

LETTERS

Uncovering the overlapping community structure of complex networks in nature and society

Gergely Palla^{1,2}, Imre Derényi², Illés Farkas¹ & Tamás Vicsek^{1,2}

Many complex systems in nature and society can be described in terms of networks capturing the intricate web of connections among the units they are made of^{1–4}. A key question is how to interpret the global organization of such networks as the co-existence of their structural subunits (communities) associated with more highly interconnected parts. Identifying these *a priori* unknown building blocks (such as functionally related proteins^{5,6}, industrial sectors⁷ and groups of people^{8,9}) is crucial to the understanding of the structural and functional properties of networks. The existing deterministic methods used for large networks find separated communities, whereas most of the actual networks are made of highly overlapping cohesive groups of nodes. Here we introduce an approach to analysing the main statistical features of the interwoven sets of overlapping communities that makes a step towards uncovering the modular structure of complex systems. After defining a set of new characteristic quantities for the statistics of communities, we apply an efficient technique for exploring overlapping communities on a large scale. We find that overlaps are significant, and the distributions we introduce reveal universal features of networks. Our studies of collaboration, word-association and protein interaction graphs show that the web of communities has non-trivial correlations and specific scaling properties.

Most real networks typically contain parts in which the nodes (units) are more highly connected to each other than to the rest of the network. The sets of such nodes are usually called clusters, communities, cohesive groups or modules^{8,10,11–13}; they have no widely accepted, unique definition. In spite of this ambiguity, the presence of communities in networks is a signature of the hierarchical nature of complex systems^{5,14}. The existing methods for finding communities in large networks are useful if the community structure is such that it can be interpreted in terms of separated sets of communities (see Fig. 1b and refs 10, 15, 16–18). However, most real networks are characterized by well-defined statistics of overlapping and nested communities. This can be illustrated by the numerous communities that each of us belongs to, including those related to our scientific activities or personal life (school, hobby, family) and so on, as shown in Fig. 1a. Furthermore, members of our communities have their own communities, resulting in an extremely complicated web of the communities themselves. This has long been understood by sociologists¹⁹ but has never been studied systematically for large networks. Another, biological, example is that a large fraction of proteins belong to several protein complexes simultaneously²⁰.

In general, each node i of a network can be characterized by a membership number m_i , which is the number of communities that the node belongs to. In turn, any two communities α and β can share $s_{\alpha,\beta}^{ov}$ nodes, which we define as the overlap size between these communities. Naturally, the communities also constitute a network,

with the overlaps being their links. The number of such links of community α can be called its community degree, d_{α}^{com} . Finally, the size s_{α}^{com} of any community α can most naturally be defined as the number of its nodes. To characterize the community structure of a large network we introduce the distributions of these four basic quantities. In particular we focus on their cumulative distribution

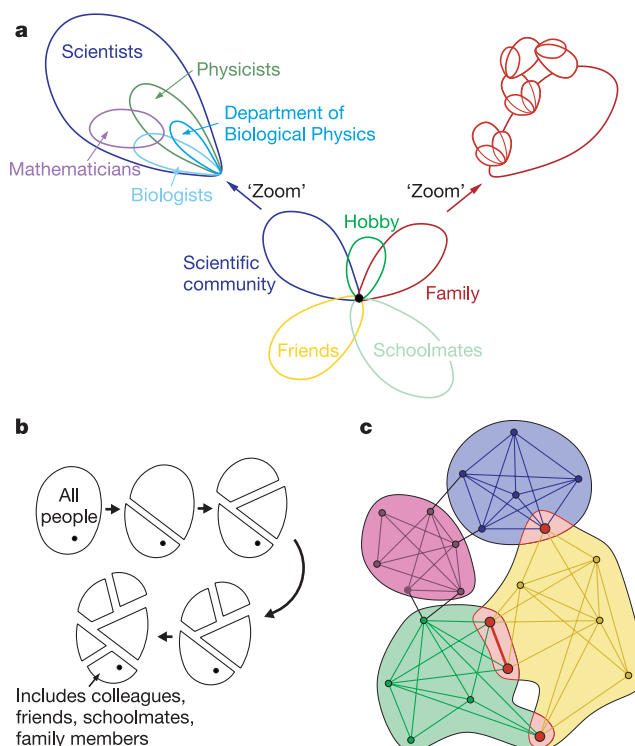


Figure 1 | Illustration of the concept of overlapping communities. **a**, The black dot in the middle represents either of the authors of this paper, with several of his communities around. Zooming in on the scientific community demonstrates the nested and overlapping structure of the communities, and depicting the cascades of communities starting from some members exemplifies the interwoven structure of the network of communities.

b, Divisive and agglomerative methods grossly fail to identify the communities when overlaps are significant. **c**, An example of overlapping k -clique communities at $k = 4$. The yellow community overlaps the blue one in a single node, whereas it shares two nodes and a link with the green one. These overlapping regions are emphasized in red. Notice that any k -clique (complete subgraph of size k) can be reached only from the k -cliques of the same community through a series of adjacent k -cliques. Two k -cliques are adjacent if they share $k - 1$ nodes.

¹Biological Physics Research Group of the Hungarian Academy of Sciences, Pázmány P. stny. 1A, H-1117 Budapest, Hungary. ²Department of Biological Physics, Eötvös University, Pázmány P. stny. 1A, H-1117 Budapest, Hungary.

functions denoted by $P(m)$, $P(s^{ov})$, $P(d^{com})$ and $P(s^{com})$. For the overlap size, for example, $P(s^{ov})$ means the proportion of those overlaps that are larger than s^{ov} . Further relevant statistical features will be introduced later.

The basic observation on which our community definition relies is that a typical community consists of several complete (fully connected) subgraphs that tend to share many of their nodes. Thus, we define a community, or more precisely a k -clique community, as a union of all k -cliques (complete subgraphs of size k) that can be reached from each other through a series of adjacent k -cliques (where adjacency means sharing $k - 1$ nodes)^{21–23}. This definition seeks to represent the fact that it is an essential feature of a community that its members can be reached through well-connected subsets of nodes. There are other parts of the whole network that are not reachable from a particular k -clique, but they potentially contain further k -clique communities. In turn, a single node can belong to several communities. All these can be explored systematically and can result in many overlapping communities (illustrated in Fig. 1c). In most cases, relaxing this definition (for example, by allowing incomplete k -cliques) is practically equivalent to decreasing k . For finding meaningful communities, the way in which they are identified is expected to satisfy several basic requirements: it cannot be too restrictive, it should be based on the density of links, it is required to be local, it should not yield any cut-node or cut-link (whose removal would disjoin the community) and, of course, it should allow overlaps. We employ the community definition specified above, because none of the others in the literature satisfy all these requirements simultaneously^{21,24}.

Although the numerical determination of the full set of k -clique communities is a polynomial problem, we use an algorithm (which can be downloaded from <http://angel.elte.hu/clustering/>) that is exponential, because it is significantly more efficient for the graphs corresponding to real data. This method is based on first locating all cliques (maximal complete subgraphs) of the network and then identifying the communities by carrying out a standard component analysis of the clique–clique overlap matrix²¹. More details about the method and its speed are given in Supplementary Information.

We use our method for binary networks (that is, with undirected and unweighted links). An arbitrary network can always be transformed into a binary one by ignoring any directionality in the links and keeping only those that are stronger than a threshold weight w^* . Changing the threshold is like changing the resolution (as in a microscope) with which the community structure is investigated: by increasing w^* the communities start to shrink and fall apart. A similar effect can be observed by changing the value of k as well: increasing k makes the communities smaller and more disintegrated but also at the same time more cohesive.

