

Geographic routing in social networks

David Liben-Nowell^{*†§}, Jasmine Novak[†], Ravi Kumar^{†¶}, Prabhakar Raghavan^{¶||}, and Andrew Tomkins^{†¶}

^{*}Department of Mathematics and Computer Science, Carleton College, 1 North College Street, Northfield, MN 55057; [†]IBM Almaden Research Center, 650 Harry Road, San Jose, CA 95120; [‡]Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA 02139; [¶]Yahoo! Research Labs, 701 First Avenue, Sunnyvale, CA 94089; and ^{||}Verity, Inc., Sunnyvale, CA 94089

Edited by Ronald L. Graham, University of California at San Diego, La Jolla, CA, and approved June 27, 2005 (received for review April 12, 2005)

We live in a “small world,” where two arbitrary people are likely connected by a short chain of intermediate friends. With scant information about a target individual, people can successively forward a message along such a chain. Experimental studies have verified this property in real social networks, and theoretical models have been advanced to explain it. However, existing theoretical models have not been shown to capture behavior in real-world social networks. Here, we introduce a richer model relating geography and social-network friendship, in which the probability of befriending a particular person is inversely proportional to the number of closer people. In a large social network, we show that one-third of the friendships are independent of geography and the remainder exhibit the proposed relationship. Further, we prove analytically that short chains can be discovered in every network exhibiting the relationship.

routing algorithms | small worlds | population networks | rank-based friendships | six degrees of separation

Ancedotal evidence that we live in a “small world,” where arbitrary pairs of people are connected through extremely short chains of intermediary friends, is ubiquitous. Sociological experiments, beginning with the seminal work of Milgram and his coworkers (1–3) and Killworth and Bernard (4), have shown that a source person can transmit a message to a target through only a small number of intermediate friends, using only scant information about the target’s geography and occupation; in other words, social networks are navigable small worlds. On average, the successful messages passed from source to target through six intermediaries; from this experiment came the popular notion of “six degrees of separation.”

As part of the recent surge of interest in networks, there has been active research exploring strategies for navigating synthetic and small-scale social networks (5–12), including routing through common membership in groups, popularity, and geographic proximity, the property on which we focus. In both the experiments by Milgram and coworkers (1–3) and a more recent e-mail-based replication (13), one sees the message geographically “zeroing in” on the target step by step as it is passed on. Furthermore, subjects report that geography and occupation are by far the two most important dimensions in choosing the next step in the chain (4), and geography tends to predominate in early steps (13). These reports lead to an intriguing question: what is the connection between friendship and geography, and to what extent can this connection explain the navigability of large-scale real-world social networks? Of course, adding nongeographic dimensions to routing strategies, especially once the chain has arrived at a point geographically close to the target, can make routing more efficient, sometimes considerably (2, 8, 9, 14). However, geography appears to be the single most valuable dimension for routing, and we are thus interested in understanding how powerful geography alone may be.

Here, we present a study that combines measurements of the role of geography in a large social network with theoretical modeling of path discovery, using the measurements to validate and inform the theoretical results. First, a simulation-based study on a 500,000-person online social network reveals that routing through geographic information alone allows people to discover short paths to a target city. Second, through empirical investigation of the rela-

tionship between geography and friendship in this network, we discover that $\approx 70\%$ of friendships are derived from geographical processes, but existing models that predict the probability of friendship solely on the basis of geographic distance are too weak to explain these friendships, rendering previous theoretical results inapplicable. [The proportion of links in a network that are between two entities separated by a particular geographic distance has been studied in a number of different contexts: the infrastructure of the Internet (15–17), small-scale e-mail networks within a company (9, 14), transportation networks (17), and wireless-radio networks (18).] Finally, we propose a density-aware model of friendship formation called rank-based friendship, relating the probability that a person befriends a particular candidate to the inverse of the number of closer candidates. We are able to rigorously prove that the presence of rank-based friendship for any population density implies that the resulting network will contain discoverable short paths to small destination regions. Rank-based friendship is then shown by measurement to be present in the large social network. Thus, we observe that a large online social network exhibits short paths under a simple geographical routing model, and we identify rank-based friendship as an important social-network property whose presence in the network implies the existence of short paths under geographic routing.

