

Structural diversity in social contagion

Johan Ugander^a, Lars Backstrom^b, Cameron Marlow^b, and Jon Kleinberg^{c,1}

^aCenter for Applied Mathematics and ^cDepartment of Computer Science, Cornell University, Ithaca, NY 14853; and ^bFacebook, Menlo Park, CA 94025

Edited by Ronald L. Graham, University of California at San Diego, La Jolla, CA, and approved February 21, 2012 (received for review October 6, 2011)

The concept of contagion has steadily expanded from its original grounding in epidemic disease to describe a vast array of processes that spread across networks, notably social phenomena such as fads, political opinions, the adoption of new technologies, and financial decisions. Traditional models of social contagion have been based on physical analogies with biological contagion, in which the probability that an individual is affected by the contagion grows monotonically with the size of his or her “contact neighborhood”—the number of affected individuals with whom he or she is in contact. Whereas this contact neighborhood hypothesis has formed the underpinning of essentially all current models, it has been challenging to evaluate it due to the difficulty in obtaining detailed data on individual network neighborhoods during the course of a large-scale contagion process. Here we study this question by analyzing the growth of Facebook, a rare example of a social process with genuinely global adoption. We find that the probability of contagion is tightly controlled by the number of connected components in an individual’s contact neighborhood, rather than by the actual size of the neighborhood. Surprisingly, once this “structural diversity” is controlled for, the size of the contact neighborhood is in fact generally a negative predictor of contagion. More broadly, our analysis shows how data at the size and resolution of the Facebook network make possible the identification of subtle structural signals that go undetected at smaller scales yet hold pivotal predictive roles for the outcomes of social processes.

social networks | systems

Social networks play host to a wide range of important social and nonsocial contagion processes (1–8). The microfoundations of social contagion can, however, be significantly more complex, as social decisions can depend much more subtly on social network structure (9–17). In this study we show how the details of the network neighborhood structure can play a significant role in empirically predicting the decisions of individuals.

We perform our analysis on two social contagion processes that take place on the social networking site Facebook: the process whereby users join the site in response to an invitation e-mail from an existing Facebook user (henceforth termed “recruitment”) and the process whereby users eventually become engaged users after joining (henceforth termed “engagement”). Although the two processes we study formally pertain to Facebook, their details differ considerably; the consistency of our results across these differing processes, as well as across different national populations (*Materials and Methods*), suggests that the phenomena we observe are not specific to any one modality or locale.

The social network neighborhoods of individuals commonly consist of several significant and well-separated clusters, reflecting distinct social contexts within an individual’s life or life history (18–20). We find that this multiplicity of social contexts, which we term structural diversity, plays a key role in predicting the decisions of individuals that underlie the social contagion processes we study.

We develop means of quantifying such structural diversity for network neighborhoods, broadly applicable at many different scales. The recruitment process we study primarily features small neighborhoods, but the on-site neighborhoods that we study in the context of engagement can be considerably larger. For small neighborhoods, structural diversity is succinctly measured by the number of connected components of the neighborhood. For larger neighborhoods, however, merely counting connected components

fails to distinguish how substantial the components are in their size and connectivity. To determine whether the structural diversity of on-site neighborhoods is a strong predictor of on-site engagement, we evaluate several variations of the connected component concept that identify and enumerate substantial structural contexts within large neighborhood graphs. We find that all of the different structural diversity measures we consider robustly predict engagement. For both recruitment and engagement, structural diversity emerges as an important predictor for the study of social contagion processes.

Results

User Recruitment. To study the spread of Facebook as it recruits new members, we require information not just about Facebook’s users but also about individuals who are not yet users. Thus, suppose that an individual A is not a user of Facebook; it is still possible to identify a set of Facebook users that A may know because these users have all imported A ’s e-mail address into Facebook. We define this set of Facebook users possessing A ’s e-mail address to be A ’s contact neighborhood in Facebook. This contact neighborhood is the subset of potential future friendship ties that can be determined from the presence of A ’s e-mail address (Fig. 1A). Whereas A may in fact know many other people on Facebook as well, such additional friendship ties remain unknown for individuals who do not choose to register and so cannot be studied as a predictor of recruitment. The e-mail contact neighborhoods we study are generally quite small, typically on the order of five or fewer nodes.