When we are interested in the community structure around a particular node, it is advisable to scan through some ranges of k and w^* and monitor how its communities change. As an illustration, in Fig. 2 we show diagrams of the communities of three selected nodes of three large networks: the social network of scientific collaborators²⁵ (Fig. 2a), the network of word associations²⁶ related to cognitive sciences (Fig. 2b) and the molecular-biological network of protein–protein interactions²⁷ (Fig. 2c). These pictures can serve as tests or validations of the efficiency of our algorithm. In particular,

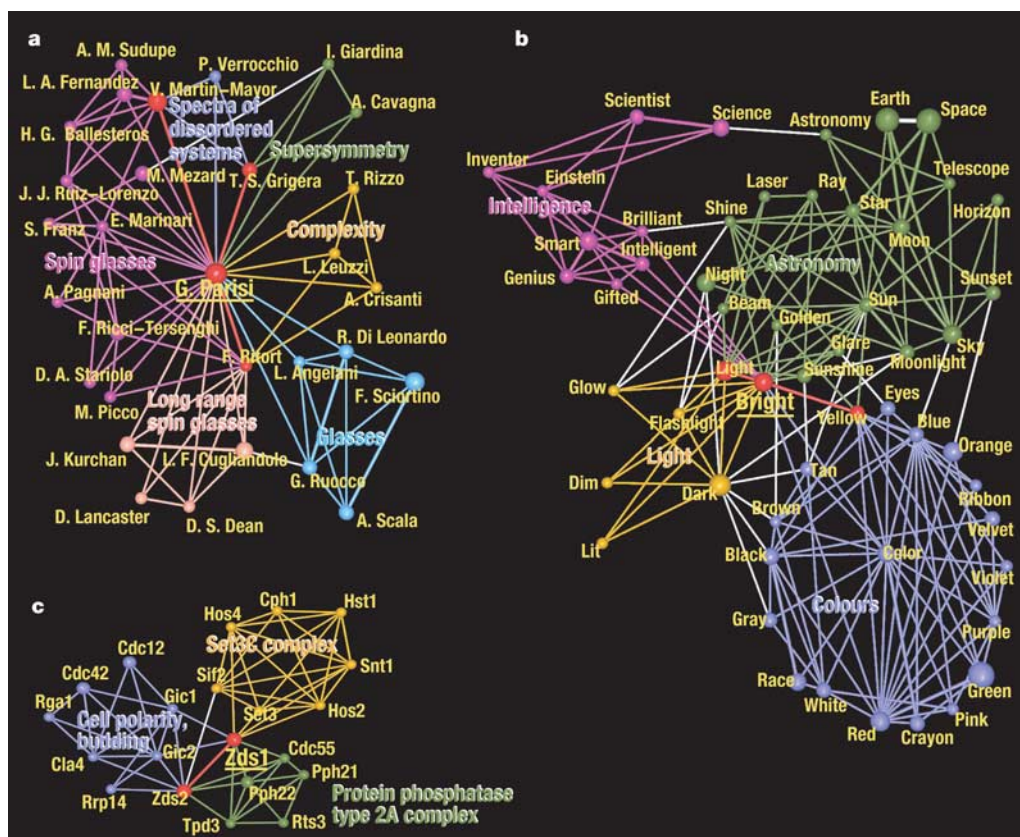


Figure 2 | The community structure around a particular node in three different networks. The communities are colour coded, the overlapping nodes and links between them are emphasized in red, and the volume of the balls and the width of the links are proportional to the total number of communities they belong to. For each network the value of k has been set to 4. **a**, The communities of G. Parisi in the co-authorship network of the Los Alamos Condensed Matter archive (for threshold weight $w^* = 0.75$) can

be associated with his fields of interest. **b**, The communities of the word 'bright' in the South Florida Free Association norms list (for $w^* = 0.025$) represent the different meanings of this word. **c**, The communities of the protein Zds1 in the DIP core list of the protein–protein interactions of *S. cerevisiae* can be associated with either protein complexes or certain functions.

the communities of G. Parisi (whose contributions in different fields of physics are well known) shown in Fig. 2a are associated with his fields of interest, as can be deduced from the titles of the papers involved. The four-clique communities of the word ‘bright’ (Fig. 2b) correspond to the various meanings of this word. An important biological application is finding the communities of proteins, based on their interactions. Indeed, most proteins in the communities shown in Figs 2c and 3 can be associated with either protein complexes or certain functions, as can be looked up by using the GO-TermFinder package²⁸ and the online tools of the *Saccharomyces* Genome Database (SGD)²⁹. For some proteins no function is yet available. Thus, the fact that they show up in our approach as members of communities can be interpreted as a prediction of their functions. One such example can be seen in the enlarged

portion of Fig. 3. For the protein Ycr072c, which is required for the viability of the cell and appears in the dark green community on the right, SGD provides no biological process (function). By far the most significant GO term for the biological process of this community is ‘ribosome biogenesis/assembly’. We can therefore infer that Ycr072c is likely to be involved in this process. In addition, new cellular processes can be predicted if as yet unknown communities are found with our method.

These examples (and further examples included in Supplementary Information) show the advantages of our approach over the existing divisive and agglomerative methods recently used for large real networks. Divisive methods cut the network into smaller and smaller pieces, and each node is forced to remain in only one community and be separated from its other communities, most of which then necessarily fall apart and disappear. This happens, for example, with the word ‘bright’ when we apply the method described in ref. 16: it tends to stay together mostly with the words of the community related to ‘light’, while most of its other communities (for example, those related to ‘colours’; see Fig. 2b) completely disintegrate (‘green’ becomes associated with the vegetables, ‘orange’ with the fruits, and so on). Agglomerative methods do the same, but in the reverse direction. For example, when we applied the agglomerative method of ref. 18, at some point ‘bright’, as a single word, joined a ‘community’ of 890 other words. In addition, such methods inevitably lead to a tree-like hierarchical rendering of the communities, whereas our approach allows the construction of an unconstrained network of communities.

The networks chosen above have been constructed in the following ways. In the co-authorship network of the Los Alamos e-print archives²⁵ each article contributes a value $1/(n-1)$ to the weight of the link between every pair of its n authors. In the South Florida Free Association norms list²⁶ the weight of a directed link from one word to another indicates the frequency with which the people in the survey associated the end point of the link with its starting point. For our purposes these directed links have been replaced by undirected ones with a weight equal to the sum of the weights of the corresponding two oppositely directed links. In the Database of Interacting Proteins (DIP) core list of the protein–protein interactions of *Saccharomyces cerevisiae*²⁷ each interaction represents an unweighted link between the interacting proteins. These networks are very large, consisting of 30,739, 10,617 and 2,609 nodes and 136,065, 63,788 and 6,355 links, respectively.

Although different values of k and w^* might be optimal for the local community structure around different nodes, we should set some global criterion to fix their values if we wish to analyse the statistical properties of the community structure of the entire network. The criterion we use is based on finding a community structure that is as highly structured as possible. In the related percolation phenomena²³ a giant component appears when the number of links is increased above some critical point. Therefore, to approach this critical point from below, for each selected value of k (typically between 3 and 6) we lower the threshold w^* until the largest community becomes twice as big as the second largest one. In this way we ensure that we find as many communities as possible, without the negative effect of having a giant community that would smear out the details of the community structure by merging many smaller communities. We denote by f^* the fraction of links stronger than w^* ,

