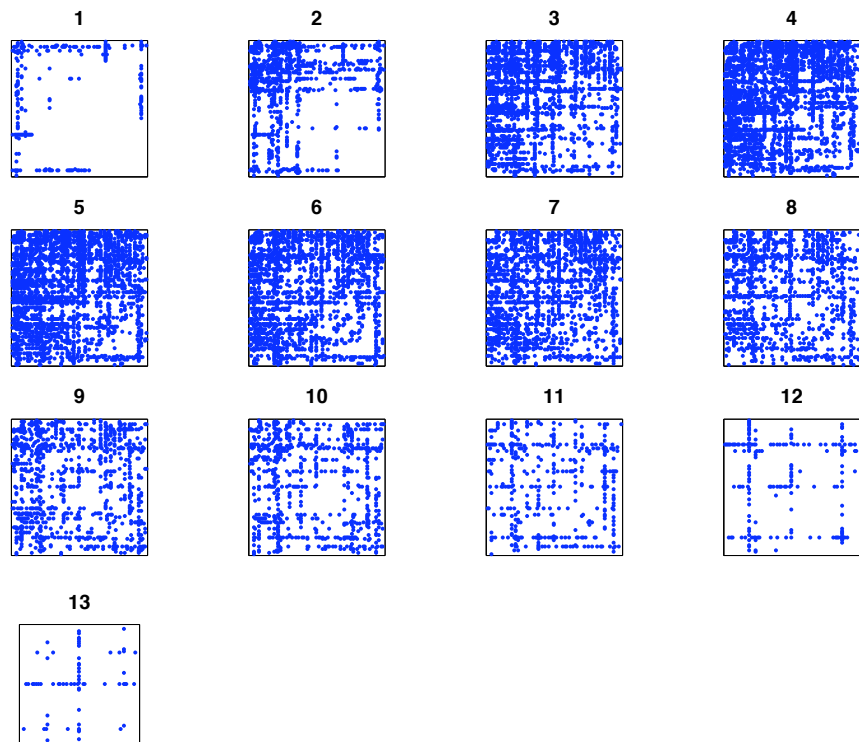


Figure 3: (Color online) Adjacency matrices for the M.I.T. telecommunication data, symmetrized and aggregated into 13 sets of 28 day windows. Each plot shows the nonzero pattern: a dot in row i and column j represents at least one telephone interaction between those individuals.

totaldegree : 21, 9, 93, 100, 32, 10, 57, 22, 49, 25, 53, 23, 6, 40, 20, 3, 2, 4, 8, 5.

In this case the overlaps between broadcast & receive, broadcast & total degree and receive & total degree contain 4, 9 and 6 nodes, respectively. Only one node appears in all three top twenty lists.

Figure 5 examines the sensitivity of the results to the parameter a . The upper pictures show how the centralities change from $a = 0.1$ to $a = 0.05$. The top twenty broadcast lists have 14 nodes in common and for the receiver lists the overlap is 16. The lower pictures show the change from $a = 0.05$ to $a = 0.01$, and in this case the top twenty overlap counts are 16 and 11. Overall, the experiments indicate that the two new measures deliver distinct information that is different from a raw degree count, and remains consistent over a range of a values.

We also found that neither of the dynamic centralities were strongly correlated with the Katz centrality for the binary aggregated matrix A^* . In this case $\text{amax}^* = 0.02$ and comparing the the $a = 0.1$ broadcast centralities with the $a = 0.015$ Katz centralities gave $\text{corr} = 0.34$ and $\tau = 0.25$. Similarly, for receive versus Katz we had $\text{corr} = 0.22$ and $\tau = 0.24$.

4.3 Email Data

We now consider a public domain data set concerning email activities of Enron employees. In [33] the static, aggregate network was analysed, but here we treat it as an evolving network. We constructed daily information representing emails between 151 Enron employees, including `to`, `cc` or `bcc`. So $A_{ij}^{[k]} = 1$ if employee i sent at least one message to employee j on day k , but because this type of communication is unidirectional, we do not automatically add the $j \mapsto i$ link. We have data over 1138 days, starting on 11th May 1999. Many of the adjacency matrices are empty, stressing the importance of the identity matrices in (4) for analysing sparse data. The upper left plot in Figure 6 shows the daily edge count.

