

## Local Structure of Directed Networks

Ginestra Bianconi,<sup>1</sup> Natali Gulbahce,<sup>2,3</sup> and Adilson E. Motter<sup>4</sup>

<sup>1</sup>*The Abdus Salam International Center for Theoretical Physics, Strada Costiera 11, 34014 Trieste, Italy*

<sup>2</sup>*Theoretical Division and CNLS, Los Alamos National Laboratory, New Mexico 87545, USA*

<sup>3</sup>*Center for Complex Networks Research and Department of Physics, Northeastern University, Boston, Massachusetts 02115, USA*

<sup>4</sup>*Department of Physics and Astronomy and NICO, Northwestern University, Illinois 60208, USA*

(Received 12 July 2007; published 20 March 2008)

Previous work on undirected small-world networks established the paradigm that locally structured networks tend to have a high density of short loops. On the other hand, many realistic networks are *directed*. Here we investigate the local organization of directed networks and find, surprisingly, that real networks often have very few short loops as compared to random models. We develop a theory and derive conditions for determining if a given network has more or less loops than its randomized counterparts. These findings carry broad implications for structural and dynamical processes sustained by directed networks.

DOI: 10.1103/PhysRevLett.100.118701

PACS numbers: 89.75.Hc, 89.75.Da, 89.75.Fb

Asymmetric interactions are widespread in natural and technological networks, particularly when the network transports a flow or underlies collective behavior [1]. The structure of such directed networks can be characterized by the statistics of loops, the building blocks of closed paths, which provides information on structural correlations [2], motifs, robustness and redundancy of pathways, and impacts dynamical as well as equilibrium critical phenomena on the network [3].

In undirected networks, the large number of short loops together with small diameter gives rise to the small-world effect encountered in many real systems [4]. Strikingly, in this Letter we show that there is a large class of directed networks for which the number of loops is strongly reduced with respect to the random hypothesis. The directed neural network of *C. elegans*, for example, has less than 50% of the short loops expected from a random ensemble with the same degree sequence, despite the well-known fact that, when regarded as an undirected network [4], it has a clustering coefficient 5.6 times larger than randomly rewired versions of the network.

Motivated by this empirical finding, we demonstrate numerically and analytically that degree correlations [5] strongly constrain the loop structure of directed networks. Moreover, we go beyond the degree-correlated picture and derive conditions for determining if a *given* network has more or less loops than its randomized counterparts. We characterize the network local organization in terms of short loops and its global organization in terms of long loops. We compare our analytical results with exact (when possible) or approximate numerical calculation of the number of loops in a class of directed networks that includes foodweb, power-grid, metabolic, neural, transcription, and WWW networks. Our findings that many directed networks are underlooped may have broad implications given that such networks exhibit, for example, improved stability in foodweb systems [8] and enhanced

synchronization [9] and transportation properties in various other systems [10].

*Short loops in random networks.*—We first derive the expected number of self-avoiding loops in directed random networks. The general way to construct random uncorrelated undirected networks is by means of the Molloy-Reed model. Given a set of nodes  $V = \{i:1, \dots, N\}$ , the construction is based on generating a sequence of degrees  $\{k^i\}$  from a given degree distribution  $P(k)$  with a structural cutoff  $K = \mathcal{O}(N^{1/2})$  [11], and randomly connecting the links. In this ensemble, the expected number  $\mathcal{N}_L$  of short loops of length  $L$  is given by [12,13]

$$E_{\text{undir}}(\mathcal{N}_L) = \frac{1}{2L} \left( \frac{\langle k(k-1) \rangle}{\langle k \rangle} \right)^L. \quad (1)$$

This formula implies that a network with diverging  $\langle k^2 \rangle$  has many more short loops than networks with finite  $\langle k^2 \rangle$ . In particular, scale-free networks with scaling exponent  $\gamma \leq 3$  have many short loops while Erdős-Rényi networks have a negligible number of short loops in the  $N \rightarrow \infty$  limit. We now show that this expression can be generalized to random *directed* networks. We again consider the Molloy-Reed construction but in this case we draw a sequence of incoming and outgoing links  $\{(k_{\text{in}}^i, k_{\text{out}}^i)\}$  from a degree distribution  $P(k_{\text{in}}, k_{\text{out}})$  for all nodes  $i \in V$ . This distribution, which is not factorisable in general, describes correlated variables  $k_{\text{in}}$  and  $k_{\text{out}}$  at any given node. For directed uncorrelated networks, the structural cutoffs for in- and out-degrees satisfy  $K_{\text{in}}K_{\text{out}} < \langle k_{\text{in}} \rangle N$ . Proceeding as in the undirected case [12], we obtain that the expected number of loops of size  $L$  in the directed network ensemble is given by

