

Fast algorithm for detecting community structure in networks

M. E. J. Newman

Department of Physics and Center for the Study of Complex Systems, University of Michigan, Ann Arbor, Michigan 48109-1120, USA

(Received 22 September 2003; revised manuscript received 22 March 2004; published 18 June 2004)

Many networks display community structure—groups of vertices within which connections are dense but between which they are sparser—and sensitive computer algorithms have in recent years been developed for detecting this structure. These algorithms, however, are computationally demanding, which limits their application to small networks. Here we describe an algorithm which gives excellent results when tested on both computer-generated and real-world networks and is much faster, typically thousands of times faster, than previous algorithms. We give several example applications, including one to a collaboration network of more than 50 000 physicists.

DOI: 10.1103/PhysRevE.69.066133

PACS number(s): 89.75.Hc, 05.10.-a, 87.23.Ge, 89.20.Hh

I. INTRODUCTION

There has in recent years been a surge of interest within the physics community in the properties of networks of many kinds, including the Internet, the world wide web, citation networks, transportation networks, software call graphs, email networks, food webs, and social and biochemical networks [1–4]. One property that has attracted particular attention is that of “community structure”: the vertices in networks are often found to cluster into tightly knit groups with a high density of within-group edges and a lower density of between-group edges. Girvan and Newman [5,6] proposed a computer algorithm based on the iterative removal of edges with high “betweenness” scores that appears to identify such structure with some sensitivity, and this algorithm has been employed by a number of authors in the study of such diverse systems as networks of email messages, social networks of animals, collaborations of jazz musicians, metabolic networks, and gene networks [5–11]. As pointed out by Newman and Girvan [6], the principal disadvantage of their algorithm is the high computational demands it makes. In its simplest and fastest form, it runs in worst-case time $O(m^2n)$ on a network with m edges and n vertices, or $O(n^3)$ on a sparse network (one for which m scales with n in the limit of large n , which covers essentially all networks of current scientific interest, with the possible exception of food webs). With typical computer resources available at the time of writing, this limits the algorithm’s use to networks of a few thousand vertices at most, and substantially less than this for interactive applications. Increasingly, however, there is interest in the study of much larger networks; citation and collaboration networks can contain millions of vertices [12,13], for example, while the world wide web numbers in the billions [14].

In this paper, therefore, we propose another algorithm for detecting community structure. The algorithm operates on different principles from that of Girvan and Newman (GN), but, as we will show, gives qualitatively similar results. The worst-case running time of the algorithm is $O((m+n)n)$, or $O(n^2)$ on a sparse graph. In practice, it runs to completion on current computers in reasonable times for networks of up to a million or so vertices, bringing within reach the study of

communities in many systems that would previously have been considered intractable.

II. THE ALGORITHM

Our algorithm is based on the idea of modularity. Given any network, the GN community structure algorithm always produces *some* division of the vertices into communities, regardless of whether the network has any natural such division. To test whether a particular division is meaningful, we define a quality function or “modularity” Q as follows [6].

Let e_{ij} be one-half of the fraction of edges in the network that connect vertices in group i to those in group j , so that the total fraction of such edges is $e_{ij} + e_{ji}$. The only exception will be the diagonal elements e_{ii} , which are equal to the fraction of edges that fall within group i (with no factor of a half). Then $\sum_i e_{ii}$ is the total fraction of edges that fall within groups. All other edges fall between groups. The maximum value of this sum is 1, and a division of the network into communities is good if this quantity is large, meaning it is of order 1. On its own, however, the sum is not a good measure of community structure, since it takes its maximal value of 1 if we put all vertices in a single group together, which is a trivial and not particularly useful form of community structure.

A more useful measure of community structure is to calculate the sum $\sum_i e_{ii}$ and then subtract from it the value that it would take if edges were placed at random. Such a measure gives a score of zero to the trivial grouping with only a single community, but nonzero scores to nontrivial groupings.

Let a_i be the fraction of all *ends* of edges that are attached to vertices in group i . We can calculate a_i straightforwardly by noting that $a_i = \sum_j e_{ij}$. If the ends of edges are connected together at random, the fraction of the resulting edges that connect vertices within group i is a_i^2 . We define the modularity to be

$$Q = \sum_i (e_{ii} - a_i^2). \quad (1)$$

If a particular division gives no more within-community edges than would be expected by random chance, this modu-

larity is $Q=0$. Values other than 0 indicate deviations from randomness, and in practice values greater than about 0.3 appear to indicate significant community structure. A number of examples are given in Ref. [6].

