

SUPPLEMENTARY INFORMATION

Table of Contents

1	Introduction	3
2	Methods	3
2.1	Link clustering	3
2.1.1	Constructing a dendrogram	3
2.1.2	Partitioning the dendrogram: partition density	4
2.2	Node clustering	5
2.3	Other methods	5
2.3.1	Clique percolation	6
2.3.2	Modularity optimization	7
2.3.3	Infomap	7
3	Properties of link communities	7
3.1	Link communities capture multiple memberships between nodes	7
3.2	Link dendrograms, node hierarchy, and overlap	9
3.3	Partition density	9
3.4	Link communities and fuzzy membership weights	9
3.5	Filtering weighted networks	11
3.6	Examples of link community structure	11
3.6.1	Biological networks	11
3.6.2	Word association networks	13
4	Generalizations and extensions of link communities	13
4.1	Networks with weighted, directed, or signed links	13
4.2	Multi-partite networks	19
4.3	Local methods	19
4.4	Partition density optimization	19
5	Testing community methods	19
5.1	Methodology	19
5.2	Measures	20
6	Network datasets	23
6.1	Overview	23
6.2	Biological networks	23
6.2.1	Protein-protein interaction	23
6.2.2	Metabolic	24
6.3	Social networks	24
6.3.1	Mobile phone	24
6.3.2	Actor	25
6.3.3	US Congress	25
6.4	Other networks	26
6.4.1	Philosopher	26
6.4.2	Word association	27
6.4.3	Amazon.com products	28
7	Validating hierarchical organization	31

7.1	Examples of hierarchical structure	31
7.2	When is hierarchical structure meaningful?	31
7.3	Dynamic dendrogram structure	33
7.3.1	Branching probability	33
7.3.2	Distributions of community sizes and node memberships	33
7.4	Revealing meaningful communities at multiple scales	34
7.4.1	Community quality as a function of cut-level	34
7.4.2	Hierarchical metadata	36
A	Tables of measures	40
A.1	Overall methods	40
A.2	Clique Percolation	41

List of Figures

1	Illustration of similarity measure between link pairs	4
2	Example node and link communities	4
3	Link communities in <i>Les Misérables</i>	6
4	Link communities identify multiple relationships between nodes	8
5	Subtle structural differences are detected by link communities	10
6	Membership and overlap statistics for link communities	11
7	Building link dendrogram intuition	12
8	Comparison of node and link dendrograms	13
9	Simultaneous overlap and hierarchy in a food web and toy model	14
10	Partition density is meaningful	14
11	Statistics for protein-protein interaction networks	15
12	Statistics for other networks	15
13	Filtering dense networks is not necessary for link communities	15
14	Link community structure for the PPI network	16
15	More link communities in the PPI network	17
16	Example link communities around Acetyl-CoA in the metabolic network	18
17	Example link communities in the word association network	18
18	The elements of composite performance	22
19	Network validation results with node clustering	23
20	Political ideology metadata in the US Congress network	26
21	Example metadata for philosopher network	27
22	Example network and metadata for Amazon.com	28
23	Examples of hierarchical structure in the word association network	32
24	Spatial hierarchy of mobile phone users	32
25	Link dendrogram branching probabilities	34
26	Branching probability as a function of window size	34
27	Community and membership distributions for various thresholds	35
28	Example of the hierarchical control	36
29	Community quality across the link dendrograms	37
30	Multi-scale metadata	38
31	Amazon.com network with hierarchical metadata	39
32	Yeast PPI networks with hierarchical metadata	39

List of Tables

1	Modularity values for the test networks as discovered using greedy optimization	7
2	Summary of the 11 network test corpus	30

1 Introduction

This document is organized as follows. Section 2 contains details regarding the implementation of link clustering, as well as the other community detection methods which were used in the main text. In Sec. 3, we discuss properties of link-partitions and important cases, such as “what happens when a link should be a member of more than one community?”, and “what happens in the case of no overlap?”. We show that the link clustering algorithm is able to successfully analyze both cases. Generalizations and extensions of link clustering are discussed in Sec. 4.

The final sections of the document focus primarily on our community validation methodology. To see how meaningful/useful link communities can be, we apply our method to a large corpus of networks, chosen specifically for their diversity and to form a representative sample of common network datasets. First, in Sec. 5, we discuss the measures we use to evaluate different community algorithms. Then, details regarding how the chosen networks were collected and curated, and any particular details regarding how to apply the various validation measures are described in Sec. 6. Section 7 focuses on studying and validating meaningful communities at multiple levels of the link dendrogram. The appendix contains raw data regarding the various quality measures.

2 Methods

Here we offer a detailed discussion of the different methods we have used in this work. In particular we offer additional results about our new link communities and we list implementation details for applying other methods, such as parameter choices. The raw (unnormalized) composite performance scores for all methods are shown in App. A.