$$E_{\text{dir}}(\mathcal{N}_L) = \frac{1}{L} \left( \frac{\langle k_{\text{in}} k_{\text{out}} \rangle}{\langle k_{\text{in}} \rangle} \right)^L, \quad (2)$$

where this approximate expression is valid for large  $N$  and

loop length satisfying  $L \ll N \langle k_{\text{in}} k_{\text{out}} \rangle^2 / \langle (k_{\text{in}} k_{\text{out}})^2 \rangle$ . For undirected networks,  $E_{\text{dir}}(\mathcal{N}_L)$  reduces to  $E_{\text{undir}}(\mathcal{N}_L)$  because the incoming connectivity is  $k$  and the outgoing connectivity at the end point of a link (on a self-avoiding loop) is  $k - 1$ . The only difference is a factor 2, which accounts for the orientation on the loops in Eq. (2).

We observe from Eq. (2) that, in directed networks, the one-point correlation between the number of incoming and outgoing links modulates the expected number of short loops. Indeed, if  $k_{\text{in}}$  and  $k_{\text{out}}$  on the same nodes are not correlated, then the number of short loops is strongly reduced as compared to the case when  $k_{\text{in}}$  and  $k_{\text{out}}$  are positively correlated. The Barabási-Albert (BA) networks [14], for example, have small degree correlations and are within the scope of  $E_{\text{dir}}(\mathcal{N}_L)$  and  $E_{\text{undir}}(\mathcal{N}_L)$  for uncorrelated random networks [15]. If we consider the undirected BA model, we find that the networks have many short loops compared to random Erdős-Rényi networks (in fact  $\langle k(k-1) \rangle \sim \log(N)$ ) [16]. In contrast, if we consider the directed version of the BA model (in which the incoming links are linked preferentially, and hence  $\langle k_{\text{in}} k_{\text{out}} \rangle = \langle k_{\text{in}} \rangle \langle k_{\text{out}} \rangle$ ), the networks have a negligible number of short loops just as the Erdős-Rényi networks in the  $N \rightarrow \infty$  limit.

*Short loops in a given network.*—A different approach is needed for counting the loops of a *specific* directed network, as required in the study of real systems. In this case, as in the case of undirected networks [16], the number of short loops can be expressed in terms of powers of the adjacency matrix. In particular, the number of (self-avoiding) loops of length  $L$  can be expressed as the total number of closed paths of length  $L$ , i.e.,  $\text{Tr} A^L / L$ , minus the closed paths of length  $L$  composed of self-intersecting loops. The number of loops of length  $L$  in a network with adjacency matrix  $A$  is then given by  $\mathcal{N}_L = \frac{1}{L} \sum_{\{L_\ell\}} c(\{L_\ell\}) \delta(L - \sum_\ell L_\ell) \sum_i \prod_\ell (A^{L_\ell})_{ii}$ , where the sequence  $\{L_\ell\}$  describes the loop composition of the paths for every correction term (for example, in the case  $L = 5$  we will find a correction term involving paths composed of  $\{L_\ell\} = \{2, 3\}$  directed loops). The coefficients  $c(\{L_\ell\})$  remain small for small  $L$ .

Starting from this general formula we derive upper and lower bounds for the number of loops in a given directed network. The upper bound is simply given by the sum of all closed paths of length  $L$ , i.e.  $\mathcal{N}_L \leq \frac{1}{L} \text{Tr} A^L = \frac{1}{L} \sum_n \lambda_n^L$ , where the sum is performed over all the eigenvalues (including multiplicities). To find a lower bound we have to express  $\mathcal{N}_L$  in terms of the eigenvalues of the adjacency matrix  $A$  and in terms of its Jordan basis. In this way, it follows that  $\mathcal{N}_L \simeq \text{Tr} A^L / L$  provided that  $\kappa_L \equiv \max_i \sum_j \sum_m \binom{L}{m} |\lambda_j^{-m} P_{ij} P_{j+m,i}^{-1}| \ll 1$ , where  $P$  is the matrix of generalized eigenvectors of  $A$  in the Jordan decomposition  $A = PJP^{-1}$  [9] and  $\sum'_m$  indicates a sum over the dimension of each Jordan block with associated eigenvalue  $\lambda_j$ , under the constraint that indices  $j$  and  $j + m$  are in the same block. If  $\kappa_L \ll 1$ , the dominant term in the expansion

of  $\mathcal{N}_L$  is the one with  $\{L_\ell\} = \{L\}$  and we have  $\mathcal{N}_L \simeq \frac{1}{L} \sum_n \lambda_n^L$ .