The social network that we consider comprises the 1,312,454 bloggers in the LiveJournal online community (www.livejournal.com), in February 2004. A blog, abbreviated from “web log,” is an online diary, often updated daily, typically containing reports on the user’s personal life, reactions to world events, and commentary on other blogs. In the LiveJournal system, each blogger also explicitly provides a profile, including his or her geographic location, topical interests, and a list of other bloggers whom he or she considers to be a friend. Of these 1.3 million bloggers, there are 495,836 in the continental United States who list a hometown and state that we find in the United States Geological Survey Geographic Names Information System (ref. 19; <http://geonames.usgs.gov>) and are thus able to map to a longitude and latitude; the resolution of our geographic data is limited to the level of towns and cities. Thus, our discussion of routing is from the perspective of reaching the home town or city of the destination individual. That is, we study the problem of “global” routing, in which the goal is to direct a message to the target’s city; once the proper locality has been reached, a “local” routing problem must then be solved to move the message from the correct city down to the correct person. There is evidence that geographic concerns predominate in the early stages of real-world message passing to solve the global-routing problem before the target individual is found by using a wide set of potential nongeographic factors, like interests or profession (1, 13).

The LiveJournal social network is defined as the $\approx 500,000$ LiveJournal users with locations in the United States in their profiles, with the “ u is a friend of v ” relationship defined by the explicit appearance of blogger u in the list of friends in the profile of blogger v . Let $d(u, v)$ denote the geographic distance between two people u and v . There are 3,959,440 friendship links in this

This paper was submitted directly (Track II) to the PNAS office.

[§]To whom correspondence should be addressed. E-mail: dlibenno@carleton.edu.

© 2005 by The National Academy of Sciences of the USA

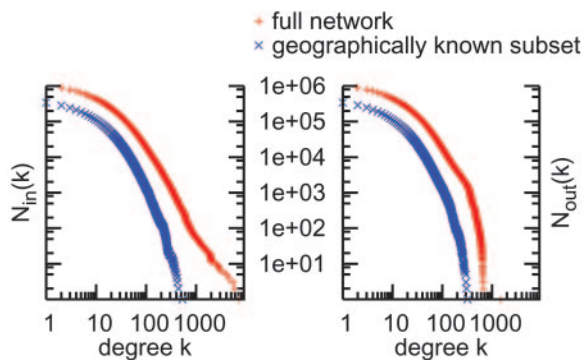


Fig. 1. In-degree (Left) and out-degree (Right) distributions in LiveJournal. For each k , the number $N_{in}(k)$ of LiveJournal users who are listed as a friend of at least k users and the number $N_{out}(k)$ of people who list at least k friends are shown, both for all 1,300,000 users and the 500,000 users who list locatable hometowns in the United States.

directed network, an average of about eight friends per user. (Although u can be listed as a friend of v without v being listed as a friend of u , we find that 80% of friendships are reciprocal.) This network exhibits many of the same important structural features observed in other social networks (20, 21). For example, 384,507 people (77.6%) form a giant component in which any two people u and v are connected by chains of friends leading from the coefficient of the network; that is, the proportion of the time that u and v are themselves friends if they have a common friend w is 0.2; this high clustering coefficient is characteristic of social networks (22). Fig. 1 shows the degree distribution, the fraction of the population with at least k friends for every k , both for the entire 1,312,454-person network and the 495,836 locatable people in the United States. The in-degree log/log plot is more linear than the out-degree plot, but both appear far more parabolic than linear; these curves provide some evidence supporting a log-normal degree distribution in social networks, instead of a power-law distribution (23–25).