But this now suggests an alternative approach to finding community structure. If a high value of Q represents a good community division, why not simply optimize Q over all possible divisions to find the best one? By doing this, we can avoid the iterative removal of edges and cut straight to the chase. The problem is that true optimization of Q is very costly. The number of ways to divide n vertices into g non-empty groups is given by the Stirling number of the second kind $S_n^{(g)}$, and hence the number of distinct community divisions is $\sum_{g=1}^n S_n^{(g)}$. This sum is not known in closed form, but we observe that $S_n^{(1)} + S_n^{(2)} = 2^{n-1}$ for all $n > 1$, so that the sum must increase at least exponentially in n . To carry out an exhaustive search of all possible divisions for the optimal value of Q would therefore take at least an exponential amount of time, and is in practice infeasible for systems larger than 20 or 30 vertices. Various approximate optimization methods are available: simulated annealing, genetic algorithms, and so forth. Here we consider a scheme based on a standard “greedy” optimization algorithm, which appears to perform well.

Our algorithm falls in the general category of agglomerative hierarchical clustering methods [15,16]. Starting with a state in which each vertex is the sole member of one of n communities, we repeatedly join communities together in pairs, choosing at each step the join that results in the greatest increase (or smallest decrease) in Q . The progress of the algorithm can be represented as a “dendrogram,” a tree that shows the order of the joins (see Fig. 2 for an example). Cuts through this dendrogram at different levels give divisions of the network into larger or smaller numbers of communities and we can select the best cut by looking for the maximal value of Q .

Since the joining of a pair of communities between which there are no edges at all can never result in an increase in Q , we need only consider those pairs between which there are edges, of which there will at any time be at most m , where m is again the number of edges in the graph. The change in Q upon joining two communities is given by

$$\Delta Q = e_{ij} + e_{ji} - 2a_i a_j = 2(e_{ij} - a_i a_j), \quad (2)$$

which can clearly be calculated in constant time. The quantities e_{ij} are initially equal to one-half of the corresponding elements of the adjacency matrix of the network, i.e., to $\frac{1}{2}$ for vertex pairs that are joined by an edge and 0 for those that are not. Following a join, some of the matrix elements e_{ij} must be updated by adding together the rows and columns corresponding to the joined communities, which takes worst-case time $O(n)$. Thus each step of the algorithm takes worst-case time $O(m+n)$. There are a maximum of $n-1$ join operations necessary to construct the complete dendrogram and hence the entire algorithm runs in time $O((m+n)n)$, or $O(n^2)$ on a sparse graph. The algorithm has the added advantage of calculating the value of Q as it goes along, making it especially simple to find the optimal community structure.

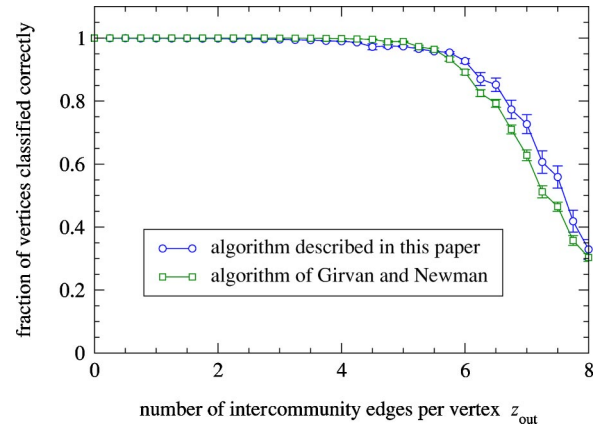


FIG. 1. The fraction of vertices correctly identified by our algorithms in the computer-generated graphs described in the text. The two curves show results for our algorithm (circles) and for the algorithm of Girvan and Newman [5] (squares). Each point is an average over 100 graphs.

It is worth noting that our algorithm can be generalized trivially to weighted networks in which each edge has a numeric strength associated with it, by making the initial values of the matrix elements e_{ij} equal to (a half of) those strengths; otherwise the algorithm is as above and has the same running time. The networks studied in this paper, however, are all unweighted.