2.1 Link clustering

2.1.1 Constructing a dendrogram

The main text has introduced a hierarchical link clustering method to classify links into topologically related groups. Here we provide further motivation for the suggested pair-wise link similarity measure. For simplicity, we limit ourselves to only *connected* pairs of links (i.e. sharing a node) since it is unlikely that a pair of disjoint links are more similar to each other than a pair of links that share a node; at the same time this choice is much more efficient. For a connected pair of links e_{ik} and e_{jk} , we call the shared node k a *keystone* node and i and j *impost* nodes.

If the only available information is the network topology, the most fundamental characteristic of a node is its neighbors. Since a link consists of two nodes, it is natural to use the neighbor information of the two nodes when we define a similarity between two links. However, since the links we are considering already share the keystone node, the neighbors of the keystone node provide no useful information. Moreover, if the keystone node is a hub, then the similarity is likely to be dominated by the keystone node’s neighbors. For instance, if the hub’s degree increases the similarity between the links connected to the hub also increases. This bias due to the keystone node’s degree also prohibits us from applying traditional methods directly to the *line graph* of the original graph, which is constructed by mapping the links into nodes. (Since a hub of degree k becomes a fully connected subgraph of size k in the line graph, the community structure can become radically different.) Thus, we neglect the neighbors of the keystone. We first define the *inclusive* neighbors of a node i as:

$$n_+(i) \equiv \{x \mid d(i, x) \leq 1\} \quad (1)$$

where $d(i, x)$ is the length of the shortest path between nodes i and x . The set simply contains the node itself and its neighbors. From this, the similarity S between links can be given by, e.g., the Jaccard

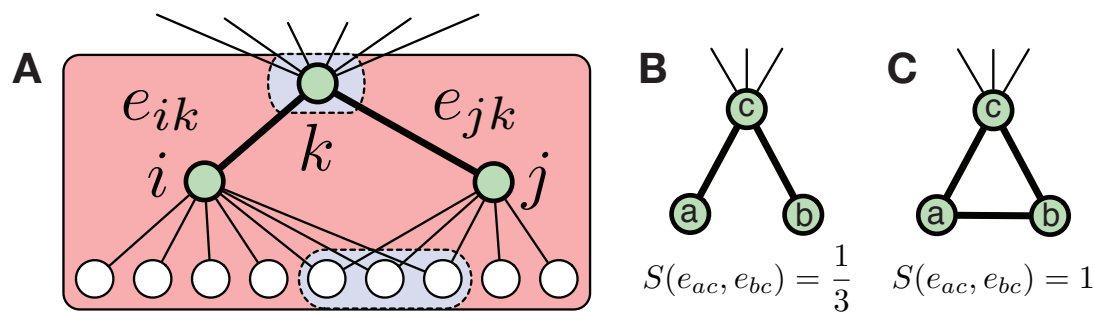


Figure 1: (A) The similarity measure $S(e_{ik}, e_{jk})$ between edges e_{ik} and e_{jk} sharing node k . For this example, $|n_+(i) \cup n_+(j)| = 12$ and $|n_+(i) \cap n_+(j)| = 4$, giving $S = 1/3$. Two simple cases: (B) an isolated ($k_a = k_b = 1$), connected triple (a, c, b) has $S = 1/3$, while (C) an isolated triangle has $S = 1$.

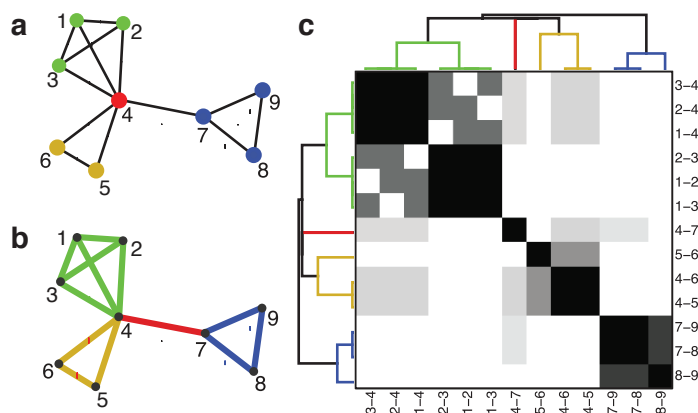


Figure 2: An example network with node communities (a) and link communities (b). (c) The resulting link similarity matrix and link dendrogram. Compare with main text Fig. 1.

index [1]:

$$S(e_{ik}, e_{jk}) = \frac{|n_+(i) \cap n_+(j)|}{|n_+(i) \cup n_+(j)|} \quad (2)$$

An example illustration of this similarity measure is shown in Fig. 1 (see Sec. 4.1 for generalizations of the similarity).