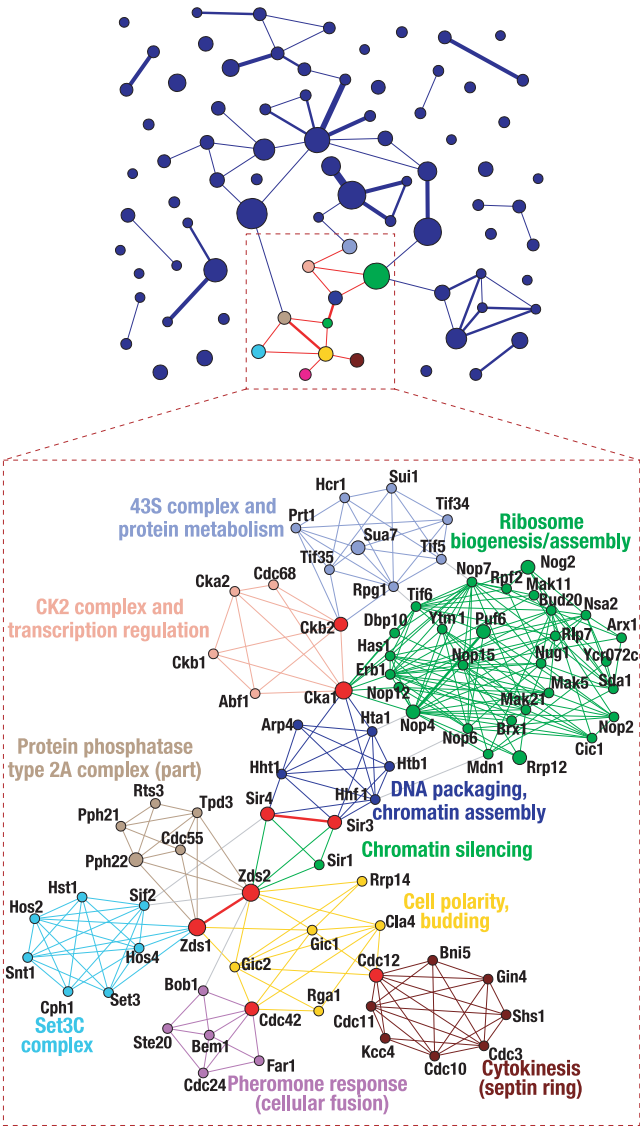


Figure 3 | Network of the 82 communities in the DIP core list of the protein–protein interactions of *S. cerevisiae* for $k = 4$. The areas of the circles and the widths of the links are proportional to the size of the corresponding communities (s_{α}^{com}) and to the size of the overlaps ($s_{\alpha\beta}^{ov}$), respectively. The coloured communities (top) are cut out and magnified to reveal their internal structure (bottom): the nodes and links of the original network have the same colour as their communities, those that are shared by more than one community are emphasized in red, and the grey links are not part of these communities. The areas of the circles and the widths of the links are proportional to the total number of communities they belong to.

Table 1 | Statistical properties of the network of communities

Network	N^{com}	$\langle d^{com} \rangle$	$\langle C^{com} \rangle$	$\langle r \rangle$
Co-authorship	2,450	12.10	0.44	0.58
Word association	670	11.33	0.56	0.72
Protein interaction	82	1.54	0.17	0.26

N^{com} is the number of communities, $\langle d^{com} \rangle$ is the average community degree, $\langle C^{com} \rangle$ is the average clustering coefficient of the network of communities, and $\langle r \rangle$ is the average fraction of shared nodes in the communities.

and use only those values of k for which f^* is not too small (not smaller than 0.5). This has led us to $k = 6$ and $k = 5$ with $f^* = 0.93$ and 0.75, respectively, for the collaboration network, and $k = 4$ with $f^* = 0.67$ for the word-association network. For the former network both sets of parameters result in very similar communities (see Supplementary Information). Because for unweighted networks no threshold weight can be set, for these we simply select the smallest value of k for which no giant community appears. For the protein interaction network this gives $k = 4$, resulting in 82 communities. Because of this relatively low number, we can depict the entire network of protein communities as in Fig. 3.

The four distributions characterizing the global community structure of these networks are shown in Fig. 4. Although the scaling of the size of non-overlapping communities has already been shown

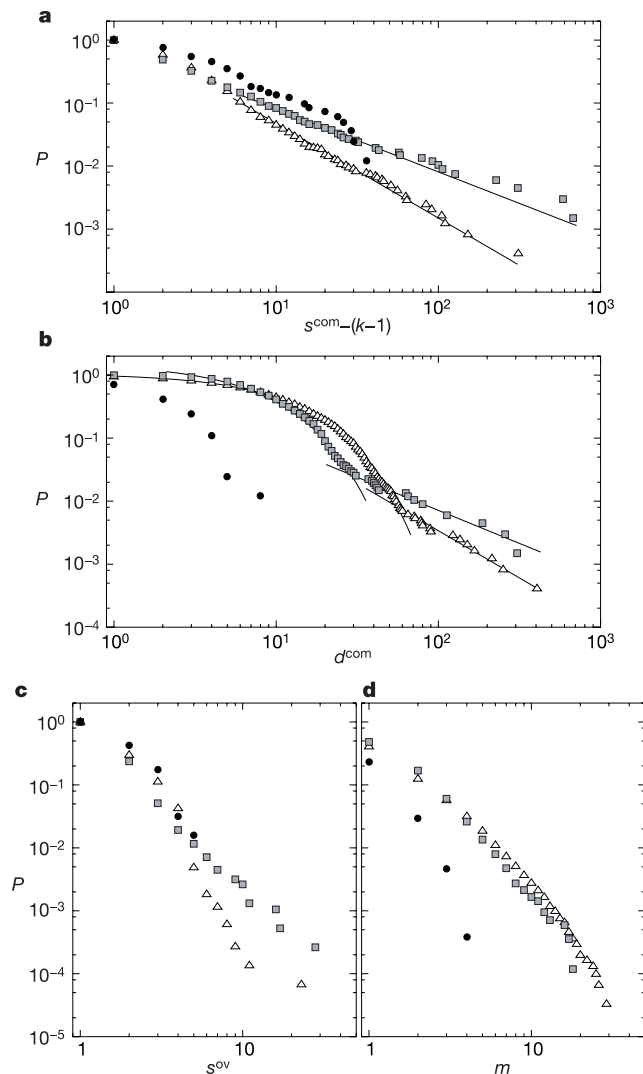


Figure 4 | Statistics of the k -clique communities for three large networks. The networks are the co-authorship network of the Los Alamos Condensed Matter archive (triangles, $k = 6$, $f^* = 0.93$), the word-association network of the South Florida Free Association norms (squares, $k = 4$, $f^* = 0.67$), and the protein interaction network of the yeast *S. cerevisiae* from the DIP database (circles, $k = 4$). **a**, The cumulative distribution function of the community size follows a power law with exponents between -1 (upper line) and -1.6 (lower line). **b**, The cumulative distribution of the community degree starts exponentially and then crosses over to a power law (with the same exponent as for the community size distribution). **c**, The cumulative distribution of the overlap size. **d**, The cumulative distribution of the membership number.

for social networks^{17,18}, it is striking to observe how this aspect of large real networks is preserved even when a more complete picture (allowing overlaps) is investigated. In Fig. 4a the power-law dependence $P(s^{\text{com}}) \propto (s^{\text{com}})^{-\tau}$ with an exponent ranging between $\tau = 1$ and $\tau = 1.6$ is well pronounced and is valid over nearly the entire range of community sizes.

It is well known²⁻⁴ that the nodes of large real networks have a power-law degree distribution. Will the same kind of distribution hold when we move to the next level of organization and consider the degrees of the communities? We find that it is not so. The community degrees (Fig. 4b) have a unique distribution, consisting of two distinct parts: an exponential decay $P(d^{\text{com}}) \propto \exp(-d^{\text{com}}/d_0^{\text{com}})$ with a characteristic community degree d_0^{com} (which is of the order of $\langle d^{\text{com}} \rangle$ shown in Table 1), followed by a power-law tail proportional to $(d^{\text{com}})^{-\tau}$. This new kind of behaviour is consistent with the community size distribution if we assume that, on average, each node of a community has a contribution δ to the community degree. The tail of the community degree distribution is therefore simply proportional to that of the community size distribution. At the first part of $P(d^{\text{com}})$, in contrast, a characteristic scale $d_0^{\text{com}} \approx k\delta$ appears, because most of the communities have a size of the order of k (see Fig. 4a) and their distribution around d_0^{com} dominates this part of the curve. Thus, the degree to which $P(d^{\text{com}})$ deviates from a simple scaling depends on k or, in other words, on the prescribed minimum cohesiveness of the communities.