III. APPLICATIONS

As a first example of the working of our algorithm, we have generated using a computer a large number of random graphs with known community structure, which we then run through the algorithm to quantify its performance. Each graph consists of $n=128$ vertices divided into four groups of 32. Each vertex has on average z_{in} edges connecting it to members of the same group and z_{out} edges to members of other groups, with z_{in} and z_{out} chosen such that the total expected degree $z_{in} + z_{out} = 16$, in this case. As z_{out} is increased from small values, the resulting community structure becomes progressively weaker and the graphs pose greater and greater challenges to the community-finding algorithm. In Fig. 1 we show the fraction of vertices correctly assigned to the four communities by the algorithm as a function of z_{out} [19]. As the figure shows, the algorithm performs well, correctly identifying more than 90% of vertices for values of $z_{out} \leq 6$. Only when z_{out} approaches the value 8 at which the number of within- and between-community edges per vertex is the same does the algorithm begin to fail. On the same plot we also show the performance of the GN algorithm and, as we can see, that algorithm performs slightly but measurably better for smaller values of z_{out} . For example, for $z_{out}=5$ our algorithm correctly identifies an average of 97.4(2)% of vertices, while the older algorithm correctly identifies 98.9(1)%. Both, however, clearly perform well.

Interestingly, for higher values of z_{out} our algorithm performs better than the older one, and we have come across a few real-world networks in which this is the case also. Nor-

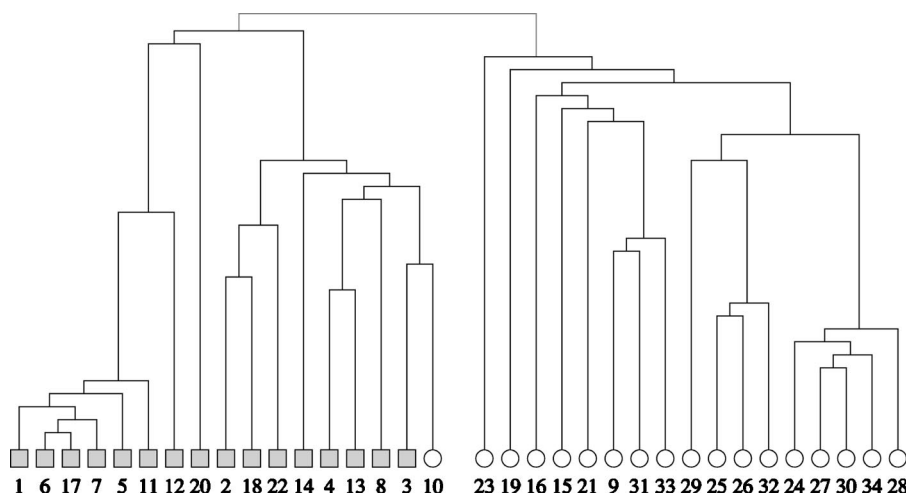


FIG. 2. Dendrogram of the communities found by our algorithm in the “karate club” network of Zachary [5,17]. The shapes of the vertices represent the two groups into which the club split as the result of an internal dispute.

mally, however, the GN algorithm seems to have the edge, and this should come as no great surprise. Our algorithm bases its decisions on purely local information about individual communities, while the GN algorithm uses nonlocal information about the entire network—information derived from betweenness scores. Since community structure is itself fundamentally a nonlocal quantity, it seems reasonable that one can do a better job of finding that structure if one has nonlocal information at one’s disposal.

For systems small enough that the GN algorithm is computationally tractable, therefore, we see no reason not to continue using it—it appears to give the best results. For systems too large to make use of this approach, however, our algorithm gives useful community structure information with comparatively little effort.

We have applied our algorithm to a variety of real-world networks also. We have looked, for example, at the “karate club” network studied in [5], which represents friendships between 34 members of a club at a U.S. university, as recorded over a two-year period by Zachary [17]. During the course of the study, the club split into two groups as a result of a dispute within the organization, and the members of one group left to start their own club. In Fig. 2 we show the dendrogram derived by feeding the friendship network into our algorithm. The peak modularity is $Q=0.381$ and corresponds to a split into two groups of 17, as shown in the figure. The shapes of the vertices represent the alignments of the club members following the dispute and, as we can see, the division found by the algorithm corresponds almost perfectly to these alignments; only one vertex, number 10, is classified wrongly. The GN algorithm performs similarly on this task, but not better—it also finds the split but classifies one vertex wrongly (although a different one, vertex 3). In other tests, we find that our algorithm also successfully detects the main two-way division of the dolphin social network of Lusseau [6,18], and the division between black and white musicians in the jazz network of Gleiser and Danon [11].