With this similarity, we use single-linkage hierarchical clustering to find hierarchical community structures. We use single-linkage mainly due to simplicity and efficiency, which enables us to apply link clustering to large-scale networks. However, it is also possible to use other options such as complete-linkage or average-linkage clustering. Each link is initially assigned to its own community; then, at each time step, the pair of links with the largest similarity are chosen and their respective communities are merged. Ties, which are common, are agglomerated simultaneously. This process is repeated until all links belong to a single cluster. The history of the clustering process is then stored in a dendrogram, which contains all the information of the hierarchical community organization. The similarity value at which two clusters merge is considered as the strength of the merged community, and is encoded as the height of the relevant dendrogram branch to provide additional information. See Fig. 2 for an example.

2.1.2 Partitioning the dendrogram: partition density

Hierarchical clustering methods repeatedly merge groups until *all* elements are members of a single cluster. This eventually forces highly disparate regions of the network into single clusters. To find meaningful communities rather than just the hierarchical organization pattern of communities, it is crucial to know where to partition the dendrogram. Modularity has been widely used for similar purposes

in node-hierarchies [2, 3], but is not easily defined for overlapping communities.¹ Thus, we introduced a new quantity, the *partition density* D , that measures the quality of a link partition (see Methods, main text). The partition density has a single global maximum along the dendrogram in almost all cases, because the value is just the average density at the top of the dendrogram (a single giant community with every link and node) and it is very small at the bottom of the dendrogram (most communities consists of a single link). This process is illustrated in Fig. 3.

The maximum of D is 1 but it can take values less than zero; $D = 1$ when every community is a fully connected clique and $D = 0$ when each community is a tree. Essentially, D measures how “clique-ish” vs. “tree-ish” each link community is. If a link community is less dense than a tree (when the community subgraph has disconnected components), then that community will give a negative contribution to D . The minimum of D_c is $-2/3$, given by one community of two disconnected edges. Since D is the average of D_c , there is a lower bound of $D = -2/3$.

2.2 Node clustering

We introduce node clustering as a control algorithm to offer a direct comparison to link clustering. In other words, if two algorithms are identical in every possible respect except that one classifies nodes and the other classifies links, how different will their performances be? The node clustering method is closely related to the method introduced in Ravasz et al. [10]. There are many ways to define a similarity between two nodes. We tried four different variations of the node similarity. The four versions are following:

- $S(i, j) = |n(i) \cap n(j)| / |n(i) \cup n(j)|$,
- $S(i, j) = |n(i) \cap n(j)| / \min(k_i, k_j)$,
- $S(i, j) = |n_+(i) \cap n_+(j)| / |n_+(i) \cup n_+(j)|$,
- $S(i, j) = |n_+(i) \cap n_+(j)| / \min(k_i, k_j)$,

where $n(i)$ means the neighbors, not inclusive neighbors, of the node i . Among those, we use the version in Eq. (3) since it finds more relevant communities across most networks we used. In addition, it is the definition most similar to link similarity. Thus, the node similarity is chosen to be

$$S(i, j) = \frac{|n_+(i) \cap n_+(j)|}{|n_+(i) \cup n_+(j)|}, \quad (3)$$

where, as in the main text, $n_+(i)$ are the inclusive neighbors of node i . To determine the node dendrogram, we use the same single linkage hierarchical clustering as we used for clustering links. This node dendrogram is cut at the point of maximum modularity [2]. Since this method is a nice control, but not necessarily applicable in the real world, we study it only in the SI.

2.3 Other methods

In order to evaluate its performance, we compare link clustering to existing, popular community detection methods. We chose three representative algorithms: the clique percolation method (CPM) [11], which is widely recognized as state-of-the-art for detecting overlapping communities; Infomap [12] which is the current state-of-the-art algorithm for detecting non-overlapping communities; and a greedy modularity optimization algorithm [13], which is widely used in the literature.

¹Several modifications of modularity that allow for “fuzzy” communities with relaxed interfaces (or overlapping nodes) to exist [4, 5, 6, 7, 8] have been suggested. However, in order to avoid the trivial optimum, where all nodes are part of all communities, each of these methods *penalize* overlap, and are therefore not suitable for networks with pervasive overlap. (See Fig. 1 of the main text)

