

Coauthorship networks and patterns of scientific collaboration

M. E. J. Newman*

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

By using data from three bibliographic databases in biology, physics, and mathematics, respectively, networks are constructed in which the nodes are scientists, and two scientists are connected if they have coauthored a paper. We use these networks to answer a broad variety of questions about collaboration patterns, such as the numbers of papers authors write, how many people they write them with, what the typical distance between scientists is through the network, and how patterns of collaboration vary between subjects and over time. We also summarize a number of recent results by other authors on coauthorship patterns.

It has long been realized that the coauthorship of articles in learned journals provides a window on patterns of collaboration within the academic community. Coauthorship of a paper can be thought of as documenting a collaboration between two or more authors, and these collaborations form a “coauthorship network,” such as that depicted in Fig. 1, in which the network nodes represent authors, and two authors are connected by a line if they have coauthored one or more papers. The structure of such networks turns out to reveal many interesting features of academic communities.

Networks are not new to bibliometrics; the field has a long history of the study of citation networks (1, 2), the networks formed by the citations between papers. These are quite distinct, however, from coauthorship networks; the nodes in a citation network are papers, not authors, and the links between them are citations, not coauthorship. The coauthorship network is as much a network depicting academic society as it is a network depicting the structure of our knowledge. And, perhaps because of this, it has received far less attention than have citation networks. Nonetheless, it has much of value to tell us, as recent work has shown.

During the 1990s (and possibly earlier), a number of authors pointed out the potential utility of coauthorship data and in some cases performed small-scale statistical analyses of such things as frequency of coauthored articles by particular authors or authors at particular institutions (3–7). But it was with the advent of comprehensive online bibliographies that construction of complete or near-complete coauthorship networks for entire fields became a realistic possibility. Starting around 2000, several researchers began the construction of large-scale networks representing research in mathematics (8–10), biology, physics, and computer science (11) and neuroscience (10).

In this paper, we look in detail at three particular networks of scientific collaborations and describe some of the patterns they reveal. The networks are:

(i) A network of coauthorships of papers in the Medline bibliographical database from 1995 to 1999, inclusive. Medline is a widely used and compendious database of papers covering biomedical research. Biomedical research accounts for the largest part of civilian scientific research by far, dwarfing research in all other subjects put together in terms of expenditure. Any study that excluded biomedicine could not claim to be representative of science as it is practiced today.

(ii) A network of coauthorships of physicists assembled from papers posted on the widely used Physics E-print Archive at Cornell University (formerly at the Los Alamos National Laboratory) between 1995 and 1999. Physics has led the way in moving from journal publication to author self-publication in online preprint databases, with preprint publication largely replacing journal publication in some subfields. Preprint databases provide a useful source of up-to-the-minute publication records, although their coverage is less complete than that of professionally maintained databases like Medline.

(iii) A collaboration network of mathematicians compiled from databases maintained by the journal *Mathematical Reviews*. Of the networks yet studied, this is probably the most complete and accurate, covering the period from 1940 to the present without any break.

Networks *i* and *ii* were constructed by the author from bibliographic data supplied by the maintainers of the corresponding databases. Network *iii* was constructed by J. Grossman and P. Ion (8) and graciously supplied by J. Grossman.

A number of other papers in this volume describe bibliometric studies of a database of papers that appeared in PNAS over the period 1997–2002. Although it would be possible to construct a coauthorship network from these data, such a network would be less satisfactory for the study of collaboration patterns than the networks studied here. Most authors publish in more than one journal, so that data on publications in a single journal would give an incomplete picture of their authorship patterns. The databases studied in this paper are more complete, although certainly they do not claim to document every paper.

The outline of this paper is as follows. In *Statistical Properties of Coauthorship Networks*, we describe a variety of results derived from analysis of our networks and highlight some differences among the three subjects studied. In *Additional Results*, we summarize some recent additional results obtained by using the same or additional data, including a number of results due to other authors. In *Conclusion*, we give our conclusions. Many of the results reported here have appeared previously in refs. 11–15, as well as a number of other papers, which are cited as appropriate.