As a demonstration of how our algorithm can sometimes miss some of the structure in a network, we take another example from Ref. [5], a network representing the schedule of games between American college football teams in a single season. Because the teams are divided into groups or

“conferences,” with intraconference games being more frequent than interconference games, we have a reasonable idea ahead of time about what communities our algorithm should find. The dendrogram generated by the algorithm is shown in Fig. 3, and has an optimal modularity of $Q=0.546$, which is a little shy of the value 0.601 for the best split reported in [5]. As the dendrogram reveals, the algorithm finds six communities. Some of them correspond to single conferences, but most correspond to two or more. The GN algorithm, by contrast, finds all 11 conferences, as well as accurately identifying independent teams that belong to no conference. Nonetheless, it is clear that our algorithm is quite capable of picking out useful community structure from the network, and of course it is much the faster algorithm. On the author’s desktop computer the algorithm ran to completion in an immeasurably small time—less than a hundredth of a second. The algorithm of Girvan and Newman took a little over a second.

A time difference of this magnitude will not present a big problem in most practical situations, but performance rapidly becomes an issue when we look at larger networks; we expect the ratio of running times to increase with the number of vertices. Thus, for example, in applying our algorithm to the 1275-node network of jazz musician collaborations mentioned above, we found that it runs to completion in about one second of CPU time. The GN algorithm by contrast takes more than three hours to reach very similar results.

As an example of an analysis made possible by the speed of our algorithm, we have looked at a network of collaborations between physicists as documented by papers posted on the widely used Physics E-print Archive at arxiv.org. The network is an updated version of the one described in Ref. [13], in which scientists are considered connected if they have coauthored one or more papers posted on the archive. We analyze only the largest component of the network, which contains $n=56\,276$ scientists in all branches of physics covered by the archive. Since two vertices that are unconnected by any path are never put in the same community by our algorithm, the small fraction of vertices that are not part of the largest component can safely be assumed to be in separate communities in the sense of our algorithm. Our algorithm takes 42 min to find the full community structure. Our best estimates indicate that the GN algorithm would take

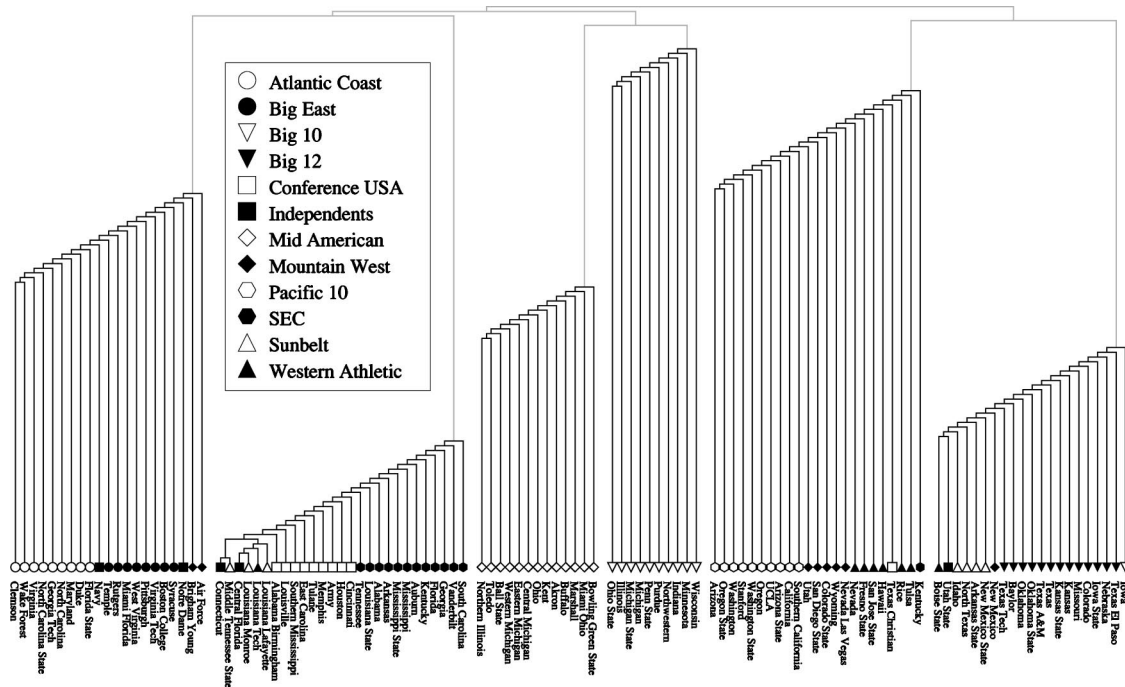


FIG. 3. Dendrogram of the communities found in the college football network described in the text. The real-world communities—conferences—are denoted by the different shapes as indicated in the legend.

somewhere between three and five years to complete its version of the same calculation.