Geographic Routing

We perform a simulated version of the message-forwarding experiment in the LiveJournal social network, using only geographic information to choose the next message holder in a chain. This simulation may be viewed as a thought experiment, with two goals. First, we seek to determine whether individuals using purely geographic information in a simple way can succeed in discovering short paths to a destination city. Second, we seek to analyze the applicability of existing theoretical models that explain the presence or absence of short discoverable paths in networks (8, 10–12). Our simulation should not be viewed as a replication of real-world experiments studying human behavior (1, 13), but rather as an investigation into what would be possible for people participating in a message-passing experiment in such a network. Our approach, using a large-scale network of real-world friendships but simulating the forwarding of messages, allows us to investigate the performance of simple routing schemes without suffering from a reliance on the voluntary participation of the people in the network. (Further, in real-world experiments, dropout rates may depend on chain length, possibly biasing results toward shorter estimates of chain length.) Our detailed information on the location of every friend of every participant then allows us to analyze in detail the underlying geographic basis of friendship in explaining these results.

In our simulation, messages are forwarded by using the geographically greedy routing algorithm GEOGREEDY (10): if a person u currently holds the message and wants to eventually reach a target t , then she considers her set of friends and chooses as the next step

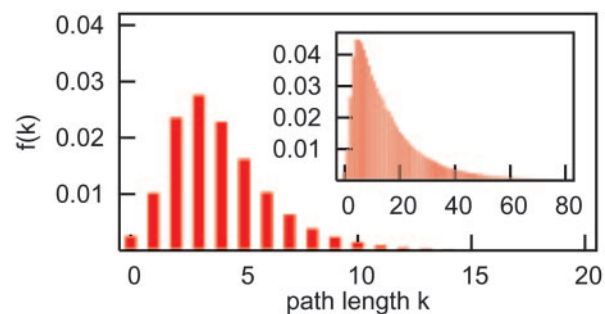


Fig. 2. Results of GEOGREEDY on LiveJournal. In each of 500,000 trials, a source s and target t are chosen randomly; at each step, the message is forwarded from the current message-holder u to the friend v of u geographically closest to t . If $d(v, t) > d(u, t)$, then the chain is considered to have failed. The fraction $f(k)$ of pairs in which the chain reaches t 's city in exactly k steps is shown (12.78% chains completed; median 4, $\mu = 4.12$, $\sigma = 2.54$ for completed chains). (Inset) For 80.16% completed, median 12, $\mu = 16.74$, $\sigma = 17.84$; if $d(v, t) > d(u, t)$ then u picks a random person in the same city as u to pass the message to, and the chain fails only if there is no such person available.

in the chain the friend in this set who is geographically closest to t . If u is closer to the target than all of her friends, then she gives up, and the chain terminates. When sources and targets are chosen randomly, we find that the chain successfully reaches the city of the target in $\approx 13\%$ of the trials, with a mean completed-chain length of slightly more than four (Fig. 2). Recall that our data set does not contain intracity geographic information, so we do not attempt to reach the target itself; we instead focus on global routing in these simulated experiments. For a target t chosen uniformly at random from the network, the average population of t 's city is 1,306 and always under 8,000; we therefore study the success of geography in narrowing the search from 500,000 users across the United States to the on-average 1,300 residents in a particular city.

A success rate of 13% with an average length of just over four in this simulated experiment shows a surface similarity to Milgram's original experiment (1), where 18% of chains were completed to the destination individual, with an average length of just under six. Our experiment, however, routes messages only to the destination city and does not suffer from problems of voluntary participation, which may explain why our completion rate is significantly higher than that of Dodds *et al.* (13). On the other hand, our simulated participants have a much narrower choice of actions, as they are restricted to friends whom they have explicitly listed in their profile and can forward only to the friend geographically closest to the target. Overall, we conclude that, even under restrictive forwarding conditions, geographic information is sufficient to perform global routing in a significant fraction of cases. This simulated experiment may be taken as a lower bound on the presence of short discoverable paths, because only the on-average eight friends explicitly listed in each LiveJournal profile are candidates for forwarding. By way of comparison, we modify the routing algorithm: an individual u who has no friend geographically closer to the target instead forwards the message to a person selected at random from u 's city. Under this modification, chains complete 80% of the time, with median length 12 and mean length 16.74 (see Fig. 2). The completion rate is not 100% because a chain may still fail by landing at a location in which no inhabitant has a friend closer to the target. This modified experiment may be taken as an upper bound on completion rate, when the simulated individuals doggedly continue forwarding the message as long as any closer friend or collocated individual exists.