Statistical Properties of Coauthorship Networks

A summary of the basic statistics of the three networks studied here is given in Table 1. The largest of the networks, not surprisingly, is the biomedical network, with >1.5 million authors over a 5-year period. Even the mathematics network, which covers a much longer period (≈ 60 years), comes nowhere close to this size. Clearly, biology dwarfs other subjects in terms not only of spending but also of manpower. The number of papers shows a similar pattern, although we do not have precise data for

This paper results from the Arthur M. Sackler Colloquium of the National Academy of Sciences, “Mapping Knowledge Domains,” held May 9–11, 2003, at the Arnold and Mabel Beckman Center of the National Academies of Sciences and Engineering in Irvine, CA.

*E-mail: mejn@umich.edu.

© 2004 by The National Academy of Sciences of the USA

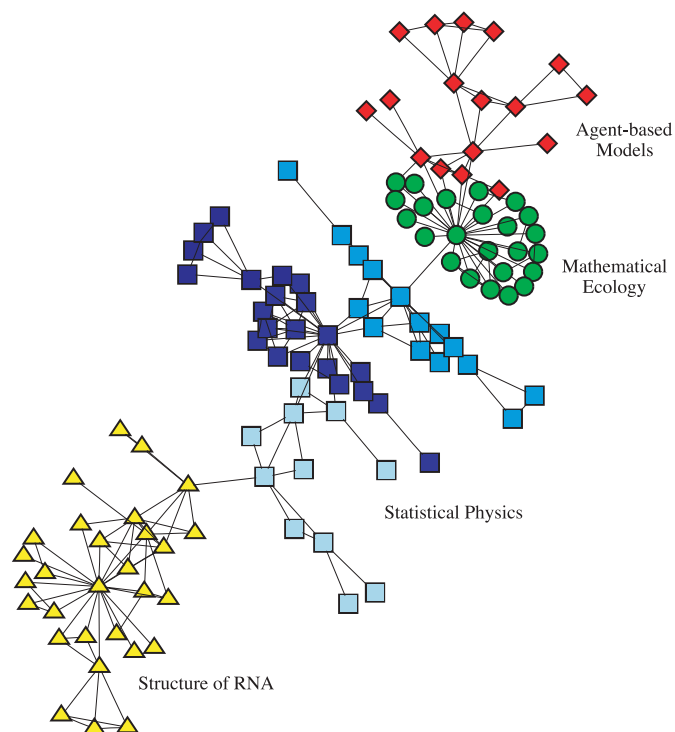


Fig. 1. An example of a small coauthorship network depicting collaborations among scientists at a private research institution. Nodes in the network represent scientists, and a line between two of them indicates they coauthored a paper during the period of study. This particular network appears to divide into a number of subcommunities, as indicated by the shapes of the nodes, and these subcommunities correspond roughly to topics of research, as discussed by Girvan and Newman (37).

the number of papers in the mathematics database. [Grossman (9) cites a figure of “about 1.6 million authored items” in the *Mathematical Reviews* database, for a slightly more recent version of the network than that studied here.]

The number of papers per author is similar across the three subject areas, between five and seven in each case. Because the mathematics database covers a longer time period, however, this may indicate that mathematicians are producing fewer papers than their more empirically minded colleagues in the sciences. Scientific productivity, measured by number of papers authored, has had a long history of study in bibliometrics, with the articles by Lotka (16) and Shockley (17) being famous early examples. Both of these authors found that the number of papers produced by scientists had a “fat-tailed” distribution, in which a small number of scientists produced a very large number of papers, a result that has since been confirmed by others (18, 19), and which is seen in our own data as well (11, 12).