Comparing these results with the result found for the random case in Eq. (2), it follows that a sufficient condition for a specific network to have less short loops of length  $L$  than its randomized versions is  $\sum_n \lambda_n^L < (\langle k_{\text{in}} k_{\text{out}} \rangle / \langle k_{\text{in}} \rangle)^L$ . Conversely, if  $\kappa_L \ll 1$ , a condition for the network to have more loops is  $\sum_n \lambda_n^L > (\langle k_{\text{in}} k_{\text{out}} \rangle / \langle k_{\text{in}} \rangle)^L$ . For loops in a certain range of values  $L \in (1, L_c)$ , it is convenient to restate these conditions as

$$\bar{\lambda} \equiv \overline{\left( \sum_n \lambda_n^L \right)^{1/L}} < \frac{\langle k_{\text{in}} k_{\text{out}} \rangle}{\langle k_{\text{in}} \rangle} \quad (3)$$

for the network to be under-shortlooped and

$$\bar{\lambda} > \frac{\langle k_{\text{in}} k_{\text{out}} \rangle}{\langle k_{\text{in}} \rangle} \quad \text{if} \quad \kappa = \max_{L \in (1, L_c)} \kappa_L \ll 1 \quad (4)$$

for the network to be over-shortlooped *on average* over loop lengths  $L \in (1, L_c)$ . The overbar indicates average over  $L \in (1, L_c)$  for  $L_c$  satisfying the condition for Eq. (2) to be valid.

*Long loops.*—The above analysis applies to short loops. Counting long loops is a difficult problem for which approximate Monte Carlo [17] and statistical mechanics methods [18] have been proposed in the undirected case. To derive a necessary condition for long *directed* loops to be present, we use percolation predictions [6] for two-point correlated networks, where the out-degree of a node is correlated (beyond the random condition) with the in-degree of the nodes at the end points of its links. These networks are expected to account for the leading correlation term that distinguishes a real network from its uncorrelated random counterparts. In networks with two-point degree correlation, the percolation condition for the largest strongly connected component (LSCC) is  $\tilde{\Lambda} > 1$ , where  $\tilde{\Lambda}$  is the largest eigenvalue of the two-point correlation matrix  $C_{\mathbf{k}', \mathbf{k}} = [k_{\text{out}} P(\mathbf{k}' | \mathbf{k})]$  [6]. A strongly connected component of a network is a set of nodes where each node can reach and be reached by all the others through directed paths. For the uncorrelated random networks of the Molloy-Reed ensemble, the largest eigenvalue of matrix  $C$  reduces to the known result  $\tilde{\Lambda} = \frac{\langle k_{\text{in}} k_{\text{out}} \rangle}{\langle k_{\text{in}} \rangle}$ . For a specific network, which is not necessarily well approximated by an uncorrelated ensemble average, we can use the approximation  $\tilde{\Lambda} \simeq \Lambda$ , where  $\Lambda$  denotes the largest eigenvalue of the adjacency matrix [19]. Consequently the percolation conditions for the real and randomized networks are respectively  $\Lambda > 1$  and  $\frac{\langle k_{\text{in}} k_{\text{out}} \rangle}{\langle k_{\text{in}} \rangle} > 1$ . Since the existence of a giant LSCC is a necessary condition for the network to have long directed loops, long loops are strongly suppressed when  $\Lambda \leq 1$ . Because percolation only provides a necessary condition for the existence of long loops, we make quantitative predictions using a modified message-passing algorithm [20] based on the belief propagation (BP) algorithm proposed in [18]. The algorithm provides

an estimation for the entropy  $\sigma(L) = \log(\mathcal{N}_L)/N$  of the loops of length  $L$ , from which we calculate  $\mathcal{N}_L$ . Within the conditions discussed in [20], namely, that the network is large and has a large number of loops, this algorithm is able to predict the maximal loop length  $L_{\max}$  reliably. However, as shown below, the BP algorithm predicts correctly the under- or over-looped nature of all networks in our database, including those with a small number of nodes or loops [21], and the results are in very good agreement with the behavior suggested by the relative values between  $\Lambda$  and  $\frac{\langle k_{in}k_{out} \rangle}{\langle k_{in} \rangle}$ .