The Geographic Basis of Friendship

Thus, geographically greedy routing in the LiveJournal social network under a restrictive model allowed 13% of paths to reach

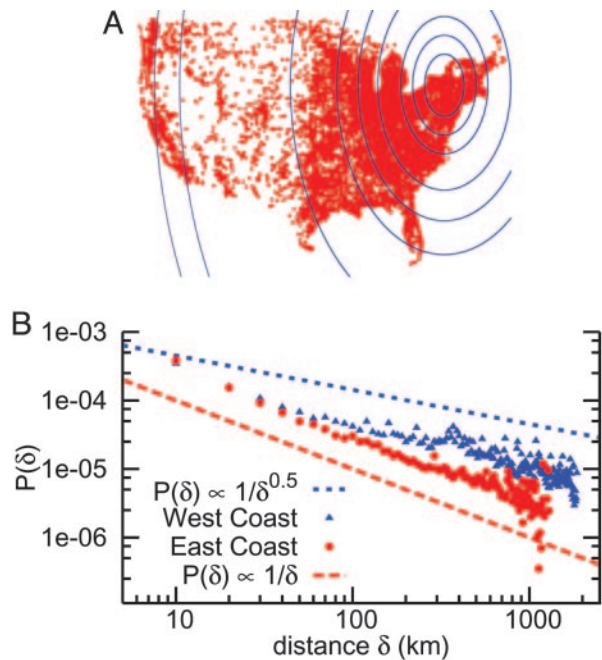


Fig. 4. Evidence of the nonuniformity of the LiveJournal population. (A) A dot is shown for every distinct United States location home to at least one LiveJournal user. The population of each successive displayed circle (all centered on Ithaca, NY) increases by 50,000 people. Note that the gap between the 350,000- and 400,000-person circles encompasses almost the entire Western United States. (B) We show the relationship between friendship probability and geographic distance, as in Fig. 3, restricted to people living on the West Coast (California, Oregon, and Washington) and the East Coast (from Virginia to Maine), respectively.

characterization of group structures in social networks, and a number of researchers have presented extensions and improved analyses of these models (29–34).

If the LiveJournal data confirm this relationship between friendship probability and geographic distance, i.e., if the probability $f[d(u, v)]$ of geographic friendship between u and v is roughly proportional to $1/d(u, v)^2$, as the Earth's surface is 2D, then the finding of short paths by GEOGREEDY will be explained. Fig. 3B explores this conjecture, showing the best fit for geographic-friendship probability as a function of geographic distance. However, the geographic-friendship probability between two people separated by distance δ is best modeled as $f(\delta) = 1/\delta^\alpha$ for $\alpha \approx 1$; Kleinberg's results (12), and those of all extensions to his model, in fact show that this exponent cannot result in a navigable social network based on a 2D mesh. Yet the LiveJournal network is clearly navigable, as shown in our simulation of GEOGREEDY.

This seeming contradiction is explained by a large variance in population density across the LiveJournal network, which is thus ill-approximated by the uniform 2D mesh of Kleinberg's model. Fig. 4 explores population patterns in more detail. Fig. 4A shows concentric circles representing bands of equal population around Ithaca, NY. Under uniform population density, the width of each band should shrink as the distance from Ithaca increases. In the LiveJournal data set, however, the distance between annuli actually gets larger instead of smaller. Furthermore, purely distance-based predictions imply that the probability of a friendship at a given distance should be constant for different people in the network. Fig. 4B explores this concern, showing a distinction in friendship probability as a function of distance for residents of the East and West coasts. Thus a geographic model of friendship must be based on more than distance alone, as no accurate uniform description of friendship as a function of distance applies throughout the network.