The number of authors per paper, by contrast, varies substantially among the subjects studied, with biology having the largest number and mathematics the smallest. This presumably reflects real differences in the way research is done in these fields, with biological research consisting often of work by large groups of laboratory scientists and mathematics consisting of theoretical work done primarily by individuals alone or by pairs of collaborators. Grossman (9) says that 66% of mathematics papers are written by a single author (although this number changes over time; see *Additional Results*). In the Medline database, the corresponding figure is 21%. These figures may offer some explanation for the possible lower productivity of mathematics in terms of papers published per unit time: with fewer coauthors

Table 1. Summary statistics for the three coauthorship networks analyzed here

	Biology	Physics	Mathematics
Number of authors	1,520,251	52,909	253,339
Number of papers	2,163,923	98,502	—
Papers per author	6.4	5.1	6.9
Authors per paper	3.75	2.53	1.45
Average collaborators	18.1	9.7	3.9
Largest component	92%	85%	82%
Average distance	4.6	5.9	7.6
Largest distance	24	20	27
Clustering coefficient	0.066	0.43	0.15
Assortativity	0.13	0.36	0.12

The statistics are, from top to bottom, total number of authors appearing in the corresponding databases; total number of papers appearing; mean number of papers published by an author; mean number of coauthors on a paper; mean number of different individuals an author collaborated with; largest connected group of individuals in the network; mean vertex-vertex distance between connected individuals in the network; largest such distance; the clustering coefficient, which is the mean probability that two coauthors will also be coauthors of one another; and the degree assortativity coefficient, which is the Pearson correlation coefficient of the degrees (i.e., number of collaborators) of adjacent vertices in the network. The material shown here is after Newman (12) and Grossman (9).

on most publications, the production of a mathematics paper involves more work per author.

A similar pattern is revealed in the average number of collaborators an individual has in the three fields, which is more than four times higher in biology than in mathematics. This again is presumably a result of different modes of research, with biology being primarily experimental, mathematics being entirely theoretical, and physics being a combination of the two. [The quintessential example of scientific experiment on an industrial scale is high-energy physics, for which it was previously found, using the SPIRES (www.slac.stanford.edu/spires) database of high-energy physics papers, that authors had an amazing 173 collaborators on average over the 5-year period from 1995 to 1999 (11).]

In addition to mean numbers of papers and coauthors, one can look at the distributions of these quantities. In Fig. 2, for instance, we show the distributions of the number of coauthors that scientists have for the three subjects. The distributions are quite similar, although the distribution for biomedicine (circles)

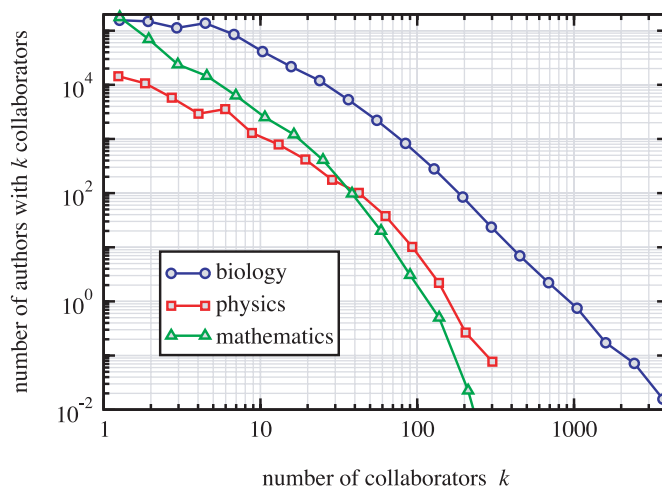


Fig. 2. Histograms of the distribution of numbers of collaborators for scientists in each of three fields studied.

has a longer tail, reflecting the higher mean number of collaborators that individuals have in that field. In each case, the distribution is fat-tailed, like the distribution of number of papers written by scientists mentioned above, with a small fraction of scientists having a very large number of collaborators, up to thousands in the case of the biology network. (Recall that the data for this network cover only a 5-year period; publishing papers with 1,000 coauthors in <2,000 days is an impressive achievement by any measure.) Unlike some other networks, such as the World Wide Web and the Internet, however, the distributions for these networks do not follow power laws; they are not “scale-free networks,” in the jargon of the field. It has been suggested that the distributions are actually power law in form with an exponential cutoff (9, 11, 20), and this appears to be a reasonable fit to the data. The cutoff may be produced by the finite time window used in the study, a hypothesis that could, in principle, be tested by varying the size of the window, although we do not do that here.