**Real networks.**—We consider several real directed networks [22]: (i) Texas power grid; (ii) food webs (Chesapeake, Mondego, Littlerock, and Seagrass regions); (iii) metabolic network of *E. coli*, where the nodes represent metabolites; (iv) Notre Dame University’s WWW; (v) *C. elegans*’ neural network; and the (vi) transcription network of *S. cerevisiae*, where the nodes correspond to regulating and regulated genes. Figure 1 shows the distributions of short loops (measured using exact enumeration [23]) for both the directed and undirected versions of four real networks along with the randomized counterparts of same number of in- and out-links in each node. The randomized networks are well approximated by the theoretical predictions in Eqs. (1) and (2), as indicated by the lines in the figure. Directed networks tend to have less loops than undirected networks, as expected. However, while real undirected networks tend to have more loops than random ones, the opposite occurs in the directed case.

Indeed, six out of the nine directed networks we analyzed are under-shortlooped, as shown in Fig. 2 and Table I. The only exceptions are the metabolic and transcription networks, which are marginally over-

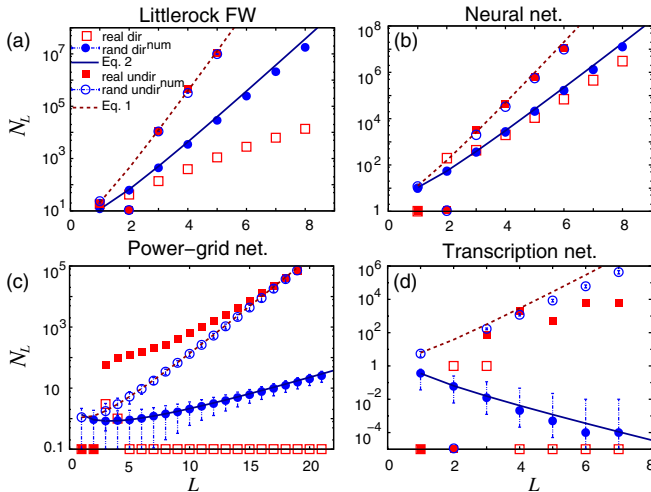


FIG. 1 (color online). Number of short directed and undirected loops in several networks, where different symbols correspond to the numerically determined values for the real and random counterparts of the networks. The lines indicate the theoretical predictions in Eqs. (1) and (2) for random networks. Points on the  $x$  axis indicate no loops.

shortlooped, and the WWW network, which is the only social network present in our database [24]. These findings are very different from what one would anticipate from previous studies on undirected networks, where highly clustered small-world networks prevail. Table I summarizes the network parameters and results for all directed networks analyzed, where  $\bar{\lambda}$  is calculated by summing over all loops up to a length cutoff  $L_c$  chosen to be 6 [24]. Our predictions compare well with direct data analysis.

**Conclusions.**—We have studied deviations in the loop statistics and provided criteria for determining if a network is underlooped or overlooped compared to its randomized counterparts. Empirical evidence coming from the study of different types of natural and technological networks shows that many of these different networks are under-shortlooped, a surprising result which is in sharp contrast with the tendency of undirected networks to be over-shortlooped. The only socio-technological network in our database, the ND WWW, contains instead very many short loops. We expect that our results will be important and further extended in the study of social, biological and technological systems. In social networks, the abundance of directed loops can be an important factor in the promotion of mutual reinforcement amongst agents [25], while in cellular and neural networks it can play a major role in information processing [26] and regulation [27]. In other systems, the reduced number of directed loops can lead to improved stability [8,9] and transportation properties [10], which we hope will stimulate other applications of our findings.