Using $a = 0.2$, in the upper right of Figure 6 we scatter plot broadcast versus receive centralities, and in the lower plots we show broadcast versus total out degree and receive versus total in degree. In this case $\text{amax} = 0.24$. As in the previous test, we see that the two new centrality measures are distinct; in particular, only two nodes appear in the overlap of top twenty broadcast and receive and it is clear that some top receivers are very poor broadcasters. The top twenty overlap between broadcast and total out degree

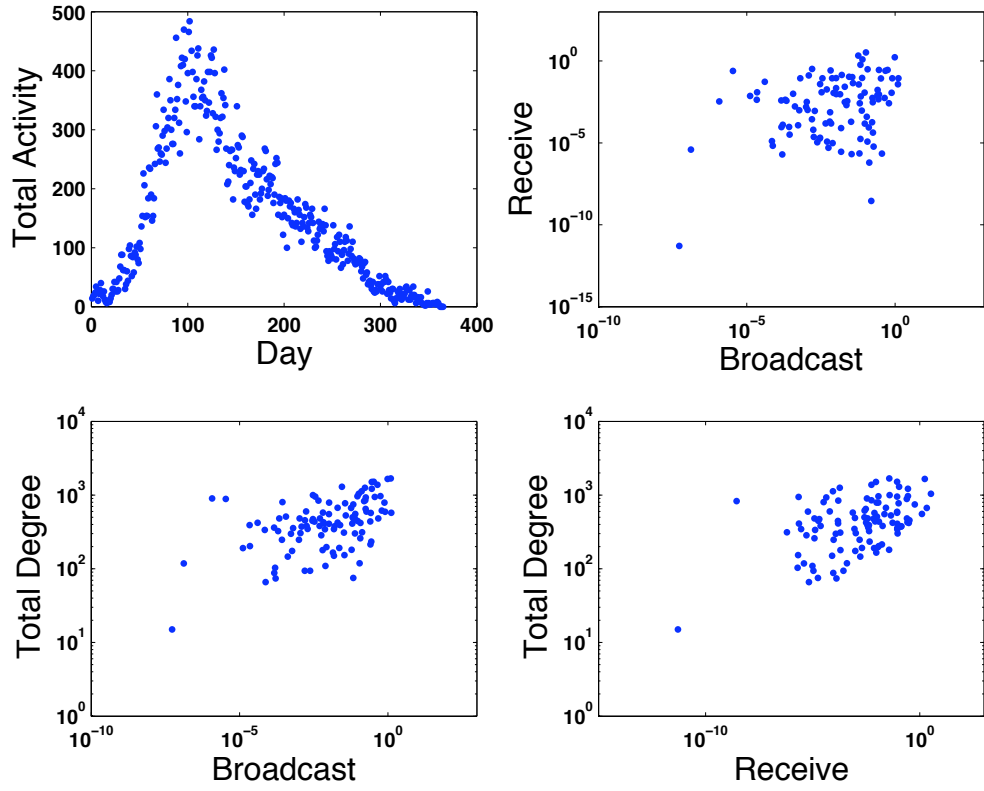


Figure 4: (Color online) Daily M.I.T. telecommunication data. Upper right: total activity per day. Upper left: Broadcast versus receive centrality; $\text{corr} = 0.14$, $\tau = 0.15$, top twenty overlap size 4. Lower left: broadcast centrality versus total degree; $\text{corr} = 0.50$, $\tau = 0.34$, top twenty overlap size 9. Lower right: Receive centrality versus total degree; $\text{corr} = 0.28$, $\tau = 0.28$, top twenty overlap size 6.