The analysis reveals that the network in question consists of about 600 communities, with a high peak modularity of $Q=0.713$, indicating strong community structure in the physics world. Four of the communities found are large, containing between them 77% of all the vertices, while the others are small—see Fig. 4, left panel. The four large communities correspond closely to subject subareas: one to astrophysics, one to high-energy physics, and two to condensed-matter physics. Thus there appears to be a strong correlation between the structure found by our algorithm and the community divisions perceived by human observers. It is precisely

correlation of this kind that makes community structure analysis a useful tool in understanding the behavior of networked systems.

We can repeat the analysis with any of the subcommunities to observe how they break up. For example, feeding the smaller of the two condensed-matter groups through the algorithm again, we find an even stronger peak modularity of $Q=0.807$ —the strongest we have yet observed in any network—corresponding to a split into about a 100 communities (Fig. 4, center panel). These communities have a broad distribution of sizes from 3 to nearly 2000. The distribution is shown in cumulative form in Fig. 5, and we observe that it is approximately power law in form with exponent about

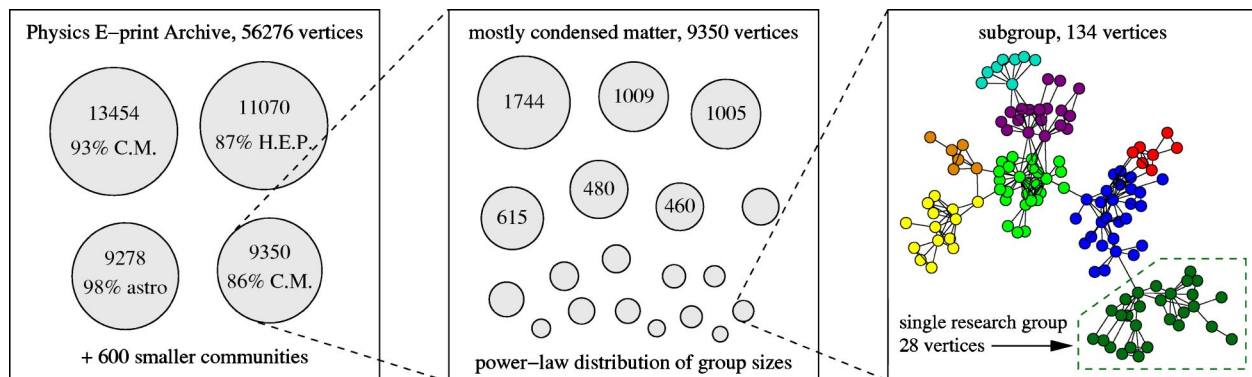


FIG. 4. Left panel: Community structure in the collaboration network of physicists. The graph breaks down into four large groups, each composed primarily of physicists of one specialty, as shown. Specialties are determined by the subsection(s) of the e-print archive in which individuals post papers: “C.M.” indicates condensed matter; “H.E.P.” indicates high-energy physics including theory, phenomenology, and nuclear physics; “astro” indicates astrophysics. Middle panel: one of the condensed matter communities is further broken down by the algorithm, revealing an approximate power-law distribution of community sizes. Right panel: one of these smaller communities is further analyzed to reveal individual research groups (different shades), one of which (in the dashed box) is the author’s own.

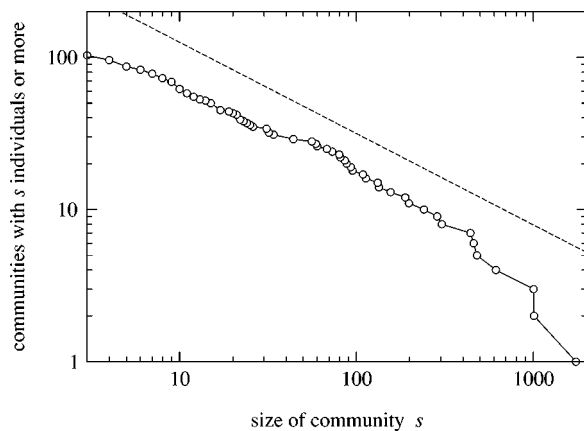


FIG. 5. Cumulative distribution function of the sizes of communities found in one of the subnetworks of the physics collaboration graph, as described in the text. The dotted line represents the slope the plot would have if the distribution followed a power law with exponent -1.6 .