The extent to which different communities overlap is also a relevant property of a network. Although the range of overlap sizes is limited, the behaviour of the cumulative overlap size distribution $P(s^{\text{ov}})$, shown in Fig. 4c, is close to a power law for each network, with a rather large exponent. We can conclude that there is no characteristic overlap size in the networks. Finally, in Fig. 4d we display the cumulative distribution of the membership number $P(m)$. These plots demonstrate that a node can belong to several communities. In the collaboration and word-association networks there seems to be no characteristic value for the membership number: the data are close to a power-law dependence, with a large exponent. However, in the protein interaction network the largest membership number is only 4, which is consistent with the also rather short distribution of its community degree. To show that the communities we find are not due to an artefact of our method, we have also determined the above distributions for 'randomized' graphs with parameters (size, degree sequence, k and f^*) the same as in our three examples but with links stochastically redistributed between the nodes. We have found that the distributions are indeed extremely truncated, signifying a complete lack of the rich community structure determined for the original data.

In Table 1 we have collected a few statistical properties of the network of communities. It should be pointed out that the average clustering coefficients $\langle C^{\text{com}} \rangle$ are relatively high, indicating that two communities overlapping with a given community are likely to overlap with each other as well, mostly because they all share the same overlapping region. The high fraction of shared nodes is yet another indication of the importance of overlaps between the communities.

The specific scaling of the community degree distribution is a hitherto undescribed signature of the hierarchical nature of the systems we study. We find that if we consider the network of communities instead of the nodes themselves, we still observe a degree distribution with a fat tail, but a characteristic scale appears, below which the distribution is exponential. This is consistent with our understanding of a complex system having different levels of organization with units specific to each level. In the present case the principle of organization (scaling) is preserved (with some specific modifications) when going to the next level, in good agreement with the recent finding of the self-similarity of many complex networks³⁰.

With recent technological advances, huge sets of data are accumulating at a tremendous pace in various fields of human activity

(including telecommunications, the Internet and stock markets) and in many areas of life and social sciences (such as biomolecular assays, genetic maps and groups of World Wide Web users). Understanding both the universal and specific features of the networks associated with these data has become a significant task. The knowledge of the community structure enables the prediction of some essential features of the systems under investigation. For example, because with our approach it is possible to 'zoom' in on a single unit in a network and uncover its communities (and the communities connected to these, and so on), we provide a tool with which to interpret the local organization of large networks and can predict how the modular structure of the network changes if a unit is removed (for example, in a gene knockout experiment). A unique feature of our method is that we can simultaneously look at the network at a higher level of organization and locate the communities that have a key role within the web of communities. Among the many possible applications is a more sophisticated approach to the spreading of infections (for example, real or computer viruses) or information in highly modular complex systems.

Received 17 January; accepted 7 April 2005.

- Watts, D. J. & Strogatz, S. H. Collective dynamics of 'small-world' networks. *Nature* **393**, 440–442 (1998).
- Barabási, A.-L. & Albert, R. Emergence of scaling in random networks. *Science* **286**, 509–512 (1999).
- Albert, R. & Barabási, A.-L. Statistical mechanics of complex networks. *Rev. Mod. Phys.* **74**, 47–97 (2002).
- Mendes, J. F. F. & Dorogovtsev, S. N. *Evolution of Networks: From Biological Nets to the Internet and WWW* (Oxford Univ. Press, Oxford, 2003).
- Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. & Barabási, A.-L. Hierarchical organization of modularity in metabolic networks. *Science* **297**, 1551–1555 (2002).
- Spirin, V. & Mirny, L. A. Protein complexes and functional modules in molecular networks. *Proc. Natl Acad. Sci. USA* **100**, 12123–12128 (2003).
- Onnela, J.-P., Chakraborti, A., Kaski, K., Kertész, J. & Kanto, A. Dynamics of market correlations: Taxonomy and portfolio analysis. *Phys. Rev. E* **68**, 056110 (2003).
- Scott, J. *Social Network Analysis: A Handbook* 2nd edn (Sage, London, 2000).
- Watts, D. J., Dodds, P. S. & Newman, M. E. J. Identity and search in social networks. *Science* **296**, 1302–1305 (2002).
- Shiffrin, R. M. & Börner, K. Mapping knowledge domains. *Proc. Natl Acad. Sci. USA* **101**, 5183–5185 (2004).
- Everitt, B. S. *Cluster Analysis* 3rd edn (Edward Arnold, London, 1993).
- Knudsen, S. *A Guide to Analysis of DNA Microarray Data* 2nd edn (Wiley-Liss, New York, 2004).
- Newman, M. E. J. Detecting community structure in networks. *Eur. Phys. J. B* **38**, 321–330 (2004).
- Vicsek, T. The bigger picture. *Nature* **418**, 131 (2002).
- Blatt, M., Wiseman, S. & Domany, E. Super-paramagnetic clustering of data. *Phys. Rev. Lett.* **76**, 3251–3254 (1996).
- Girvan, M. & Newman, M. E. J. Community structure in social and biological networks. *Proc. Natl Acad. Sci. USA* **99**, 7821–7826 (2002).
- Radicchi, F., Castellano, C., Cecconi, F., Loreto, V. & Parisi, D. Defining and identifying communities in networks. *Proc. Natl Acad. Sci. USA* **101**, 2658–2663 (2004).
- Newman, M. E. J. Fast algorithm for detecting community structure in networks. *Phys. Rev. E* **69**, 066133 (2004).
- Faust, K. in *Models and Methods in Social Network Analysis* (eds Carrington, P., Scott, J. & Wasserman, S.) 117–147 (Cambridge Univ. Press, New York, 2005).
- Gavin, A. C. et al. Functional organization of the yeast proteome by systematic analysis of protein complexes. *Nature* **415**, 141–147 (2002).
- Everett, M. G. & Borgatti, S. P. Analyzing clique overlap. *Connections* **21**, 49–61 (1998).
- Batagelj, V. & Zaversnik, M. Short cycles connectivity. *arXiv cs.DS/0308011* (<http://arxiv.org/abs/cs/0308011>) (2003).
- Derényi, I., Palla, G. & Vicsek, T. Clique percolation in random networks. *Phys. Rev. Lett.* **94**, 160202 (2005).
- Kosub, S. in *Network Analysis* (eds Brandes, U. & Erlebach, T.) 112–142 (Lecture Notes in Computer Science 3418, Springer, Berlin, 2005).
- Warner, S. E-prints and the Open Archives Initiative. *Library Hi Tech* **21**, 151–158 (2003).
- Nelson, D. L., McEvoy, C. L. & Schreiber, T. A. The University of South Florida word association, rhyme, and word fragment norms. (<http://www.usf.edu/FreeAssociation/>).
- Xenarios, I. et al. DIP: the Database of Interacting Proteins. *Nucleic Acids Res.* **28**, 289–291 (2000).
- Boyle, E. I. et al. GO:TermFinder—open source software for accessing Gene Ontology information and finding significantly enriched Gene Ontology terms associated with a list of genes. *Bioinformatics* **20**, 3710–3715 (2004).
- Cherry, J. M. et al. Genetic and physical maps of *Saccharomyces cerevisiae*. *Nature* **387**, 675–735 (1997).
- Song, C., Havlin, S. & Makse, H. A. Self-similarity of complex networks. *Nature* **433**, 392–395 (2005).

Supplementary Information accompanies the online version of the paper at www.nature.com/nature.