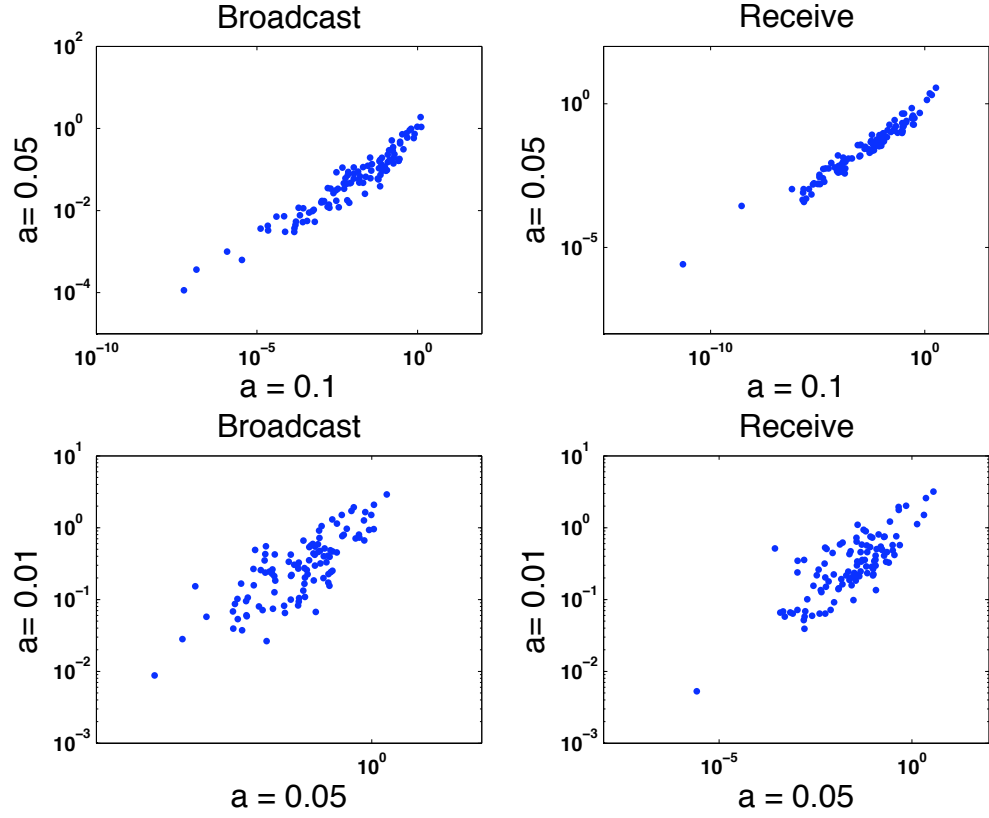


Figure 5: (Color online) Daily M.I.T. telecommunication data. Upper: $a = 0.1$ versus $a = 0.05$; for broadcast $\text{corr} = 0.93$ and $\tau = 0.82$, for receive $\text{corr} = 0.98$ and $\tau = 0.87$. Respective top twenty overlap sizes are 14 and 16. Lower: $a = 0.05$ versus $a = 0.01$; for broadcast $\text{corr} = 0.82$ and $\tau = 0.57$, for receive, $\text{corr} = 0.81$ and $\tau = 0.54$. Respective top twenty overlap sizes are 16 and 11.

is 11 and between receive and total in degree in 6, showing that the new measures do not simply reflect aggregate connectivity.

For the binarized aggregate network we have $\text{amax}^* = 0.05$ and comparing the $a = 0.2$ broadcast centralities with the $a = 0.04$ Katz centralities gave $\text{corr} = 0.10$ and $\tau = 0.37$. Similarly, for receive versus Katz we had $\text{corr} = 0.02$ and $\tau = 0.21$.

The upper plots of Figure 7 show how the new centralities change when a is reduced from 0.2 to 0.1, indicating robustness in this parameter regime. The lower plots show the effect of symmetrizing the data, so that $j \mapsto i$ whenever $i \mapsto j$, in the case $a = 0.1$. We then have $\text{amax} = 0.12$. We see that the new dynamic centralities are relatively insensitive to this transformation of the data, suggesting that the dominant asymmetry is caused by time.