We can now study an individual’s decision to join Facebook as follows. Facebook provides a tool through which its users can e-mail friends not on Facebook to invite them to join; such an e-mail invitation contains not only a presentation of Facebook and a profile of the inviter, but also a list of the other members of the individual’s contact neighborhood. We analyze a corpus of 54 million such invitation e-mails, and the fundamental question we consider is the following: How does an individual’s probability of accepting an invitation depend on the structure of his or her contact neighborhood?

Traditional hypotheses suggest that this probability should grow monotonically in the size of the contact neighborhood (3, 9, 10). What we find instead, however, is a striking stratification of acceptance probabilities by the number of connected components in the contact neighborhood (Fig. 1B–D and Fig. S1). When going beyond component count, one may suspect that edge density has a significant impact on the recruitment conversion rate: Among the single-component neighborhoods of a given size, there is a considerable structural difference between neighborhoods connected as a tree and those connected as a clique. However, within the controlled conditional datasets of

Author contributions: J.U., L.B., C.M., and J.K. designed research; J.U., L.B., C.M., and J.K. performed research; J.U., L.B., C.M., and J.K. contributed new reagents/analytic tools; J.U., L.B., C.M., and J.K. analyzed data; and J.U. and J.K. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

¹To whom correspondence should be addressed. E-mail: kleinber@cs.cornell.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1116502109/-DCSupplemental.

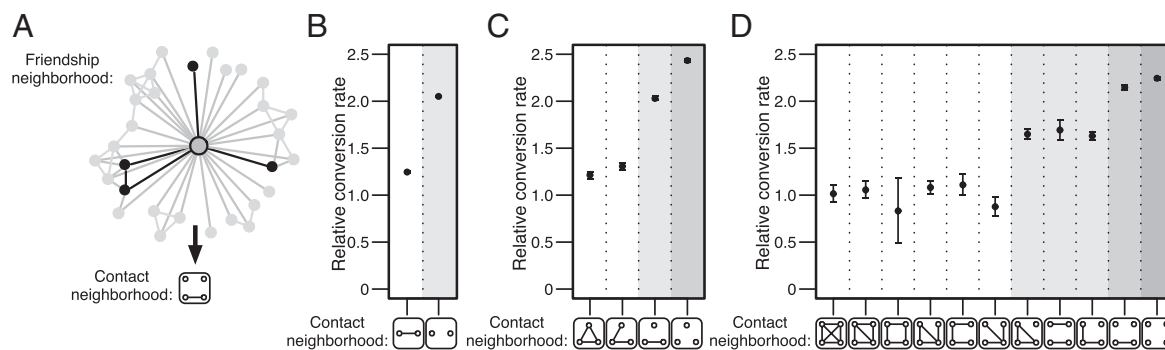


Fig. 1. Contact neighborhoods during recruitment. (A) An illustration of a small friendship neighborhood and a highlighted contact neighborhood consisting of four nodes and three components. (B–D) The relative conversion rates for two-node, three-node, and four-node contact neighborhood graphs. Shading indicates differences in component count. For five-node neighborhoods, see Fig. S1. Invitation conversion rates are reported on a relative scale, where 1.0 signifies the conversion rate of one-node neighborhoods. Error bars represent 95% confidence intervals and implicitly reveal the relative frequency of the different topologies.

one-component neighborhoods of sizes 4–6, we see that edge density has no discernible effect (Fig. 2A).