Acknowledgements We thank A.-L. Barabási and P. Pollner for discussions, and B. Kovács and G. Szabó for help with visualization and software support. This research was supported by the Hungarian Research Grant Foundation (OTKA).

Author Information Reprints and permissions information is available at npg.nature.com/reprintsandpermissions. The authors declare no competing financial interests. Correspondence and requests for materials should be addressed to T.V. (vicsek@angel.elte.hu).

Uncovering the overlapping community structure of complex networks in nature and society

Supplementary Information

1 The k -clique-community finding algorithm

Our community definition is based on the observation that a typical member in a community is linked to many other members, but not necessarily to all other nodes in the community. In other words, a community can be interpreted as a union of smaller complete (fully connected) subgraphs that share nodes. In the mathematical literature, such complete subgraphs are called k -cliques, where k refers to the number of nodes in the subgraph. Therefore, we define a k -clique-community as the union of all k -cliques that can be reached from each other through a series of *adjacent k -cliques*, where two k -cliques are said to be adjacent if they share $k - 1$ nodes. Using k -clique adjacency we can define a k -clique chain as the union of a sequence of adjacent k -cliques, and introduce the concept of k -clique connectedness: two k -cliques are k -clique-connected if they are parts of a k -clique chain. Our k -clique-communities are equivalent to the k -clique connected components of the network.

An illustration of these communities can be given by “ k -clique template rolling”. A k -clique template can be thought of as an object that is isomorphic to a complete graph of k nodes. Such a template can be placed onto any k -clique of the network, and rolled to an adjacent k -clique by relocating one of its nodes and keeping its other $k - 1$ nodes fixed. Thus, the k -clique-communities of a graph are all those subgraphs that can be fully explored by rolling a k -clique template in them but cannot be left by this template.

The k -clique-communities of a network at $k = 2$ are equivalent to the connected components, since a 2-clique is simply an edge and a 2-clique-community is the union of those edges that can be reached from each other through a series of shared nodes. Similarly, a 3-clique-community is given by the union of triangles that can be reached from one another through a series of shared edges. As we increase k , the k -clique-communities shrink, but on the other hand become more cohesive since their member nodes have to be part of at least one k -clique.

Our experience shows that in real networks complete subgraphs of size between 10 and 100 can easily occur. Such a large complete subgraph of size s contains $\binom{s}{k}$ different k -cliques, therefore, an algorithm that tries to locate the k -cliques individually and examine the adjacency between them would be extremely slow when analysing real networks. However, a complete subgraph of size s is obviously a k -clique connected subset for any $k \leq s$, since for any pair of included smaller k -cliques, a series of adjacent k -cliques linking them can be trivially found. Furthermore, two large complete subgraphs that share at least $k - 1$ nodes form one k -clique connected component as well. This implies that instead of searching for k -cliques, it is a far better strategy to locate the large complete subgraphs in the network first, and then look for the k -clique connected subsets of given k (the k -clique-communities) by studying the overlap between them.

1.1 The method

1.1.1 From cliques to k -clique-communities

To be more precise, our algorithm first extracts all complete subgraphs of the network that are not parts of larger complete subgraphs. (The details of this procedure are discussed in Sect. 1.1.2.) These maximal complete subgraphs are simply called *cliques*, and the difference between k -cliques and cliques is that k -cliques can be subsets of larger complete subgraphs. Once the cliques are located, the clique-clique overlap matrix is prepared [1]. In this symmetric matrix each row (and column) represents a clique and the matrix elements are equal to the number of common nodes between the corresponding two cliques, and the diagonal entries are equal to the size of the clique. (Note that the intersection of two cliques is always a complete subgraph.) The k -clique-communities for a given value of k are equivalent to such connected clique components in which the neighbouring cliques are linked to each other by at least $k - 1$ common nodes. These components can be found by erasing every off-diagonal entry smaller than $k - 1$

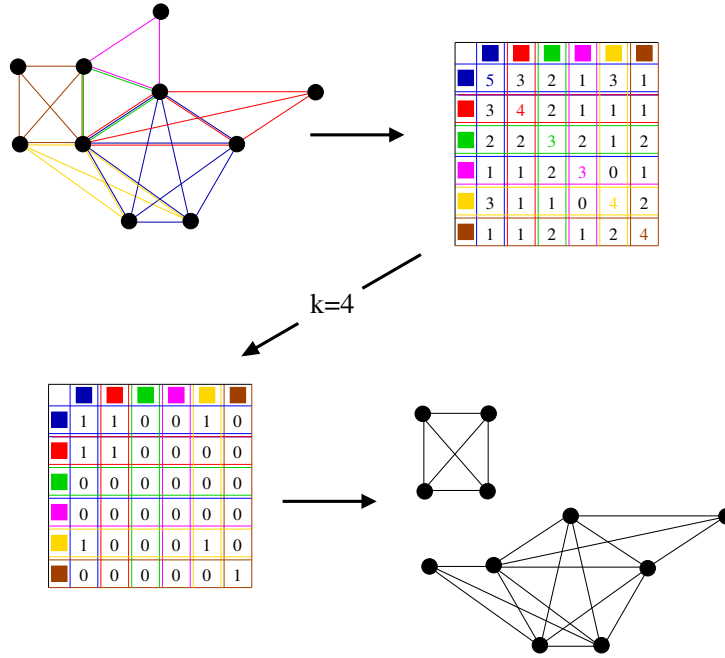


Figure 1: A simple illustration of the extraction of the k -clique-communities at $k = 4$ using the clique-clique overlap matrix. Top left picture shows the graph in which the different cliques are marked by different colours. The according clique-clique overlap matrix is shown in the top right corner. To obtain the k -clique-communities at $k = 4$, we delete the off-diagonal elements that are smaller than 3 and also the diagonal elements that are smaller than 4, resulting in the matrix shown in the bottom left of the figure. The connected components (the k -clique-communities) corresponding to this matrix are shown in the bottom right.

and every diagonal element smaller than k in the matrix, replacing the remaining elements by one, and then carrying out a component analysis of this matrix. The resulting separate components are equivalent to the different k -clique-communities. A simple illustration of the above is given in Fig. 1.

Another advantage of this method is that the clique-clique overlap matrix encodes all information necessary to obtain the communities for any value of k , therefore once the clique-clique overlap matrix is constructed, the k -clique-communities for all possible values of k can be obtained very quickly. In contrast to this, in a simple k -clique finding approach the search for the k -cliques would have to be restarted from the beginning for every single value of k .

1.1.2 Locating the cliques

As discussed in the previous section, in contrast to the k -cliques, cliques cannot be subsets of larger cliques, therefore they have to be located in a decreasing order of their size. The largest possible clique size in the studied graph is determined from the degree-sequence. Starting with this clique size, our algorithm repeatedly chooses a node, extracts every clique of this size containing that node, then deletes the node and its edges. (The deletion of the already examined nodes inhibits the finding of the same clique multiple times). When no nodes are left, the clique size is decreased by one and the clique finding procedure is restarted on the original graph. The already found cliques influence the further search since the yet unrevealed (smaller) cliques cannot be subsets of them.