5 Discussion

The new centrality measures introduced here can be computed at any point in time, and hence they may be used to monitor network behaviour dynamically. A practical problem with evolving networks is that of an observer who may be able make some kind of intervention; for example by injecting some information (marketing content, rumours, propaganda, misinformation) at key nodes at some instant, or by isolating or even removing a node. This raises the issue of predicting future network behavior. We will briefly discuss an approach based on the observer's expectation of the future communicability: an estimate of \mathcal{Q} going forwards.

Suppose we have a stochastic model for the evolution of the network based on historical data and some specific knowledge. More precisely, suppose we have $P(A^{[p+1]}|H_p)$, the conditional distribution for the adjacency matrix at the next time step given its entire history up to and including step p , so $H_p = \{A^{[p]}, A^{[p-1]}, A^{[p-2]}, \dots\}$. Then applying this model iteratively we obtain the conditional distribution for $A^{[p']}$ for any $p' > p$; that is, $P(A^{[p']}|H_p)$.

Let us write $E(A^{[p']}|H_p)$ to denote the corresponding expected value of the future adjacency matrix, given H_p .

Now suppose we have observed the network up to and including some time step, say $p = 0$ for convenience. Then from (6) and (7) we can calculate estimates for the expectation of the communicability over the current and

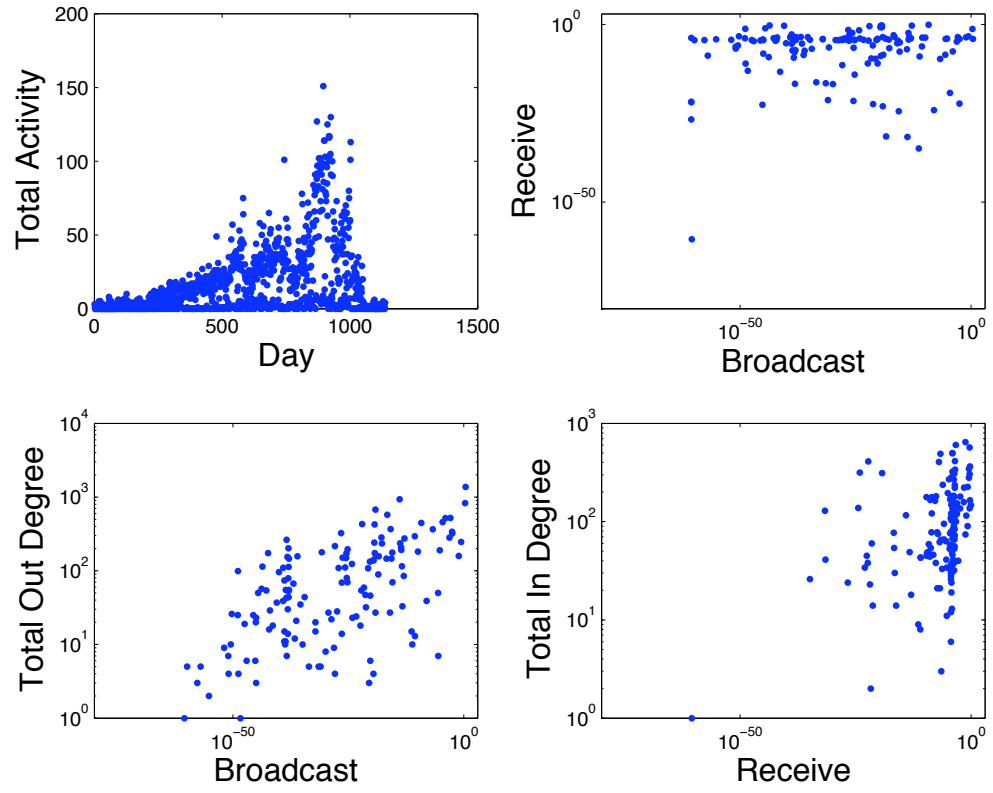


Figure 6: (Color online) Results for Enron email data. Upper left: total number of edges per day. Upper right: Scatter plot of broadcast and receive centralities; $\text{corr} = 0.00$, $\tau = 0.05$, top twenty overlap size 2. Lower left: Scatter plot of broadcast centrality and total out degree; $\text{corr} = 0.62$, $\tau = 0.46$, top twenty overlap size 11. Lower right: Scatter plot of receive centrality and total in degree; $\text{corr} = 0.28$, $\tau = 0.31$, top twenty overlap size 6.