Moreover, we see that once component count is controlled for (Fig. 2B), neighborhood size is largely a negative indicator of conversion. In effect, it is not the number of people who have invited you, nor the number of links among them, but instead the number of connected components they form that captures your probability of accepting the invitation. Note that this analysis has been performed in aggregate and thus unavoidably reflects the decisions of different individuals. The ability to reliably estimate acceptance probabilities as a function of something as specific as the precise topology of the contact neighborhood is possible only because the scale of the dataset provides us with sufficiently many instances of each possible contact neighborhood topology (up through size 5).

We view the component count as a measure of “structural diversity,” because each connected component of an individual’s contact neighborhood hints at a potentially distinct social context in that individual’s life. Under this view, it is the number of distinct social contexts represented on Facebook that predicts the probability of joining. We show that the effect of this structural diversity persists even when other factors are controlled for. In particular, the number of connected components in the contact neighborhood remains a predictor of invitation acceptance even when restricted to individuals whose neighborhoods are demographically homogeneous (in terms of sex, age, and nationality; Fig. S2), thus controlling for a type of demographic diversity that is potentially distinct from structural diversity. The component count also remains a predictor of acceptance even when we compare neighborhoods that exhibit precisely the same mixture of “bridging” and “embedded” links (Fig. S3), the key distinction in sociological arguments based on information novelty (19, 20).

For contact neighborhoods consisting of two nodes, we observe that the probability an invitation is accepted is much higher when the two nodes in the neighborhood are not connected by a link (hence forming two connected components, Fig. 1B) compared with when they are connected (forming one component). Is there a way to identify cases where people are likely to know each other, even if they are not linked on Facebook? The photo tagging feature on Facebook suggests such a mechanism. Photographs uploaded to Facebook are commonly annotated by users with “tags” denoting the people present in the photographs. We can use these tags to deduce whether two unlinked nodes in a contact neighborhood have been jointly tagged in any photos, a property we refer to as “co-tagging,” which serves as an indication of a social tie through copresence at an event (21).

Using photo co-tagging, we find strong effects even in cases where the presence of a friendship tie is only implicit. If a contact neighborhood consists of two unlinked nodes that have

nevertheless been co-tagged in a photo, then the invitation acceptance probability drops to approximately what it is for a neighborhood of two linked nodes (Fig. 2C). In other words, being co-tagged in a photo indicates roughly the same lack of diversity as being connected by a friendship link. We interpret this result as further evidence that diverse endorsement is key to predicting recruitment. Meanwhile, when the two nodes are friends, co-tags offer a proxy for tie strength, and we see that if the two nodes have also been co-tagged, then the probability of an accepted invitation decreases further. From this we can interpret tie strength as an