To summarize, we have shown that any model of friendship that is based solely on the distance between people is insufficient to explain the geographic nature of friendships in the LiveJournal network. The model of Watts *et al.* (8) naturally captures individuals' geographic similarity by approximating their Euclidean distance by their distance in some hierarchy, assigning friendship probability based on a function of this distance as long as geography remains the most proximate coordinate. Thus current models do not take into account the organization of people into cities of arbitrary location and population, and they cannot explain the success of the simulated message-passing experiment. We therefore seek a network model that reconciles the linkage patterns in real-world networks with the success of GEOGREEDY on these networks. Such a model must be based on something beyond distance alone.

Population Networks and Rank-Based Friendship

We explore the idea that a simple model of the probability of friendship that combines distance and density may apply uniformly over the network. Consider a person u and a person v who lives 500 m away from u . In rural Iowa, say, u and v are probably next-door neighbors and very likely know each other; in Manhattan, there may be 10,000 people who live closer to u than v does, and the two are unlikely to have ever even met. This discrepancy suggests why geographic distance alone is insufficient as the basis for a geographical model. Instead, our model uses rank as the key geographic notion: when examining a friend v of u , the relevant quantity is the number of people who live closer to u than v does. Formally, we define the rank of v with respect to u as

$$\text{rank}_u(v) = |\{w : d(u, w) < d(u, v)\}|.$$

Under the rank-based friendship model, we model the probability that u and v are geographic friends by

$$\Pr[u \rightarrow v] \propto \frac{1}{\text{rank}_u(v)}.$$

Under this model, the probability of a link from u to v depends only on the number of people within distance $d(u, v)$ of u and not on the geographic distance itself; thus the nonuniformity of LiveJournal population density fits naturally into this framework. Although either distance- or rank-based models may be appropriate in some contexts, we will show that (i) analytically, rank-based friendship implies that GEOGREEDY will find short paths in any social network; and (ii) empirically, the LiveJournal network exhibits rank-based friendship.

We model a geographic n -person social network as follows. Consider a 2D N as a model of the 2D surface of the Earth. The grid divides the Earth's surface into small squares; we may take N to represent 1° -by- 1° squares centered at the intersection of integral lines of longitude and latitude, for example. At every point $(x, y) \in N$, we have a population $p(x, y)$ denoting the number of people who live in the square centered at (x, y) , with $\sum_{x,y} p(x, y) = n$ and $p(x, y) > 0$. The condition that $p(x, y) > 0$ is imposed to guarantee that a routing algorithm can always make some progress toward any target at every step of the chain. We refer to the combination of the grid N and the population p as a population network. Building on the navigable small-world model of Kleinberg (11, 10), we model linkage in population networks as follows. Each person u in the network has an arbitrarily chosen neighbor in each of the four adjacent grid points: north, east, south, and west. In addition to these four neighbors, person u has a long-range link to a fifth person chosen according to rank-based friendship, that is, the probability that u chooses v as her long-range link is inversely proportional to $\text{rank}_u(v)$.

The notion of adding links with probability inversely proportional to the number of closer candidates is implicit in Kleinberg's work

P , then the expected number of steps before the message reaches a person within distance $d(s, t)/2$ of t is at most $c \cdot \log^2 n$. After the distance to t is halved $\log n$ times, the message will have arrived at its destination. Thus the expected number of steps to reach t is at most $c \cdot \log^3 n$. To prove the claim, we show that the probability that a person forwards the message to someone within the small $d(s, t)/2$ -radius neighborhood around t is at least $1/\log n$ times the relative densities of the small neighborhood around t compared with a larger neighborhood containing both s and t . By taking expectations over t and appropriately approximating the densities of these neighborhoods, we show that the expected number of steps before the message reaches the small neighborhood around t is at most $c \cdot \log^2 n$.

There is significant evidence from real-world message-passing experiments that an effective routing strategy typically begins by making long geography-based hops as the message leaves the source and ends by making hops based on attributes other than geography (1, 13). Thus there is a transition from geography-based to nongeography-based routing at some point in the process. We can extend our theorem to show that short paths are constructible by using GEOGREEDY through any level of resolution in which rank-based friendship holds.