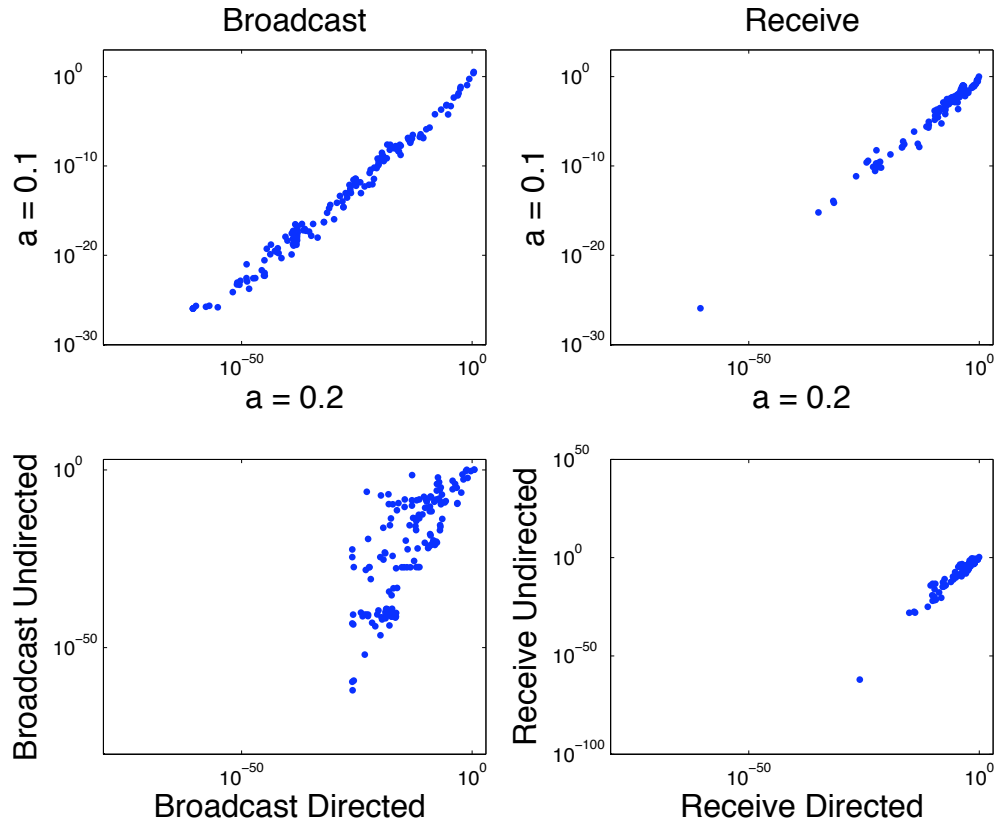


Figure 7: (Color online) Results for Enron email data. Upper left: Scatter plot of broadcast centralities for $a = 0.2$ and $a = 0.1$; $\text{corr} = 1.00$, $\tau = 0.93$, top twenty overlap size 18. Upper right: Scatter plot of receive centralities for $a = 0.2$ and $a = 0.1$; $\text{corr} = 0.97$, $\tau = 0.88$, top twenty overlap size 14. Lower left: Scatter plot of broadcast centralities for directed and undirected networks for $a = 0.1$; $\text{corr} = 0.82$, $\tau = 0.58$, top twenty overlap size 12. Lower right: Scatter plot of receive centralities for directed and undirected networks for $a = 0.1$; $\text{corr} = 0.76$, $\tau = 0.81$, top twenty overlap size 13.

future time steps. We have the small a approximations

$$E(\mathcal{S}(\mathcal{Q})|H_0) = I + a \sum_{p=0}^M E(A^{[p]}|H_0) + O(a^2),$$

and

$$2E(\mathcal{AS}(\mathcal{Q})|H_0) = a^2 \sum_{p=0}^M \sum_{p'=p+1}^M E([A^{[p]}, A^{[p']}])|H_0) + O(a^3).$$