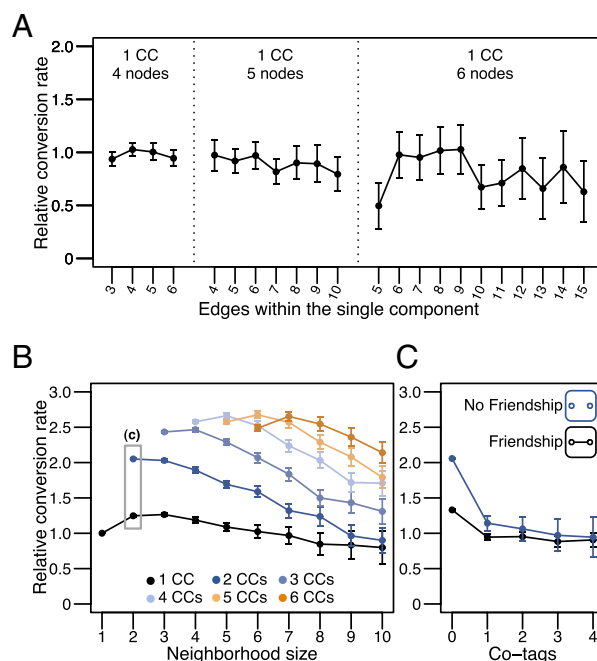


Fig. 2. Recruitment contact neighborhoods and component structure. (A) Conversion as a function of edge count neighborhoods with one connected component (1 CC) with four to six nodes, where variations in edge count predict no meaningful difference in conversion. (B) Conversion as a function of neighborhood size, separated by CC count. When component count is controlled for, size is a negative indicator of conversion. (C) Conversion as a function of tie strength in two-node neighborhoods, measured by photo co-tags, a negative indicator of predicted conversion. Recruitment conversion rates are reported on a relative scale, where 1.0 signifies the conversion rate of one-node neighborhoods. Error bars represent 95% confidence intervals.

extension of context, because two strongly tied nodes plausibly constitute an even less diverse endorsement neighborhood.

Finally, we study the position of the inviter within the neighborhood topologies. When studying recruitment, one might suspect that the structural position of the inviter—the person who extended the invitation—might signify differences in tie strength with the invitee and therefore might significantly affect the predicted conversion rate. We find that inviter position figures only slightly in the conversion rate (Fig. 3), with invitations stemming from a high-degree position in the contact neighborhood predicting only a slightly higher conversion rate than if the inviter is a peripheral node.

User Engagement. Participation in a social system such as Facebook is built upon a spectrum of social decisions, beginning with the decision to join (recruitment) and continuing on to decisions about how to choose a level of engagement. We now show how structural diversity also plays an analogous role in this latter type of decision process, studying long-term user engagement in the Facebook service. Whereas recruitment is a function of the complex interplay between multiple acts of endorsement, engagement is a function of the social utility a user derives from the service. Our study of engagement focuses on users who registered for Facebook during 2010, analyzing the diversity of their social neighborhoods 1 week after registration as a basis for predicting whether they will become highly engaged users 3 months later.

Users are considered engaged at a given time point if they have interacted with the service during at least 6 of the last 7 days. Facebook had 845 million monthly active users on December 31, 2011, and during the month of December 2011, an average of 360 million users were active on at least 6 out of the last 7 days. We define engagement on a weekly timescale to stabilize the considerable weekly variability of user visits. Our goal is therefore to predict whether a newly registered user will visit Facebook at least 6 of 7 days per week 3 months after registration.

Friendship neighborhoods on Facebook are significantly larger than the e-mail contact neighborhoods from our recruitment study. We focus our engagement study on a population of ~10 million users who registered during 2010 and had assembled neighborhoods consisting of exactly 10, 20, 30, 40, or 50 friends 1 week after registration. For social network neighborhoods of this size, we find that a neighborhood containing a large number of connected components primarily indicates a large number of one-node components, or “singletons”, and as such, it is not an accurate reflection of social context diversity.

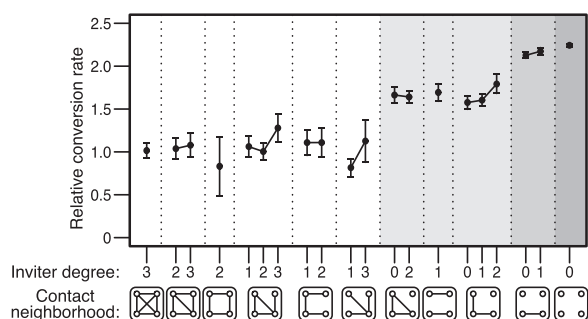


Fig. 3. Inviter position during recruitment. Shown is recruitment conversion as a function of neighborhood graph topology and inviter position in neighborhoods of size 4. The position of the inviter within the neighborhood graph is described exactly (up to symmetries) by node degree. Shading indicates differences in component count. Recruitment conversion rates are reported on a relative scale, where 1.0 signifies the conversion rate of one-node neighborhoods. Error bars represent 95% confidence intervals.