Table 1 also gives the size of the largest component in each of the networks. A component is a set of network nodes connected via coauthorship, such that any node in the set can be reached from any other by traversing a suitable path of intermediate collaborators. For each of the networks studied here, the largest component fills most of the network, occupying 82–92% of the network in the three cases. Thus a large portion of each of these communities is connected in a kind of linked research enterprise rather than working separately in isolation. Overall, this seems a promising picture; intellectual isolation from the mainstream of one’s research area cannot often be a good thing. Most scientists who do not belong to the largest component are members of small disconnected components containing only a handful of others.

Many recent studies of networks of various types have focused on network distance between nodes. This distance is defined as the number of “hops” along links in the network that one needs to make to move from one given node to another. A pair of individuals who have coauthored a paper, for instance, are distance 1 apart, whereas a pair who have not done so but who share a common coauthor are distance 2 apart, and so forth. In the late 1950s, Kochen and Pool (21) speculated on mathematical grounds that networks might show surprisingly small typical distances between pairs of nodes, and in a famous experiment some years later, Milgram (22, 23) demonstrated that this was the case for acquaintance networks, at least in the U.S. Our coauthorship networks appear to follow the same pattern. We calculate the distance between all pairs of individuals in a network using a breadth-first search or “burning” algorithm (24) and then take the average to give the figures shown in Table 1. For each of the networks, the result is very small, at least compared with the size of the network. Mathematics has the largest mean distance, possibly again as a result of the relative sparsity of mathematics collaborations, but even its value of 7.6 is tiny compared with the quarter of a million mathematicians in the network. This appears to indicate a close-knit cohesive community in which most people are connected not only by some path through the network but also by a short one.

A certain amount of attention has focused also on the distances from particular scientists to others in the coauthorship network. Mathematicians have long discussed the “Erdős number,” the distance through the mathematics network from a given mathematician to Paul Erdős, an influential Hungarian number theorist of the 20th century who was renowned for his prolific publication and collaboration. Erdős numbers have been studied in depth by a number of authors by using the *Mathematical Reviews* data (8, 9, 25, 26). It is found, for instance, that the mean distance from Paul Erdős to other mathematicians is much lower than the mean distance in the network as a whole, taking a value ≈ 4.7 . [Mean distances from other individuals, most of which are significantly higher than Erdős’ mean, are sometimes called “Doe numbers”

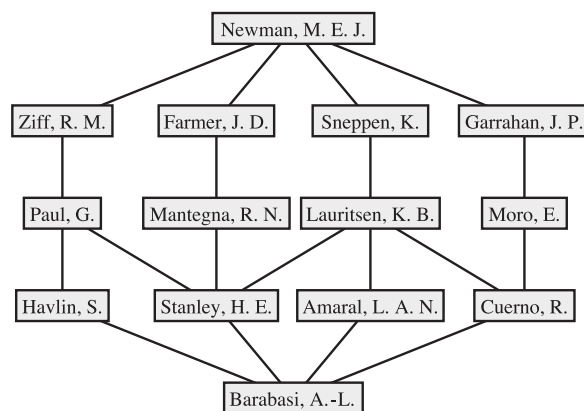


Fig. 3. The shortest paths through the collaboration network of physics papers from the author of this paper to A.-L. Barabási, who also publishes on networks.

(9).] The largest distance in a network, which is called the “diameter,” is also occasionally of interest; it is between 20 and 30 for each of the networks studied here.

It is straightforward to create a computer algorithm to find the shortest path between two particular scientists, again using breadth-first search, and it has been suggested by Kautz *et al.* (27) that such algorithms could be of use for providing “referral chains,” links of acquaintances that individuals could use to establish contact with other scientists. In the simplest case, for example, it might be useful to know that one shared a common collaborator with another scientist if one wished to arrange an introduction. Note that the shortest path between two individuals need not be unique, and in fact, it happens quite frequently that there are two or more shortest paths of equal length. Fig. 3 shows shortest paths in the network of the Physics E-print Archive between the present author and A.-L. Barabási of the University of Notre Dame, who also publishes on networks. As Fig. 3 shows, there are several different paths from one scientist to the other, all with length four. This particular case is interesting, because it shows that scientists working in the same field need not be linked through others in their field. The shortest paths in this case are established via my collaborations with J. D. Farmer, J. P. Garrahan, K. Sneppen, and R. M. Ziff, only the last of which collaborations involved work on networks (and then only peripherally).