The cliques of size s containing a given node v can be found by examining the interrelations of the neighbours of v . In our algorithm this is implemented in the following way: First, a set \mathcal{A} is constructed

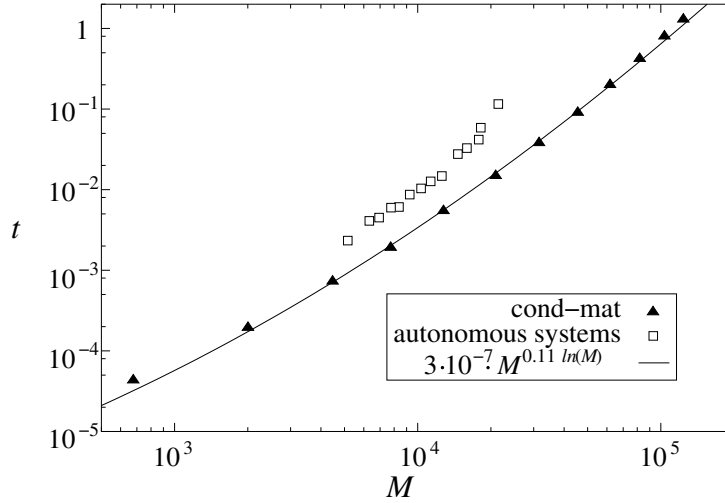


Figure 2: The time in hours on a PC needed to locate the communities as a function of the system size in the number of edges for the cond-mat archive (triangles) and for the graph of autonomous systems (squares). The former dataset is fitted with $3 \cdot 10^{-7} M^{0.11 \ln(M)}$ (solid curve).

that contains nodes all linked to each other. Initially \mathcal{A} consists of v only and our goal is to enlarge this set to the actual clique-size s . Another disjunct set \mathcal{B} is also determined as the set of nodes that are linked to each node in \mathcal{A} , but not necessarily to the nodes in \mathcal{B} . Initially set \mathcal{B} consists of the neighbours of v .

Set \mathcal{A} can be enlarged transferring nodes from \mathcal{B} . This is accomplished in a recursive way in order to check every possible combination of the nodes being transferred. (To avoid finding the same clique multiple times, the nodes have to be transferred from \mathcal{B} to \mathcal{A} in a decreasing/increasing order of their indices.) When a node w from \mathcal{B} is placed into \mathcal{A} , the nodes that are not neighbours of w are removed from \mathcal{B} . (This is done in order to preserve the property that the members of \mathcal{B} are all linked to each member of \mathcal{A}).

If \mathcal{B} runs out of nodes before \mathcal{A} reaches size s , or if the union of the sets \mathcal{A} and \mathcal{B} can be included in an already found (larger) clique, the recursion is stepped back to check other possibilities. Whenever the size of \mathcal{A} reaches s , a new clique is found. After recording the clique, the algorithm is stepped back again to check the remaining possible combinations of the neighbours indices.

1.2 Efficiency of the algorithm

The determination of the full set of cliques of a graph is widely believed to be non-polynomial problem. In spite of this, our algorithm proves to be very efficient when applied to the graphs of the investigated real systems. Our experience shows that the required CPU time depends on the structure of the input data very strongly, therefore in general no closed formula can be given even to estimate the system size dependence. As an illustration of the computational speed, however, we note that a complete analysis of a co-authorship network with 127000 links takes less than 2 hours on a PC.

In Fig. 2 we display the time it took to explore the community structure (using a PC) as a function of the system size in case of the co-authorship network of the Los Alamos Condensed Matter e-print archive [2, 3] at the optimal threshold for $k = 6$ and the network of autonomous systems [4]. (In both cases the graphs of different size correspond to the state of the system at different times). As it can be seen in the figure, the curves can be fitted with $t = AM^{B \ln(M)}$ where t denotes the time needed by our algorithm, M stands for the number of edges, and A and B are fitting parameters.

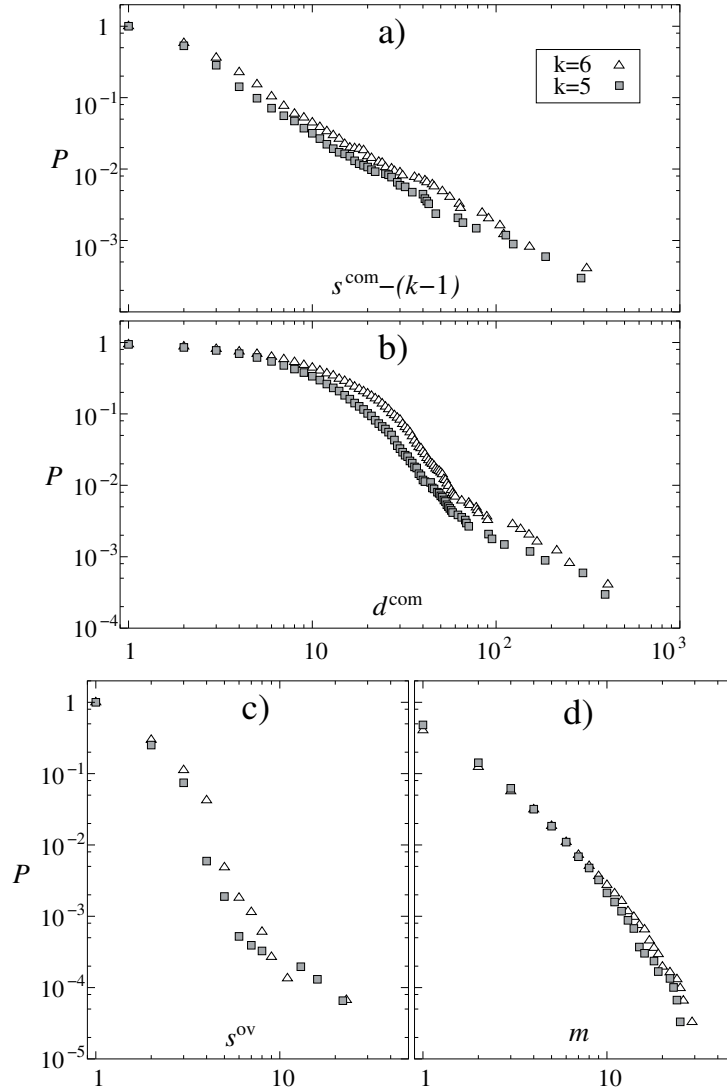


Figure 3: Statistics of the k -clique-communities for the Los Alamos Condensed Matter e-print archive at $k = 5$ (squares) and $k = 6$ (triangles). (a) the cumulative distribution function of the k -clique-community size (b) the cumulative distribution function of the k -clique-community degree (the degree distribution of the graph of communities), (c) the cumulative distribution function of the overlap size, and (d) the cumulative distribution function of the membership number of nodes.

2 Community statistics at different values of k

Our method can be directly applied to binary (undirected, unweighted) networks only. Therefore, when analysing an arbitrary system, the directionality of the links has to be ignored and if the connections are weighted, a threshold weight w^* can be introduced to prune weak links and keep those that are stronger than w^* . (If we want to keep all links, w^* is simply set to zero). If the threshold weight is increased, the number of edges is decreased and the communities shrink, however they consist of stronger links on average. Similarly, if k is increased at fixed threshold weight, the communities become smaller and more disintegrated, but at the same time also more cohesive (since every member in a community has to be part of a larger complete subgraph).

The criterion we used to fix the optimal k and w^* values is based on finding a community structure

as highly structured as possible. Usually a lower threshold weight is accompanied by a larger number of communities as more edges are left in the network. However, at a certain critical point a giant community appears which smears out the details of the community structure. Thus, for each selected value of k we adjusted the weight threshold to the point where the largest community becomes twice as big as the second largest one (just below the critical point). The restriction for the value of k we used was that at least half of the links should remain for the optimal threshold.

In case of the network representing the Los Alamos Condensed Matter e-print archive the criterions for the global k and w^* values could be matched at both $k = 5$ and $k = 6$. (In the former case the fraction f^* of the connections being kept during the application of our method was equal to $f^* = 0.75$, whereas in the latter case it turned out to be $f^* = 0.93$). In Fig. 3 we compare the relevant distributions characterising the community structure for the two values of k . In Fig. 3a the two scaling cumulative community size distributions are almost on top of each other. In case of the community degree (Fig. 3b) the scaling tails of the distribution functions are parallel similarly to the previous case. However the two distributions differ slightly at their exponential part, namely the characteristic community degree is a bit higher for $k = 6$ than for $k = 5$. There is a small difference between the two overlap size distributions as well at the middle part of the distributions (Fig. 3c). Finally, the two membership number distributions displayed in Fig. 3d match each other very well.