To address this, we evaluate three distinct parametric generalizations of component count. First, we measure diversity simply by considering only components over a certain size k . Second, we measure diversity by the component count of the k -core of the neighborhood graph (22), the subgraph formed by repeatedly deleting all vertices of degree less than k . Third, we define a measure that isolates dense social contexts by removing edges according to their *embeddedness*, the number of common neighbors shared by their two endpoints; intuitively this is an analog, for edges, of the type of node removal that defines the k -core. Adapting earlier work on embeddedness by Cohen (23), we define the k -brace of a graph to be the subgraph formed by repeatedly deleting all edges of embeddedness less than k and then deleting all single-node connected components. (Cohen's work was concerned with a definition equivalent to the largest connected component of the k -brace; because we deal with the full subgraph of all nontrivial components, it is useful to adapt the definitions as needed.) Examples of these three measures applied to a neighborhood graph are shown in Fig. 4 *A* and *B*, illustrating the

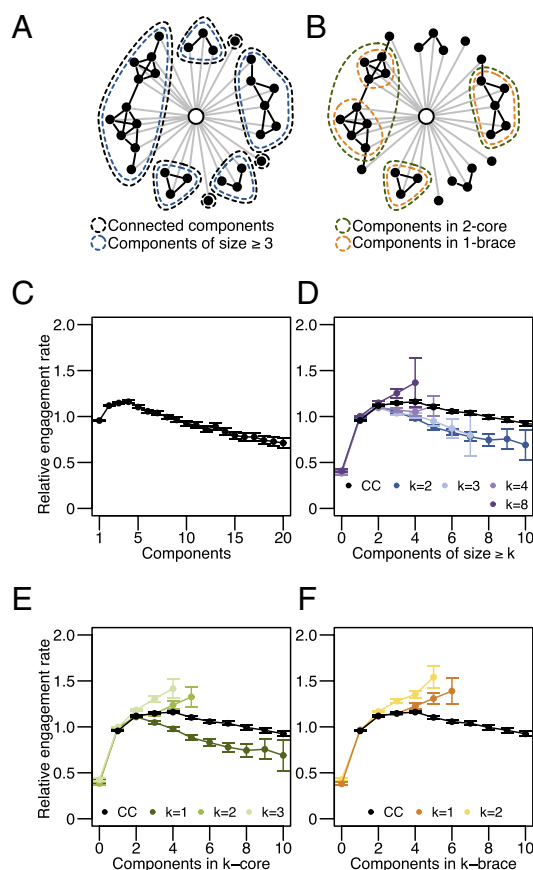


Fig. 4. Engagement and structural diversity for 50-node friendship neighborhoods. (A) Illustration of the connected components in a friendship neighborhood, delineating connected components and components of size ≥ 3 . (B) Illustration of the k -core and the k -brace, delineating the connected components of the 2-core and the 1-brace. (C) Engagement as a function of connected component count. (D) Engagement as a function of the number of components of size $\geq k$, for $k = 2, 3, 4, 8$, with connected component (CC) count shown for comparison. (E) Engagement as a function of k -core component count for $k = 1, 2, 3$, with CC count shown for comparison. (F) Engagement as a function of k -brace component count for $k = 1, 2$, with CC count shown for comparison. Engagement rates are reported on a relative scale, where 1.0 signifies the average conversion rate of all 50-node neighborhoods. All error bars are 95% confidence intervals. For other neighborhood sizes, see Fig. S4.

the inviter invited at most 20 e-mail addresses on the date of the invitation. This conditioning is meant to omit invitation batches where the inviter opted to "select all" within the contact import tool and focuses our investigation on socially selective invitations.

Invitations were sent during an 11-week period spanning July 12, 2010 to September 26, 2010. An e-mail address was considered to have converted to a registered user account if the address was registered for an account within 14 days of the invitation, counting both individuals who signed up via links provided in the invitation e-mail and users who signed up by visiting the Facebook website directly within 14 days. Only contact import events that occurred before the invitation event are considered. Likewise, only friendship edges that existed before the invitation event are considered to be part of the neighborhood.