Another interesting network measure related to shortest paths has been suggested by S. H. Strogatz (personal communication), who asks how many of the shortest paths from a particular individual to others pass through each of their collaborators. Is it the case that most of our connections to others are via just one or two of our best-connected collaborators, or are they distributed evenly among our collaborators? For the networks studied here, it turns out that the former is the case, as is evident in Fig. 4, which shows for the physics network what percentage of shortest paths pass through each of a scientist’s coauthors, on average. Thus, Fig. 4 reveals that on average $\approx 64\%$ of an individual’s shortest paths to others pass through the best-connected of their collaborators, and most of the remainder pass through the next-best connected. This may indicate that a small number of scientists are playing the role of broker for communications among others. (See also the discussion of betweenness centrality in *Additional Results*.)

Two other quantities of interest, both previously studied for many networks, are also given in Table 1. The first is the “clustering coefficient” (28), which measures network “clustering” or “transitivity,” the probability that two of a scientist’s coauthors have themselves coauthored a paper. In topological terms, it is a measure of the density of triangles in a network, a

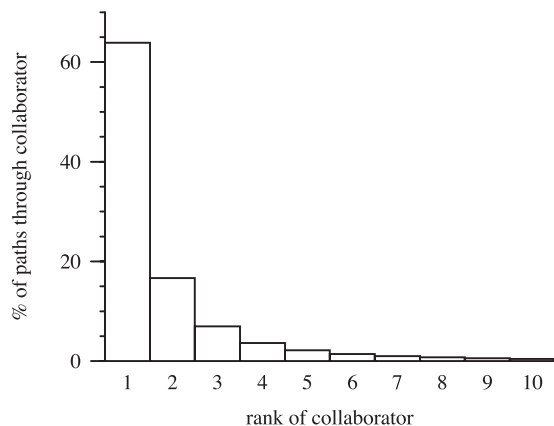


Fig. 4. The average percentage of paths from other scientists to a given scientist that pass through each collaborator of that scientist, ranked in decreasing order. The plot is for the physics network, although similar results are found for the others. [Reproduced with permission from ref. 13 (Copyright 2001, American Physical Society)].

triangle being formed every time two of one's collaborators collaborate with each other. The clustering coefficient is highest for physics (43%) and lowest for biology (7%), and it is unclear why there is so much variation among fields. Presumably the numbers reflect substantial differences among collaboration patterns in the sciences, but what these differences are is far from obvious. Part of the clustering in each network can be accounted for by papers with three or more coauthors. Such papers introduce triangles of collaborating authors and hence increase the clustering coefficient. This effect can account for only about one-half of the clustering seen in coauthorship networks, however (29); the rest must be due to sociological or organizational effects of some kind.

An alternative way to measure the clustering effect is to look at the time evolution of a network. Among social networks, coauthorship networks are unusual in having well-documented time evolution. Because each paper comes with a date of publication or submission, we can say approximately when each connection was added to the network, and so we can reconstruct the order in which the network grew. This allows us to ask the probability of two scientists coauthoring a paper, given that they have a third mutual collaborator and have not collaborated in the past. By studying only scientists who have not previously collaborated, we eliminate any bias introduced by papers with three or more coauthors. In ref. 14, we showed that scientists with a single mutual collaborator are ≈ 45 times more likely to coauthor a paper than those with no mutual collaborators. Those with two are >100 times as likely to coauthor a paper.

The last line in Table 1 gives the "assortativity coefficient" or degree correlation for the networks (15). This is the correlation coefficient for the number of collaborators that coauthors have. It lies in the range of -1 to 1 , with positive values indicating that people with many collaborators tend to collaborate with others who have many collaborators and negative values indicating the reverse. The coefficient is positive for all networks studied, indicating that the most gregarious scientists tend to be connected to each other. Again, it is an open question why this should be the case.