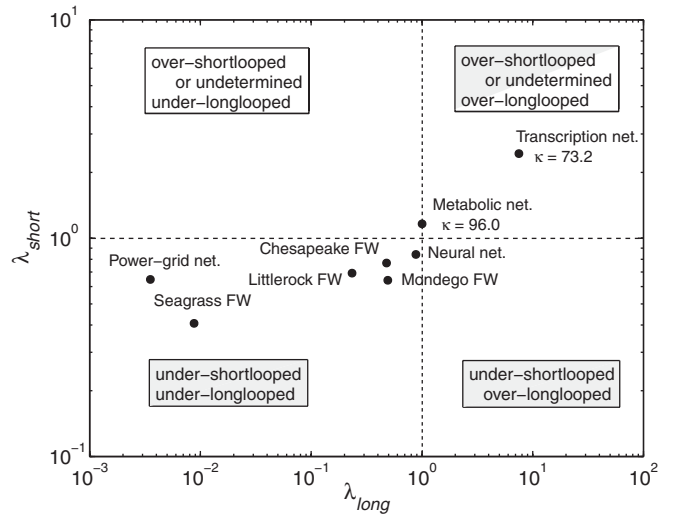


FIG. 2. Underlooped, overlooped, and undetermined regions in the  $\lambda_{\text{short}} \equiv \bar{\lambda}/(\langle k_{in}k_{out} \rangle/\langle k_{in} \rangle)$  vs  $\lambda_{\text{long}} \equiv L_{\max}^{\text{real}}/L_{\max}^{\text{rand}}$  diagram, where  $L_{\max}^{\text{real}}$  and  $L_{\max}^{\text{rand}}$  are predicted using the BP algorithm. The points correspond to the predictions for both short and long loops for the networks in Table I, except for the ND WWW, which is over-shortlooped and is not shown because it is difficult to calculate its  $\lambda_{\text{long}}$ . The actual counting of the loops confirms the predictions (Table I).

TABLE I. Properties of real directed networks: number of nodes  $N$  and links  $M$ , eigenvalue  $\Lambda$ ,  $\langle k_{in}k_{out} \rangle / \langle k_{in} \rangle$ , and spectral quantity  $\bar{\lambda}$  (left-hand side columns); loop structure and percolation properties (right-hand side columns). The values of  $\kappa$  for the undetermined and over-shortlooped cases are 96.0, 73.2, and 0.2 for the metabolic, transcription, and WWW network, respectively.

Network	N	Network parameters				$\bar{\lambda}$	Short loops <sup>b</sup>	Prediction/Actual <sup>a</sup>	
		M	$\Lambda$	$\frac{\langle k_{in}k_{out} \rangle}{\langle k_{in} \rangle}$	Percolation/LSCC <sup>c</sup>			Long loops <sup>d</sup>	
Littlerock FW	183	2494	7.00	11.47	7.93	und/und	p-p/(12 vs 92)	und/und <sup>e</sup>	
Chesapeake FW	39	177	2.85	3.12	2.40	und/und	p-p/(41 vs 76)	und/und	
Mondego FW	46	400	8.95	9.14	5.86	und/und	p-p/(76 vs 92)	und/und <sup>e</sup>	
Seagrass FW	48	226	1.00	4.05	1.65	und/und	np-p/(0 vs 75)	und/und	
Metabolic net.	532	596	2.85	2.58	3.00	undet/over	p-p/(82 vs 94)	undet/undet <sup>e</sup>	
Power-grid net.	4889	5855	1.00	1.36	0.88	und/und	np-p/(0.1 vs 33)	und/und	
ND WWW	325 729	1 497 135	152.00	43.14	153.32	over/over	p-p/(17 vs 41)	...	
Neural net.	306	2359	9.15	10.49	8.84	und/und	p-p/(78 vs 86)	und/und <sup>e</sup>	
Transcription net.	688	1079	1.32	0.36	0.88	undet/over	p-np/(0.4 vs 0.3)	over/over	

<sup>a</sup>Underlooped (und), overlooped (over), undetermined (undet), not determined numerically (...

<sup>b</sup>From left to right: predicted and actual values determined by averaging over the directed loops up to length  $L_c = 6$  ( $L_c = 3$  for the ND WWW).

<sup>c</sup>From left to right: predicted percolating ( $p$ ) or nonpercolating ( $np$ ) LSCC in real and random networks together with the actual percentage of nodes in the LSCC of the real vs random networks.

<sup>d</sup>From left to right: prediction for long loops obtained using the BP algorithm [20] to estimate  $L_{max}$  and the actual result obtained using exhaustive ...

<sup>e</sup>or partial enumeration of the loops.

The authors thank Dong-Hee Kim and Marian Boguñá for providing feedback on the manuscript. This work was supported by IST STREP GENNETEC Contract No. 034952 (G.B.), DOE LANL Contract No. DE-AC52-06NA25396 (N.G.), and NSF Grant No. DMS-0709212 (A. E. M.).