Many of the findings we investigate are governed by complex nonlinear effects, which make traditional regression controls generally inadequate. In an attempt to control for confounding signals in our data, several parallel observation groups were maintained, against which all findings were validated. As a means of capturing potential artifacts from duplicitous private/business e-mail address use, a first such validation group was constructed by conditioning upon e-mail invitations sent to a small set of common and commonly private e-mail providers: Hotmail, Yahoo!, Gmail, AOL, and Yahoo! France. As a means of observing any differences between already established and growing Facebook markets, two parallel validation groups were constructed to observe established markets (United States) and emerging Facebook markets (Brazil, Germany, Japan, and Russia), classified by the most recently resolved country of login for the inviting Facebook account. Whereas invitation conversion rates were generally higher in emerging markets, none of the conditional datasets were observed to deviate from the complete dataset with regard to internal structural findings.

Highly sparse neighborhoods were a very common occurrence in these data, owing to the fact that the neighborhoods we study here are only partial observations of an individual's actual connection to Facebook. We are able to infer links only to those site users who have used the contact importer tool and maintain active e-mail communication with the e-mail address in question, criteria that induce a sampled subgraph that we then observe. The probability of sampling an edge uniformly at random in any neighborhood

with low edge density is therefore quite low, and the probability that all sampled nodes come from the same cluster within a clustered neighborhood is lower still. From the perspective of communication multiplexity (31), we should in fact expect that our randomly induced subgraph sample is biased toward strongly connected ties that tend to communicate on multiple mediums, but this expectation is not at issue with our results. The real matter of the fact is that contact neighborhoods where the induced subgraph consists of a single connected component are likely to come from very tightly connected neighborhood graphs.

Although the contact importer tool and invitation tool are prominently featured as part of the new user experience on Facebook, they are also heavily used by experienced users of the site: The median site age of an inviter in our dataset was 262 days. Although e-mail invitations constitute only a small portion of Facebook's growth, they provide a valuable window into the otherwise invisible growth process of the Facebook product.

For the analysis of photo co-tags, only co-tags since January 1, 2010 were considered.

Engagement Data Collection. We consider users *engaged* at a given time point if they have interacted with the application during at least 6 of the last 7 days. As with any measure of user behavior, this metric is a heuristic merely meant to approximate a broader notion of involvement on the site. Highly engaged users who do not access the Internet on weekends will never qualify as "six-plus engaged," whereas users who simply log in on a daily basis to check their messages will qualify. Our analysis is restricted to the population level, so such confounders are not a problem.

Due to the technical nature of how engagement data are stored at Facebook, it is impractical to retrieve six-plus engagement measures for dates exactly 3 months after registration. As an appropriate surrogate, we consider the six-plus engagement of users on the first day of their third calendar month as users.

ACKNOWLEDGMENTS. We thank M. Macy, J. Fowler, D. Watts, and S. Strogatz for comments. This research has been supported in part by a MacArthur Foundation Fellowship and National Science Foundation Grants IIS-0705774, IIS-0910664, CCF-0910940, and IIS-1016099.