It is also possible to extract from coauthorship data a measure of the strength of the collaboration between pairs of individuals. The simplest such measure would be just a count of the frequency with which two scientists have coauthored papers, the number of coauthored papers over a given interval, for instance. However, this fails to take into account the number of other coauthors on each paper. Presumably two authors who collab-

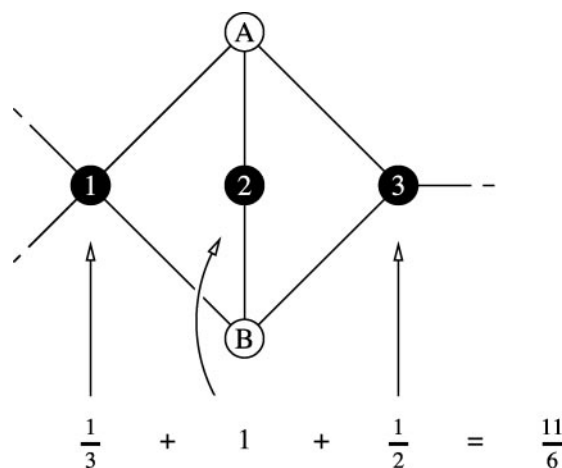


Fig. 5. Authors A and B have coauthored three papers, labeled 1, 2, and 3, which had, respectively, four, two, and three authors. The tie between A and B accordingly accrues weight $\frac{1}{3}$, 1 , and $\frac{1}{2}$ from the three papers, for a total weight of $\frac{11}{6}$. [Reproduced with permission from ref. 13 (Copyright 2001, American Physical Society)].

orate on a 10-author paper are, in general, working less closely with one another than two who produce a two-author paper with no other help. To account for this effect, we proposed in ref. 13 the measure of collaboration strength illustrated in Fig. 5. Each paper coauthored by a given author pair adds an amount $1/(n-1)$ to the strength of their collaboration, where n is the total number of authors on the paper. The rationale behind this choice is that an author divides his/her time between the $n-1$ other authors with whom he/she works on a paper, and hence the strength of the connection to each of them varies inversely as $n-1$. As an example of this measure, we show in Fig. 6 the coauthors of G. Barkema, one of the author's most frequent collaborators, with line thickness used to indicate the strength of the connections. Clearly, there is considerable variation in connection strength. Over the entire physics network connection, strengths range from a maximum of 34.0 to a minimum of ≈ 0.01 .

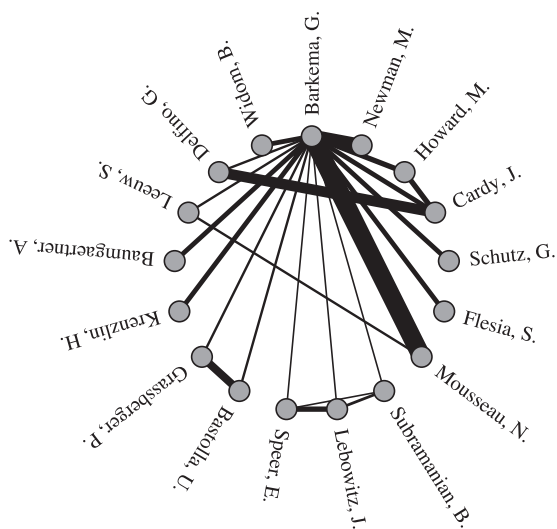


Fig. 6. G. Barkema and collaborators, with lines representing collaborations whose thickness is proportional to the estimate of collaboration strength defined in the text and illustrated in Fig. 5.

Additional Results

The network of collaborations of mathematicians compiled by Grossman and Ion (8) and studied here has also been analyzed extensively by Grossman and collaborators (8, 9, 25) and by others (26). Grossman (9) gives a number of results about the time evolution of the collaboration network. He notes that the rate of publication has increased slightly over the last 50 years or so, but there has been a much more striking increase in the level of collaboration. From the start of the period covered by the *Mathematical Reviews* data in 1940 until the end of the 1950s, less than one-half of all mathematicians had ever coauthored a paper with another writer; nowadays, virtually all of them have. Presumably this trend reflects some combination of changes in the social organization of the mathematics community, better communications, and possibly changes in the types of problems studied and approaches used, making modern mathematics more amenable to collaborative investigation.