These estimates may be accessible in practice, depending on the complexity and memory dependence of the model. For example, suppose we make the dramatically simplifying assumption that our model is a symmetric, edge independent Markov process. Letting α_{ij} and ω_{ij} denote the stepwise birth and death rates for the evolution of the (i, j) th edge, we have $A^{[p]} \rightarrow A^{[\infty]}$ as $p \rightarrow \infty$, where $A_{ij}^{[\infty]} = \alpha_{ij}/(\alpha_{ij} + \omega_{ij})$. In this Markovian case we can also replace the history, H_p , by the single previous step $A^{[p]}$. Then considering time steps 0 up to M we have

$$E(\mathcal{Q}|A^{[0]}) = I + a (R_M \circ (A^{[0]} - A^{[\infty]}) + (M + 1)A^{[\infty]}) + O(a^2),$$

where R_M is the symmetric matrix given by $(R_p)_{ij} = (1 - (1 - \alpha_{ij} + \omega_{ij})^{M+1})/(\alpha_{ij} + \omega_{ij})$, and \circ denotes componentwise multiplication. This quantifies the relative contributions to \mathcal{Q} made by the initial condition and the long term expected *equilibrium* value for each edge. So if the observer wishes to intervene based on the dominance of some large row or column sums of \mathcal{Q} , we can see that this may require such action sooner or later depending on the current state of the network and the longer term expectation.

So, overall, we believe that the new class of walk-based centrality measures introduced here offers great potential as a computationally and analytically attractive means to treat time-stamped network sequences, both for summarizing existing data sets and real-time actioning.

Acknowledgement This work was supported by the Engineering and Physical Sciences Research Council and the Research Councils UK Digital Economy programme.

References

- [1] E. Estrada, M. Fox, D. J. Higham, and G.-L. Oppo, editors. *Network Science: Complexity in Nature and Technology*. Springer, Berlin, to appear, 2010.
- [2] M. E. J. Newman. The structure and function of complex networks. *SIAM Review*, 45(2):167–256, 2003.
- [3] Rowland R. Kao. Networks and models with heterogeneous population structure in epidemiology. In Ernesto Estrada, Desmond J. Higham, Maria Fox, and Gian-Luca Oppo, editors, *Network Science: Complexity in Nature and Technology*, pages 51–84. Springer, 2010.
- [4] J. Tang, S. Scellato, M. Musolesi, C. Mascolo, and V. Latora. Small-world behavior in time-varying graphs. *Physical Review E*, 81:05510, 2010.
- [5] Jure Leskovec and Eric Horvitz. Planetary-scale views on a large instant-messaging network. In *Proceeding of the 17th international conference on World Wide Web*, pages 915–924, Beijing, China, 2008. ACM.
- [6] Aurelien Gautreau, Alain Barrat, and Marc Barthelemy. Microdynamics in stationary complex networks. *Proc. Nat. Acad. Sci.*, 106:8847–8852, 2009.
- [7] K. Berman. Vulnerability of scheduled networks and a generalization of Menger’s Theorem. *Networks*, 28:125–134, 1996.
- [8] L. McNamara, C. Mascolo and L. Capra. Media Sharing based on Colocation Prediction in Urban Transport. In *Proc. of ACM 14th International Conference on Mobile Computing and Networking (Mobicom08)*, pages 58–69, San Francisco, CA, September 2008.
- [9] John Tang, Mirco Musolesi, Cecilia Mascolo, and Vito Latora. Characterising temporal distance and reachability in mobile and online social networks. *SIGCOMM Comput. Commun. Rev.*, 40:118–124, January 2010.
- [10] Evrim Acar, Daniel M. Dunlavy, and Tamara G. Kolda. Link prediction on evolving data using matrix and tensor factorizations. In *LDMTA ’09*:

- Proceeding of the ICDM'09 Workshop on Large Scale Data Mining Theory and Applications*, pages 262–269. IEEE Computer Society Press, December 2009.
- [11] Peter Grindrod and Desmond J. Higham. Evolving graphs: Dynamical models, inverse problems and propagation. *Proceedings of the Royal Society, Series A*, 466:753–770, 2010.
 - [12] John Tang, Mirco Musolesi, Cecilia Mascolo, and Vito Latora. Temporal distance metrics for social network analysis. In *Proceedings of the 2nd ACM SIGCOMM Workshop on Online Social Networks (WOSN09)*, Barcelona, 2009.
 - [13] John Tang, Mirco Musolesi, Cecilia Mascolo, Vito Latora, and Vincenzo Nicosia. Analysing information flows and key mediators through temporal centrality metrics. In *SNS '10: Proceedings of the 3rd Workshop on Social Network Systems*, pages 1–6, New York, NY, USA, 2010. ACM.
 - [14] L. Katz. A new index derived from sociometric data analysis. *Psychometrika*, 18:39–43, 1953.
 - [15] David Kempe, Jon Kleinberg, and Amit Kumar. Connectivity and inference problems for temporal networks. *J. Comput. Syst. Sci.*, 64:820–842, 2002.
 - [16] Petter Holme. Network reachability of real-world contact sequences. *Physical Review E (Statistical, Nonlinear, and Soft Matter Physics)*, 71(4), 2005.
 - [17] Gueorgi Kossinets, Jon Kleinberg, and Duncan Watts. The structure of information pathways in a social communication network. In *Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '08, pages 435–443, New York, NY, USA, 2008. ACM.
 - [18] Peter J. Mucha, Thomas Richardson, Kevin Macon, Mason A. Porter, and Jukka-Pekka Onnela. Community structure in time-dependent, multiscale, and multiplex networks. *Science*, 328:876–878, 2010.
 - [19] A.-L. Barabási and R. Albert. Emergence of scaling in random networks. *Science*, 286(5439):509–12, 1999.

- [20] S. Wasserman and K. Faust. *Social Network Analysis: Methods and Applications*. Cambridge University Press, Cambridge, 1994.
- [21] S. P. Borgatti. Centrality and network flow. *Social Networks*, 27:55–71, 2005.
- [22] S. P. Borgatti and M. G. Everett. A graph theoretic perspective on centrality. *Social Networks*, 28:466–484, 2006.
- [23] M. E. J. Newman. *Networks an Introduction*. Oxford University Press, Oxford, 2010.
- [24] L. C. Freeman, S. P. Borgatti, and D. R. White. Centrality in valued graphs: A measure of betweenness based on network flow. *Social Networks*, 13:141–154, 1991.
- [25] Ernesto Estrada and Naomichi Hatano. Communicability in complex networks. *Phys. Rev. E*, 77(3):036111, Mar 2008.
- [26] J. J. Crofts and D. J. Higham. A weighted communicability measure applied to complex brain networks. *J. Royal Society Interface*, 33, 2009.
- [27] J. J. Crofts and D. J. Higham. Googling the brain: Discovering hierarchical and asymmetric network structures, with applications in neuroscience. Technical Report Mathematics and Statistics Research Report 35, University of Strathclyde, 2010.
- [28] Ernesto Estrada and Desmond J. Higham. Network properties revealed through matrix functions. *SIAM Review*, 52:696–671, 2010.
- [29] E. Estrada, D. J. Higham, and N. Hatano. Communicability betweenness in complex networks. *Physica A*, 388:764–774, 2009.
- [30] E. Estrada and J. A. Rodríguez-Velázquez. Subgraph centrality in complex networks. *Phys. Rev. E*, 71:056103, 2005.
- [31] E. Estrada and N. Hatano. Statistical-mechanical approach to subgraph centrality in complex networks. *Chemical Physics Letters*, 439:247–251, 2007.

- [32] Nathan Eagle, Alex S. Pentland, and David Lazer. Inferring friendship network structure by using mobile phone data. *Proceedings of the National Academy of Sciences*, 106(36):15274–15278, September 2009.
- [33] Anurat Chapanond, Mukkai Krishnamoorthy, and Bülent Yener. Graph theoretic and spectral analysis of Enron email data. *Computational & Mathematical Organization Theory*, 11:265–281, 2005.