- Pastor-Satorras R, Vespignani A (2001) Epidemic spreading in scale-free networks. *Phys Rev Lett* 86:3200–3203.
- Newman ME, Watts DJ, Strogatz SH (2002) Random graph models of social networks. *Proc Natl Acad Sci USA* 99(Suppl 1):2566–2572.
- Dodds PS, Watts DJ (2004) Universal behavior in a generalized model of contagion. *Phys Rev Lett* 92:218701.
- Backstrom L, Huttenlocher D, Kleinberg J, Lan X (2006) Group formation in large social networks: Membership, growth, and evolution. *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, eds Eliassi-Rad T, Ungar LH, Craven M, Gunopulos D (Association for Computing Machinery, New York), pp 44–54.
- Kearns M, Suri S, Montfort N (2006) An experimental study of the coloring problem on human subject networks. *Science* 313:824–827.
- Watts DJ, Dodds PS (2007) Influentials, networks, and public opinion formation. *J Consum Res* 34:441–458.
- Christakis NA, Fowler JH (2007) The spread of obesity in a large social network over 32 years. *N Engl J Med* 357:370–379.
- Sun E, Rosenn I, Marlow C, Lento T (2009) Gesundheit! Modeling contagion through Facebook news feed. *Proceedings of the AAAI International Conference on Weblogs and Social Media*, eds Adar E, et al. (Association for the Advancement of Artificial Intelligence, Menlo Park, CA), pp 146–153.
- Schelling T (1971) Dynamic models of segregation. *J Math Sociol* 1:143–186.
- Granovetter M (1978) Threshold models of collective action. *Am J Sociol* 83:1420–1443.
- Burt R (1987) Social contagion and innovation: Cohesion versus structural equivalence. *Am J Sociol* 92:1287–1335.
- Kossinets G, Watts DJ (2006) Empirical analysis of an evolving social network. *Science* 311:88–90.
- Centola D, Eguiluz V, Macy M (2007) Cascade dynamics of complex propagation. *Physica A* 374:449–456.
- Centola D, Macy M (2007) Complex contagions and the weakness of long ties. *Am J Sociol* 113:702–734.
- Palla G, Barabási AL, Vicsek T (2007) Quantifying social group evolution. *Nature* 446:664–667.
- Aral S, Muchnik L, Sundararajan A (2009) Distinguishing influence-based contagion from homophily-driven diffusion in dynamic networks. *Proc Natl Acad Sci USA* 106:21544–21549.
- Fowler JH, Christakis NA (2010) Cooperative behavior cascades in human social networks. *Proc Natl Acad Sci USA* 107:5334–5338.
- Simmel G (1955) *Conflict and the Web of Group Affiliations*, eds trans Wolff K, Bendix R (Free Press, Glencoe, IL).
- Granovetter M (1973) The strength of weak ties. *Am J Sociol* 78:1360–1380.
- Burt R (1992) *Structural Holes: The Social Structure of Competition* (Harvard Univ Press, Cambridge, MA).
- Crandall D, et al. (2010) Inferring social ties from geographic coincidences. *Proc Natl Acad Sci USA* 107:22436–22441.
- Bollobás B (2001) *Random Graphs* (Cambridge Univ Press, Cambridge, UK), 2nd Ed, p 150.
- Cohen JD (2008) Trusses: Cohesive subgraphs for social network analysis. *National Security Agency Technical Report* (National Security Agency, Fort Meade, MD).
- Luczak T (1991) Size and connectivity of the k-core of a random graph. *Discrete Math* 91:61–68.
- Janson S, Luczak MJ (2007) A simple solution to the k-core problem. *Random Struct Algo* 30:50–62.
- Alvarez-Hamelin JI, Dall'Astra L, Barrat A, Vespignani A (2006) Large scale networks fingerprinting and visualization using the k-core decomposition. *Adv Neural Inf Process Syst* 18:41–50.
- Carmi S, Havlin S, Kirkpatrick S, Shavitt Y, Shir E (2007) A model of Internet topology using k-shell decomposition. *Proc Natl Acad Sci USA* 104:11150–11154.
- Liggett T (1985) *Interacting Particle Systems* (Springer, Berlin).
- Durrett R (1995) *Ten Lectures on Particle Systems* (Springer, Berlin).
- Mossel E, Roch S (2007) On the submodularity of influence in social networks. *Proceedings of the ACM Symposium on Theory of Computing*, eds Johnson DS, Feige U (Association for Computing Machinery, New York), pp 128–134.
- Haythornthwaite C, Wellman B (1998) Work, friendship, and media use for information exchange in a networked organization. *J Am Soc Inf Sci* 49:1101–1114.