The time evolution of coauthorship networks has also been investigated by the present author (14) and others (10) in the context of tests of the “preferential attachment” hypothesis. Price (30) and later Barabási and Albert (31) have suggested that networks grow by the addition of connections in such a way that the probability of an individual gaining a new connection is proportional to the number they already have. Because coauthorship networks are well time-resolved, as discussed in the preceding section, one can test this hypothesis by measuring the probability that a newly published paper contributes new connections to an individual, as a function of the number of connections that individual already has. Refs. 10 and 14 tackle this measurement in slightly different ways, but both conclude that preferential attachment of an approximately linear variety is indeed taking place in collaboration networks.

A number of authors have also looked at “betweenness centrality” in coauthorship networks (13, 32–34). The betweenness centrality of a node A in a network is defined to be the number of shortest paths between other pairs of nodes that pass through A (35). It is regarded as a measure of the influence that individuals have over information flow between others. Individuals who act as brokers for information flow between their colleagues will have high betweenness scores. In ref. 13, it was shown that betweenness scores vary widely from one individual to another in coauthorship networks, with a few having much higher scores than the majority. Later work by Goh *et al.* (33) showed that in fact the distribution of betweenness scores approximately follows a power law. This appears to indicate that collaboration networks contain a small number of influential individuals and many peripheral actors, a conclusion bolstered by the findings of Holme *et al.* (32), who showed that collaboration networks are highly susceptible to the removal of the

individuals with highest betweenness scores. One need only remove a few such individuals from the network, it turns out, to break the connection between others and fracture the network into disconnected parts. In a recent paper, Goh *et al.* (34) have extended their investigation of betweenness to the correlation between the betweenness scores of collaborators. They find that there is very little such correlation, implying that influential scientists are not collaborating preferentially with other influential scientists to any significant extent; the probability of one’s collaborator having a high betweenness appears not to be significantly greater if one has a high betweenness than if one does not.

Conclusion

In this paper, we have discussed the structure of three networks of scientific collaborations, as deduced from the pattern of coauthorships of papers. The networks cover biomedical research, physics, and mathematics, respectively, and the results reveal both similarities and differences among the different fields. All fields appear to have a broad distribution of the number of coauthors that an individual has, with most individuals having only a few coauthors, whereas a few have many, hundreds or even thousands in some cases. Biological scientists tend to have significantly more coauthors than mathematicians or physicists, a result that reflects the labor-intensive, predominantly experimental direction of current biology. Other differences are less easily explained. In biology, for instance, it is far less likely than in mathematics that two of one’s coauthors will also be coauthors of one another, a result that has yet to receive a clear explanation.

Coauthorship networks provide a copious and meticulously documented record of the social and professional networks of scientists. The results reported here represent only a small portion of what could be done with these data. Possible future directions for study might include tests for community structure or “invisible colleges” within the networks (36, 37) or further investigations of changes in collaboration patterns over time (9, 14), as well as other measurements not yet thought of. Coauthorship data represent a superb resource for the pursuit of questions such as these, and I look forward to future developments with interest.

I thank particularly Paul Ginsparg for help in obtaining the data used for this study. The data were generously made available by Oleg Khovayko, David Lipman, and Grigoriy Starchenko (Medline); Paul Ginsparg and Geoffrey West (Physics E-print Archive); and Jerry Grossman (*Mathematical Reviews*). I also thank Steve Strogatz for suggesting the “funneling effect” calculation of *Statistical Properties of Coauthorship Networks* and László Barabási, Paul Ginsparg, Jon Kleinberg, Sidney Redner, Steven Strogatz, and Duncan Watts for useful comments and suggestions. This work was funded in part by the James S. McDonnell Foundation, by the Intel Corporation, and by the U.S. National Science Foundation under Grants DMS-0109086 and DMS-0234188.