Geographic Linking in the LiveJournal Social Network

We return to the LiveJournal social network to show that rank-based friendship holds in a real network. The relationship between $\text{rank}_v(u)$ and the probability that u is a friend of v shows an approximately inverse linear fit for ranks up to $\approx 100,000$ (Fig. 5A). Because the LiveJournal data contain geographic information limited to the level of towns and cities, our data do not have sufficient resolution to distinguish between all pairs of ranks. (Specifically, an average person in the network lives in a city of population 1,306.) Thus in Fig. 5B we show the same data, where the probabilities are averaged over a range of 1,306 ranks. (Because of the logarithmic scale of the rank axis, the sliding window may appear to apply broad smoothing; however, smoothing by a window of size 1,300 causes a point to be influenced by $<0.3\%$ of the closest points on the curve.) This experiment validates that the LiveJournal social network does exhibit rank-based friendship, which thus yields a sufficient explanation for the experimentally observed navigability properties.

Fig. 6 displays the same data as in Fig. 5B, restricted to the East and West coasts. The slopes of the lines for the two coasts are nearly the same, and they are much closer together than the distance/friendship-probability slopes shown in Fig. 4B, confirming that probabilities based on ranks are a more accurate representation than distance-based probabilities.

In summary, the LiveJournal social network displays a surprising and variable relationship between geographic distance and probability of friendship, which is inconsistent with earlier theoretical models. Further, the network evinces short paths discoverable by using geography alone, even though existing models predict the opposite. We present rank-based friendship, a core mechanism for geographically biased friendship formation that may be embedded into a wide variety of broader models. Rank-based friendship is unique in simultaneously providing two desirable properties: (i) it matches our experimental observations regarding the relationship between geography and friendship; and (ii) it admits a mathematical proof that networks exhibiting rank-based friendship will contain discoverable short paths. As a validation of this theorem, the LiveJournal network exhibits rank-based friendship and does indeed contain discoverable short paths. Thus, we nominate rank-based friendship as a mechanism that has been empirically observed in real networks and theoretically guarantees small-world properties.

In fact, rank-based friendship explains geographic routing to a destination city; our data do not allow conclusions about routing within a city. Watts *et al.* (8) suggest that multiple independent dimensions play a role in message routing, and our results confirm this viewpoint: on average about one-third of LiveJournal friendships are independent of geography and may derive from other dimensions, like occupation and interests. These edges may play a role in local routing within the destination city and may also supplement the geographic links in global routing to the city; characterizing these geography-independent friendships is an interesting area for future work.

We have shown that the natural mechanisms of friendship formation result in rank-based friendship: people in aggregate have formed relationships with almost exactly the connection between friendship and rank that is required to produce a navigable small world. In a lamentably imperfect world, it is remarkable that people form friendships so close to the perfect distribution for navigating their social structures.