-1.6 , although this conclusion should be treated with caution as there is significant deviation from a perfect power law [20].

Narrowing our focus still further to the particular one of these communities that contains the present author, we find the structure shown in the right panel of Fig. 4. Feeding this one last time through the algorithm, it breaks apart into communities that correspond closely to individual institutional research groups, the author's group appearing in the corner of the figure, highlighted by the dashed box. One could pur-

sue this line of analysis further, identifying individual groups, iteratively breaking them down, and looking, for example, at the patterns of collaboration between them, but we leave this for later studies.

IV. CONCLUSIONS

In this paper we have described an algorithm for extracting community structure from networks, which has a considerable speed advantage over previous algorithms, running to completion in a time that scales as the square of the network size. This allows us to study much larger systems than has previously been possible. Among other examples, we have applied the algorithm to a network of collaborations between more than 50 000 physicists, and found that the resulting community structure corresponds closely to the traditional divisions between specialties and research groups in the field.

We believe that our method will not only allow for the extension of community structure analysis to some of the very large networks that are now being studied for the first time, but will also provide a useful tool for visualizing and understanding the structure of these networks, whose daunting size has hitherto made many of their structural properties obscure.

ACKNOWLEDGMENTS

The author thanks Leon Danon, Pablo Gleiser, David Lusseau, and Douglas White for providing network data used in the examples. This work was supported in part by the National Science Foundation under Grant No. DMS-0234188.

-
- [1] S. H. Strogatz, *Nature (London)* **410**, 268 (2001).
 - [2] R. Albert and A.-L. Barabási, *Rev. Mod. Phys.* **74**, 47 (2002).
 - [3] S. N. Dorogovtsev and J. F. F. Mendes, *Evolution of Networks: From Biological Nets to the Internet and WWW* (Oxford University Press, Oxford, 2003).
 - [4] M. E. J. Newman, *SIAM Rev.* **45**, 167 (2003).
 - [5] M. Girvan and M. E. J. Newman, *Proc. Natl. Acad. Sci. U.S.A.* **99**, 7821 (2002).
 - [6] M. E. J. Newman and M. Girvan, *Phys. Rev. E* **69**, 026113 (2004).
 - [7] D. Wilkinson and B. A. Huberman, *Proc. Natl. Acad. Sci. U.S.A.* **101**, 5241 (2004).
 - [8] P. Holme, M. Huss, and H. Jeong, *Bioinformatics* **19**, 532 (2003).
 - [9] R. Guimerà, L. Danon, A. Díaz-Guilera, F. Giralt, and A. Arenas, *Phys. Rev. E* **68**, 065103 (2003).
 - [10] J. R. Tyler, D. M. Wilkinson, and B. A. Huberman, in *Proceedings of the First International Conference on Communities and Technologies*, edited by M. Huysman, E. Wenger, and V. Wulf (Kluwer, Dordrecht, 2003).
 - [11] P. Gleiser and L. Danon, *Adv. Complex Syst.* **6**, 565 (2003).
 - [12] S. Redner, *Eur. Phys. J. B* **4**, 131 (1998).
 - [13] M. E. J. Newman, *Proc. Natl. Acad. Sci. U.S.A.* **98**, 404 (2001).
 - [14] J. Kleinberg and S. Lawrence, *Science* **294**, 1849 (2001).
 - [15] B. Everitt, *Cluster Analysis* (John Wiley, New York, 1974).
 - [16] J. Scott, *Social Network Analysis: A Handbook*, 2nd ed. (Sage, London, 2000).
 - [17] W. W. Zachary, *J. Anthropol. Res.* **33**, 452 (1977).
 - [18] D. Lusseau, *Proc. R. Soc. London, Ser. B* **270**, S186 (2003).
 - [19] The criterion for deciding correct classification is as follows. We find the largest set of vertices that are grouped together by the algorithm in each of the four known communities. If the algorithm puts two or more of these sets in the same group, then all vertices in those sets are considered incorrectly classified. Otherwise, they are considered correctly classified. All other vertices not in the largest sets are considered incorrectly classified. This criterion is quite harsh—there are cases in which one might consider some of the vertices to have been identified correctly, where this method would not. Even with this harsh definition, however, our algorithm performs well, and a laxer definition would only make its performance more impressive.
 - [20] This power law is different from the one observed in an email network by Guimerà *et al.* [9]. They studied the histogram of community sizes over all levels of the dendrogram; we are looking only at the single level corresponding to the maximum value of Q .