It can be seen from the distributions at $m = 1$ that the fraction of nodes belonging to at least one community is somewhere between 25% and 50%. The majority of the rest of the nodes fall out simply because their degree is less than $k - 1$. Nevertheless, after identifying the communities, most of these weakly connected nodes can be associated with the communities to which they are most strongly connected.

Besides this very good agreement between the relevant statistical distributions, the communities themselves show great similarities in the two cases: 44 % of the 6-clique-communities are present amongst the 5-clique-communities, and for 70 % of the 6-clique-communities one can find a corresponding 5-clique-community that differs in less than 10 % of the members. The good agreement between the results obtained for different values of k signals that the fundamental properties of the observed community structure are characteristic to the system itself and are largely independent of k .

3 Further examples

In this section we present a few more examples from the results of our community finding method. These concern both the global statistical properties of the communities determined for two additional data sets, (the Hungarian synonyms and the variables of the source code of the ftp program under Linux), as well as the local community structure around further vertices in the word association graph and in the network of the ftp program.

3.1 Community statistics

Similarly to Fig. 4 in the manuscript, the four major distributions characterising the global community structure of two further systems are plotted in Fig. 4. The triangles correspond to the network of the wu-ftp program under Linux [5] and the squares refer to the Hungarian synonym graph obtained from the OpenOffice word processor [6]. In the former network the nodes correspond to variables in the source code and are assumed to be connected if they appear together in an expression or function call, whereas in the second network two words are linked if they are synonyms of each other. The number of nodes N and links M are given by $N = 1886, 20139$ and $M = 6001, 100427$ for the network of the ftp program and the synonyms respectively. In both cases, our criterions for the global choice of the k -clique size can be matched only at $k = 5$.

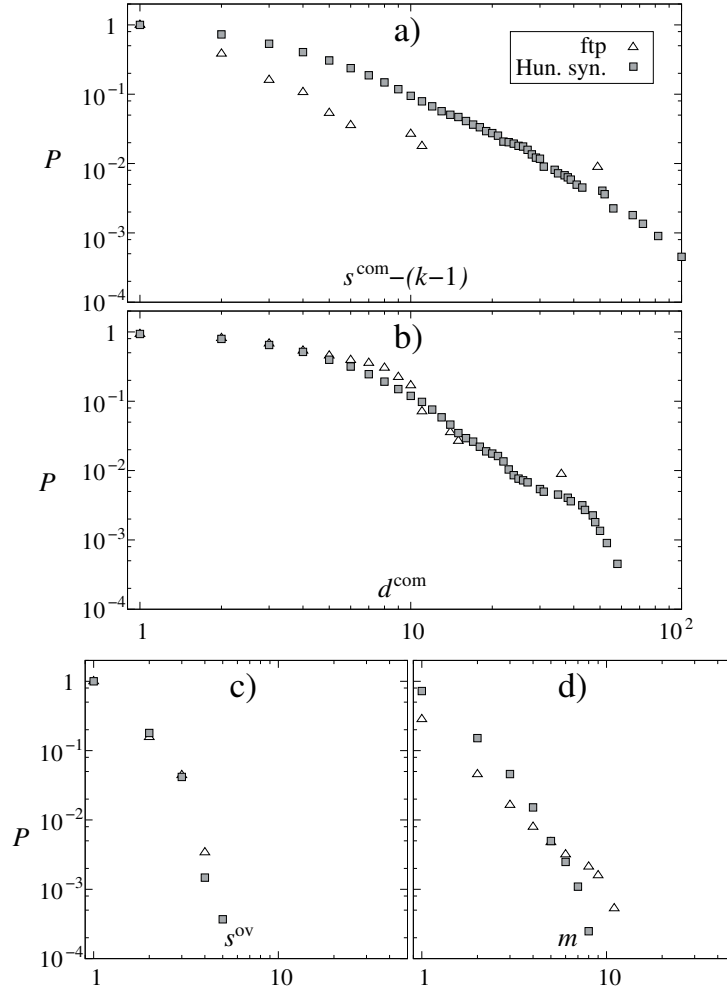


Figure 4: Statistics of the k -clique-communities for the wu-ftp program under Linux (triangles, $k=4$) and the graph of the Hungarian synonyms obtained from the OpenOffice word processor (squares, $k = 4$). (a) The cumulative distribution of the community size, (b) the cumulative distribution of the community degree, plot (c) is the cumulative distribution of the overlap size and (d) is that of the membership number.

Although our results for the two new data sets resemble those obtained for the data in the manuscript, there are also some deviations. In Fig. 4a the tails of the community size distributions are power-law like (however, not over such a wide range as, *i.e.*, in case of the co-authorship network). The lower part of the community degree distributions is exponential (Fig. 4b), but the extra power-law like tail present in case of the co-authorship network and the word association network is much less pronounced here. Due to the relatively small system size there is only one outstanding community degree in case of the ftp program, whereas the tail of the community degree distribution of the synonyms is somewhat like staircase. The community overlap distributions (Fig. 4c) are rather truncated, the maximal overlap size reaches just the k -clique size for the synonyms and is equal to $k - 1$ for the ftp program. In Fig. 4d, the membership number distributions decay somewhat faster than in case of the co-authorship network or the word association network.

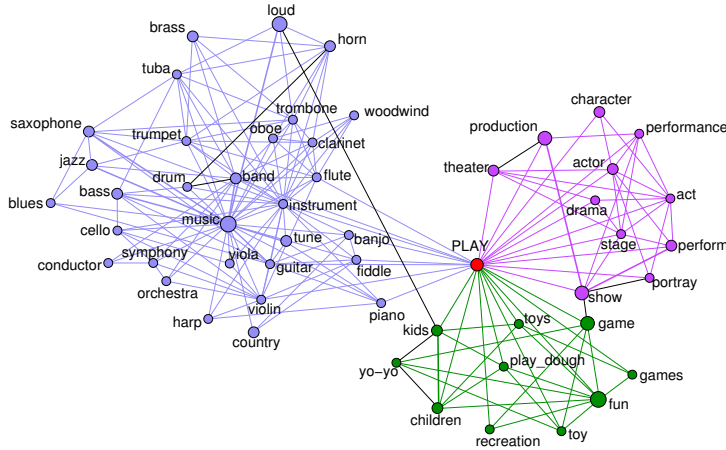


Figure 7: The k -clique communities of the word *play* in the South Florida Free Association norm list for $w^* = 0.025$ and $k = 4$. The blue community is associated with music, the purple one is related to theatre and the yellow community can be associated with children.

welfare respectively. In Fig. 6 the communities of the word *day* are shown. The green community can be associated with work days. *Thursday* has only two neighbours (*Wednesday* and *Wed*) even in the original (unpruned) network, therefore it is missing from this community, whereas *Saturday* and *Sunday* are in another community related to weekend containing *Friday*, *night*, *week* and *weekend* itself as well. The purple community of Fig. 6 consists of day times, the yellow community contains common adjectives of *day* related to weather, and the blue community can be associated with the calendar. Fig. 7 displays the three communities of the word *play*: the blue one is related to music, the purple one to theatre and the green one can be associated with children.

In Fig. 8 of we show a component from the community graph of the *wu-ftp* program at $k = 5$ in a fashion similar to Fig. 3 in the manuscript. The name of each node consists of two parts: the first one is specific to the variable represented by the node and the second part (separated by '@') is specific to the scope of the variable (typically a function). The names ending in '@glb' denote global variables. Since these variables have global scope, (and therefore are visible in the entire program), they may appear in several function calls and expressions throughout the entire source code. Thus, in the corresponding network the vertices representing these variables are candidates for community overlaps. Indeed, in Fig. 8, the majority of the communities are related to functions in the source code, and several community overlaps are provided by vertices representing global variables.