- Milgram, S. (1967) *Psychol. Today* **1**, 61–67.
- Travers, J. & Milgram, S. (1969) *Sociometry* **32**, 425–443.
- Korte, C. & Milgram, S. (1970) *J. Pers. Soc. Psychol.* **15**, 101–118.
- Killworth, P. & Bernard, H. (1978) *Soc. Networks* **1**, 159–192.
- Adamic, L. A., Lukose, R. M., Puniyani, A. R. & Huberman, B. A. (2001) *Phys. Rev. E* **64**, 046135.
- Kim, B. J., Yoon, C. N., Han, S. K. & Jeong, H. (2002) *Phys. Rev. E* **65**, 027103.
- Adamic, L. A., Lukose, R. M. & Huberman, B. A. (2002) *Local Search in Unstructured Networks* (Wiley, New York).
- Watts, D. J., Dodds, P. S. & Newman, M. E. J. (2002) *Science* **296**, 1302–1305.
- Adamic, L. A. & Adar, E. (2003) e-Print Archive, <http://www.arxiv.org/abs/cond-mat/0310120>.
- Kleinberg, J. M. (2000) *Nature* **406**, 845.
- Kleinberg, J. (2000) in *Proceedings of the ACM Symposium on Theory of Computing*, ed. Yao, F. (ACM Press, New York), pp. 163–170.
- Kleinberg, J. (2001) *Adv. Neural Inform. Process.* **14**, 431–438.
- Dodds, P. S., Muhamad, R. & Watts, D. J. (2003) *Science* **301**, 827–829.
- Adamic, L. A. & Huberman, B. A. (2004) in *Complex Networks*, eds. Ben-Naim, E., Frauenfelder, H. & Toroczkai, Z. (Springer, New York), No. 650, pp. 371–398.
- Yook, S.-H., Jeong, H. & Barabási, A.-L. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 13382–13386.
- Gorman, S. P. & Kulkarni, R. (2004) *Environ. Plan. B* **31**, 273–296.
- Gastner, M. T. & Newman, M. E. J. (2004) e-Print Archive, <http://www.arxiv.org/abs/cond-mat/0407680>.
- Miller, L. E. (2001) *J. Res. Natl. Inst. Stan.* **106**, 401–412.
- United States Geological Survey (2000) *2000 Census* (United States Geological Survey, Reston, VA).
- Newman, M. E. J. (2003) *SIAM Rev.* **45**, 167–256.
- Wasserman, S. & Faust, K. (1994) *Social Network Analysis* (Cambridge Univ. Press, Cambridge, U.K.).
- Watts, D. J. & Strogatz, S. H. (1998) *Nature* **393**, 440–442.
- Barabási, A.-L. & Albert, R. (1999) *Science* **286**, 509–512.
- Barabási, A.-L., Albert, R. & Jeong, H. (1999) *Physica A* **272**, 173–187.
- Mitzenmacher, M. (2004) *Internet Math.* **1**, 226–251.
- Caimcross, F. (1997) *The Death of Distance* (Harvard Business School Press, Boston).
- Kiesler, S. & Cummings, J. N. (2002) in *Distributed Work*, eds. Hinds, P. & Kiesler, S. (MIT Press, Cambridge, MA), pp. 57–82.
- Alon, N., Karp, R. M., Peleg, D. & West, D. (1995) *SIAM J. Comp.* **24**, 78–100.
- Barrière, L., Fraigniaud, P., Kranakis, E. & Krizanc, D. (2001) in *Proceedings of the International Conference on Distributed Computing*, ed. Welch, J. (Springer, London), pp. 270–284.
- Slivkins, A. (2005) in *Proceedings of the Symposium on Principles of Distributed Computing*, ed. Aspnes, J. (ACM Press, New York), pp. 41–50.
- Fraigniaud, P., Gavoille, C. & Paul, C. (2004) in *Proceedings of the Symposium on Principles of Distributed Computing*, ed. Kutten, S. (ACM Press, New York), pp. 169–178.
- Fraigniaud, P. (2005) *A New Perspective on the Small-World Phenomenon: Greedy Routing in Tree-Decomposed Graphs*, Technical Report LRI-1397 (University Paris-Sud, Paris).
- Martel, C. & Nguyen, V. (2004) in *Proceedings of the Symposium on Principles of Distributed Computing*, ed. Kutten, S. (ACM Press, New York), pp. 179–188.
- Nguyen, V. & Martel, C. (2005) in *Proceedings of the Symposium on Discrete Algorithms*, ed. Buchsbaum, A. (Society for Industrial and Applied Mathematics, Philadelphia), pp. 311–320.
- Newman, M. E. J. & Watts, D. J. (1999) *Phys. Rev. E* **60**, 7332–7342.
- Csányi, G. & Szendrői, B. (2004) *Phys. Rev. E* **70**, 016122.
- Kumar, R., Liben-Nowell, D., Novak, J., Raghavan, P. & Tomkins, A. (2005) *Theoretical Analysis of Geographic Routing in Social Networks*, Technical Report MIT-LCS-TR-990 (MIT Press, Cambridge, MA).