- [1] K. Klemm and S. Bornholdt, Proc. Natl. Acad. Sci. U.S.A. **102**, 18 414 (2005); T. Galla, J. Phys. A **39**, 3853 (2006).
- [2] R. Pastor-Satorras *et al.*, Phys. Rev. Lett. **87**, 258701 (2001).
- [3] S.N. Dorogovtsev *et al.*, arXiv:0705.0010v5 [Rev. Mod. Phys. (to be published)].
- [4] D. J. Watts and S. H. Strogatz, Nature (London) **393**, 440 (1998).
- [5] We consider degree correlations [2,6,7] of two types: between the in-out degrees of a given node and between the degrees of connected nodes.
- [6] M. Boguñá and M. A. Serrano, Phys. Rev. E **72**, 016106 (2005).
- [7] A. Capocci and F. Colaiori, Phys. Rev. E **74**, 026122 (2006).
- [8] A.-M. Neutel *et al.*, Nature (London) **449**, 599 (2007).
- [9] T. Nishikawa and A. E. Motter, Physica (Amsterdam) **224D**, 77 (2006); Phys. Rev. E **73**, 065106 (2006).
- [10] Z. Toroczkai *et al.*, Nature (London) **428**, 716 (2004); Phys. Rev. E **74**, 046114 (2006).
- [11] Z. Burda and A. Krzywicki, Phys. Rev. E **67**, 046118 (2003); M. Boguñá *et al.*, Eur. Phys. J. B **38**, 205 (2004).
- [12] G. Bianconi and M. Marsili, J. Stat. Mech. (2005) P06005.
- [13] Z. Burda *et al.*, Phys. Rev. E **70**, 026106 (2004).
- [14] A.-L. Barabási and R. Albert, Science **286**, 509 (1999).
- [15] This is not necessarily the case for other growing net-

works. See, for example, J. D. Noh, arXiv:0707.0560v2.

- [16] G. Bianconi and A. Capocci, Phys. Rev. Lett. **90**, 078701 (2003).
- [17] H. D. Rozenfeld *et al.*, J. Phys. A **38**, 4589 (2005).
- [18] E. Marinari *et al.*, Europhys. Lett. **73**, 8 (2006).
- [19] Because  $\text{Tr}C^L \approx \text{Tr}A^L$ , if the two matrices are expanders we have  $\bar{\Lambda} \approx \Lambda$ .
- [20] G. Bianconi and N. Gulbahce, arXiv:0709.1446v1.
- [21] For example,  $L_{max}^{real} / \langle L_{max}^{rand} \rangle = 0.48$  (BP) and 0.19 (exact) for the Chesapeake network while  $L_{max}^{real} / \langle L_{max}^{rand} \rangle < 0.004$  (BP) and  $< 0.08$  (exact) for the power grid, consistently indicating that long loops are underrepresented in these networks.
- [22] The power-grid data set was provided by Ken Werley and the metabolic network was generated using flux balance analysis [28] on the reconstructed model iJE660 available at <http://gcrq.ucsd.edu/organisms/>; We consider the metabolic network without the metabolites  $\text{CO}_2$ ,  $\text{NH}_3$ ,  $\text{PP}_i$ ,  $\text{P}_i$ , ATP, ADP, NAD, NADP, and NADH, as in D. A. Fell and A. Wagner, Nat. Biotechnol. **18**, 1121 (2000); The other network data are available at <http://vlado.fmf.unilj.si/pub/networks/data/>, [www.cosinproject.org/](http://www.cosinproject.org/), <http://cdg.columbia.edu/cdg/>, [www.nd.edu/~networks/](http://www.nd.edu/~networks/) and [www.weizmann.ac.il/mcb/UriAlon/](http://www.weizmann.ac.il/mcb/UriAlon/).
- [23] R. Tarjan, SIAM J. Comput. **2**, 211 (1973).
- [24] For the ND WWW, due to its size, we use  $L_c = 3$  and we estimate  $\bar{\lambda}$  and  $\kappa$  from the 100 largest eigenvalues and corresponding eigenvectors.
- [25] C. Troutman *et al.* (unpublished).
- [26] A. Roxin *et al.*, Phys. Rev. Lett. **92**, 198101 (2004); S. Song *et al.*, PLoS Biol. **3**, e68 (2005).
- [27] R. Milo *et al.*, Science **298**, 824 (2002); S. Krishna *et al.*, Nucleic Acids Res. **34**, 2455 (2006).
- [28] J. S. Edwards and B. O. Palsson, Proc. Natl. Acad. Sci. U.S.A. **97**, 5528 (2000).