4 Random community statistics

The non-trivial aspects of the distributions presented in Fig. 4 of the manuscript naturally give rise to the question whether the community statistics of a random graph would significantly differ from those studied in the manuscript. In other words, what happens with the community structures if the links of the networks studied in the manuscript are reshuffled in a random way?

We calculated the major statistical distributions for two types of random graphs corresponding to the three systems studied in the manuscript. In the first case, *the degree sequences of the original graphs were preserved* during the randomisation process. We implemented this by link randomisation [8]: in each step two links were selected randomly, and then one of the endpoints of the links were swapped. This process was repeated until on average about a dozen relocations per link was reached. The other type of random graphs we tested were simple Erdős-Rényi random graphs [9] with the same number of

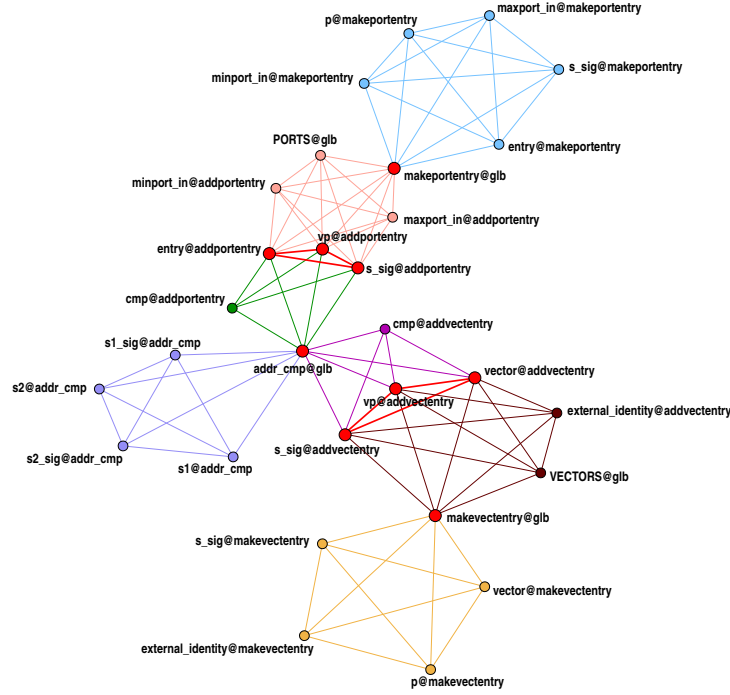


Figure 8: A component in the community graph of the wu-ftp program. Most of these communities are related to functions (sub routines) in the source code. The nodes with a name ending in '@glb' represent global variables. These are likely to appear in several function calls in the source code, hence they are likely to be members in several communities at the same time.

nodes and links as the co-authorship network at $f^* = 0.93$, the word association graph at $f^* = 0.67$ or the protein interaction graph. (The degree sequences in these cases are different from the original ones).

We have found that except for the link-randomised word association graph, *cliques of size larger than three were totally absent in the random networks, therefore, naturally, no k -clique communities for $k > 3$ can exist at all in them.* In comparison the largest clique sizes are 12, 8, 9 and $k = 6, 4, 4$ in the original co-authorship network, word association network and protein interaction network respectively. In Fig. 9. we show the four major statistical distributions for the link-randomised word association network (triangles) compared to the original system (squares, the same as in Fig. 4 in the manuscript). In the randomised system the maximal community size is five (Fig. 9a), the maximal community degree is two (Fig. 9b), the maximal overlap size is one (Fig. 9c), and the maximal membership number is two (Fig. 9d), therefore the corresponding distributions are very truncated compared to the original ones.

In conclusion, we can say that *randomisation severely (in some cases entirely) destroys the observed community structure.* The fact that randomisation can lead to complete loss of communities also implies that they are present in the original system entirely due to specific correlations.

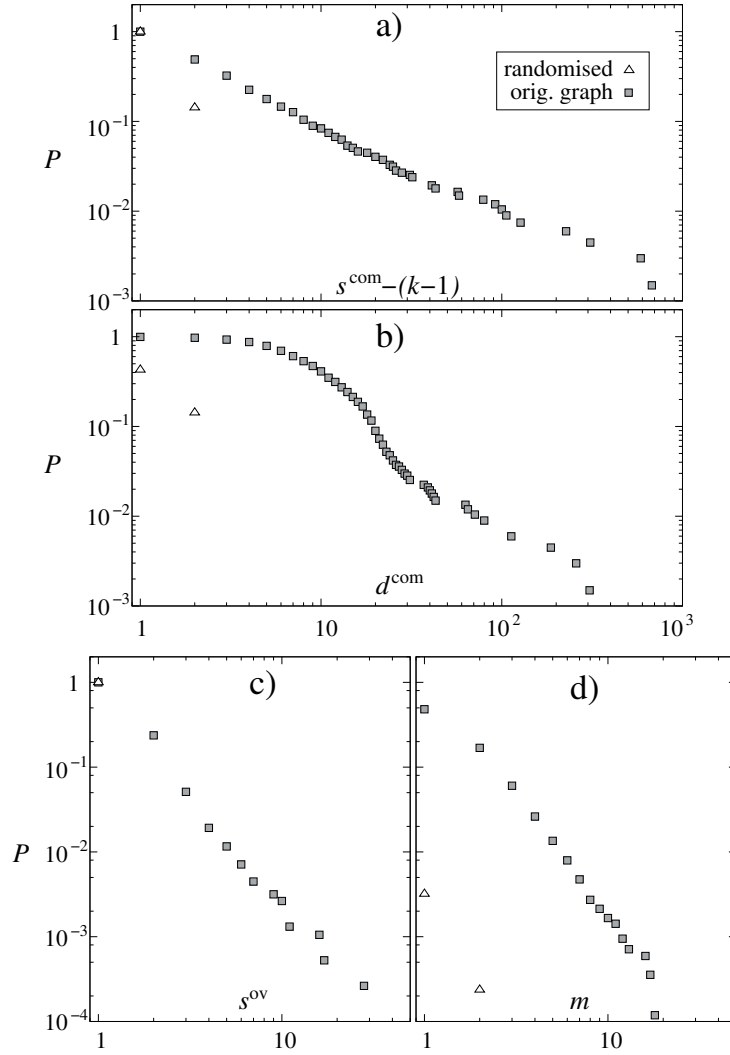


Figure 9: Statistics of the 4-clique-communities for the link-randomised word association network of the South Florida Free Association norm list at $f^* = 0.67$ (triangles), plotted together with the distributions of the original system (squares). The degree sequence was preserved during the randomisation process. (a) The cumulative distribution of the community size, (b) the cumulative distribution of the community degree, plot (c) is the cumulative distribution of the overlap size and (d) is that of the membership number.

References

- [1] M. G. Everett and S. P. Borgatti, Analyzing clique overlap. *Connections* **21**, 49–61 (1998).
- [2] S. Warner, E-prints and the Open Archives Initiative. *Library Hi Tech* **21**, 151–158 (2003).
- [3] <http://arxiv.org/> The co-authorship data were kindly provided by Simeon Warner.
- [4] The data concerning the time evolution of the network of autonomous systems was downloaded from <http://www.cosin.org/extra/data/internet/nlanr.html>.
- [5] <http://www.wu-ftp.d.org>
- [6] <http://www.openoffice.org/>

- [7] Nelson, D. L., McEvoy, C. L., & Schreiber, T. A. *The University of South Florida word association, rhyme, and word fragment norms*. <http://www.usf.edu/FreeAssociation/>.
- [8] S. Maslov and K. Sneppen: *Science* **296**, 910 (2002)
- [9] P. Erdős and A. Rényi, *Publ. of the Math. Inst. of the Hung. Acad. of Sci.* **5**, 17-61 (1960).