1. Price, D. J. (1965) *Science* **149**, 510–515.
2. Egghe, L. & Rousseau, R. (1990) *Introduction to Informetrics* (Elsevier, Amsterdam).
3. Kretschmer, H. (1994) *Scientometrics* **30**, 363–369.
4. Persson, O. & Beckmann, M. (1995) *Scientometrics* **33**, 351–366.
5. Melin, G. & Persson, O. (1996) *Scientometrics* **36**, 363–377.
6. Ding, Y., Foo, S. & Chowdhury, G. (1999) *Int. Inform. Lib. Rev.* **30**, 367–376.
7. Bordens, M. & Gómez, I. (2000) in *The Web of Knowledge: A Festschrift in Honor of Eugene Garfield*, eds. Atkins, H. B. & Cronin, B. (Information Today, Medford, NJ).
8. Grossman, J. W. & Ion, P. D. F. (1995) *Congressus Numerantium* **108**, 129–131.
9. Grossman, J. W. (2002) *Congressus Numerantium* **158**, 202–212.
10. Barabási, A.-L., Jeong, H., Ravasz, E., Néda, Z., Schuberts, A. & Vicsek, T. (2002) *Physica A* **311**, 590–614.
11. Newman, M. E. J. (2001) *Proc. Natl. Acad. Sci. USA* **98**, 404–409.
12. Newman, M. E. J. (2001) *Phys. Rev. E Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.* **64**, 016131.
13. Newman, M. E. J. (2001) *Phys. Rev. E Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.* **64**, 016132.
14. Newman, M. E. J. (2001) *Phys. Rev. E Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.* **64**, 025102.
15. Newman, M. E. J. (2002) *Phys. Rev. Lett.* **89**, 208701.
16. Lotka, A. J. (1926) *J. Wash. Acad. Sci.* **16**, 317–323.
17. Shockley, W. (1957) *Proc. IRE* **45**, 279–290.
18. Voos, H. (1974) *J. Am. Soc. Inf. Sci.* (July–August 1974), 270–272.
19. Pao, M. L. (1986) *J. Am. Soc. Inf. Sci.* (January 1986), 26–33.
20. Fennel, T., Levene, M. & Loizou, G. (2002) cond-mat/0209463 (preprint).
21. Pool, I. de S. & Kochen, M. (1978) *Soc. Networks* **1**, 1–48.
22. Milgram, S. (1967) *Psychol. Today* **2**, 60–67.
23. Travers, J. & Milgram, S. (1969) *Sociometry* **32**, 425–443.
24. Cormen, T. H., Leiserson, C. E., Rivest, R. L. & Stein, C. (2001) *Introduction to Algorithms* (MIT Press, Cambridge, MA), 2nd Ed.
25. de Castro, R. & Grossman, J. W. (1999) *Math. Intell.* **21**, 51–63.
26. Batagelj, V. & Mrvar, A. (2000) *Soc. Networks* **22**, 173–186.
27. Kautz, H., Selman, B. & Shah, M. (1997) *Comm. ACM* **40**, 63–65.
28. Watts, D. J. & Strogatz, S. H. (1998) *Nature* **393**, 440–442.
29. Newman, M. E. J., Strogatz, S. H. & Watts, D. J. (2001) *Phys. Rev. E Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.* **64**, 026118.

30. Price, D. J. (1976) *J. Am. Soc. Inform. Sci.* **27**, 292–306.
31. Barabási, A.-L. & Albert, R. (1999) *Science* **286**, 509–512.
32. Holme, P., Kim, B. J., Yoon, C. N. & Han, S. K. (2002) *Phys. Rev. E Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.* **65**, 056109.
33. Goh, K.-I., Oh, E., Jeong, H., Kahng, B. & Kim, D. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 12583–12588.
34. Goh, K.-I., Oh, E., Kahng, B. & Kim, D. (2003) *Phys. Rev. E Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.* **67**, 017101.
35. Freeman, L. C. (1977) *Sociometry* **40**, 35–41.
36. Crane, D. (1972) *Invisible Colleges: Diffusion of Knowledge in Scientific Communities* (Univ. of Chicago Press, Chicago).
37. Girvan, M. & Newman, M. E. J. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 7821–7826.