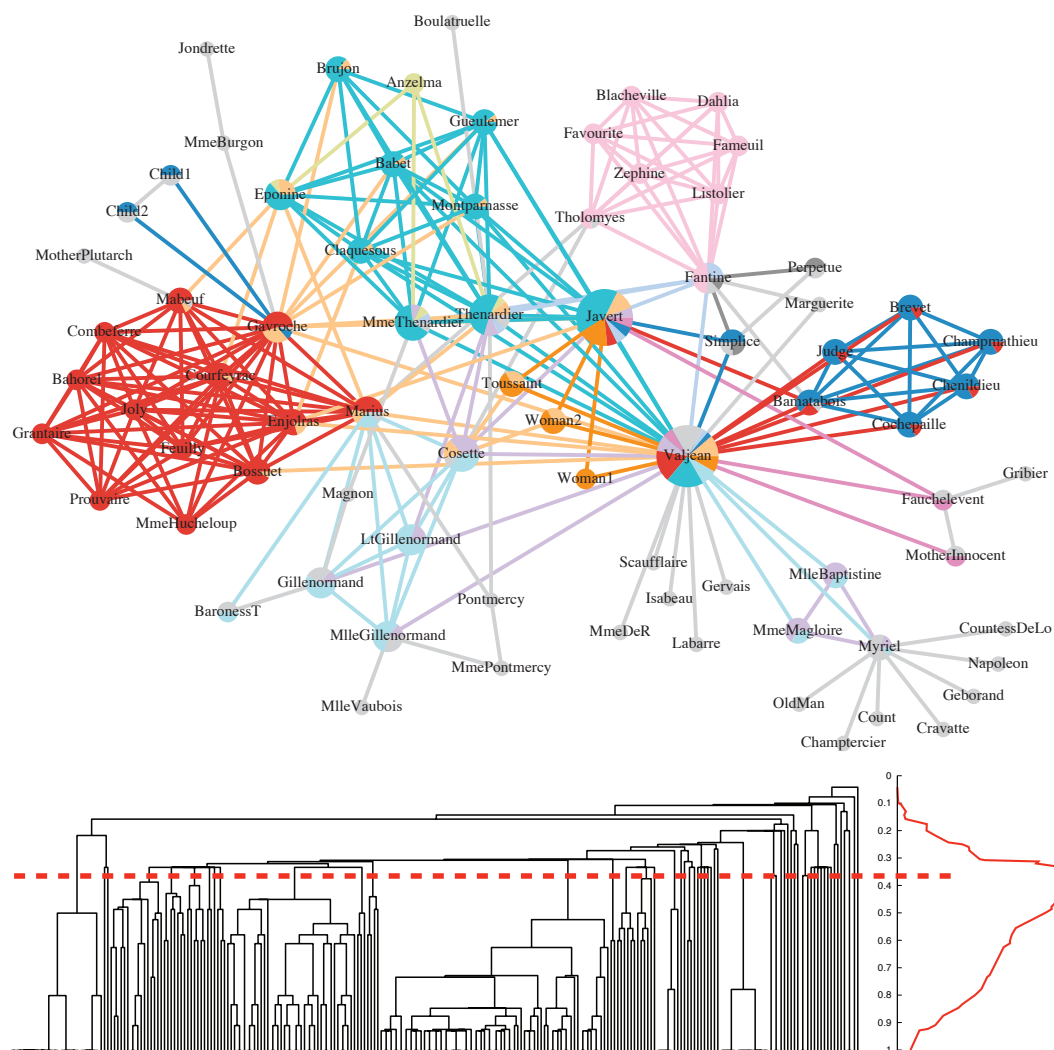


Figure 3: Link communities for the coappearance network of characters in the novel *Les Misérables* [9]. (**Top**) the network with link colors indicating the clustering, with grey indicating single-link clusters. Each node is depicted as a pie-chart representing its membership distribution. The main characters have more diverse community membership. (**Bottom**) the full link dendrogram (left) and partition density (right). Note the internal blue community in the large blue and red clique containing Valjean. Link clustering is able to unveil hierarchical structure even inside of cliques.

2.3.1 Clique percolation

Clique percolation [11, 15] provides an elegant and highly useful method to uncover overlapping community structure [16]. It is currently the most popular and most successful tool available for this task. A particularly interesting feature of this method is that it presents the experimenter with a “knob” k , the clique size, which can be used to tune the result between high coverage, low community quality (sparse communities) and low coverage, high community quality (dense communities). For some networks, such as the mobile phone network, a precedent exists for the choice of k , which we follow. Whenever that is not the case, we have computed the composite performance for a range of k 's and chosen the k which results in the optimum overall performance². This weighs coverage and quality equally, however, and it remains at the discretion of the researcher to decide if this is optimal for his or her application. See Appendix A.2.

²For some of the very large or very dense networks, we were not able to run clique percolation for large values of k with the fastest existing software (even on a machine with 32 Gb of RAM), using the fast algorithm developed by Kumpala et al. [17].

Network	Modularity Q
Metabolic	0.360562
PPI (Y2H)	0.733042
PPI (AP/MS)	0.722658
PPI (LC)	0.864972
PPI (all)	0.728056
Phone	0.652382
Actor	0.867364
US Congress	0.275167
Philosopher	0.454025
Word Assoc.	0.343629
Amazon.com	0.889058

Table 1: The modularity values for the test networks studied in the main text, found using greedy modularity optimization [14]. Many values are very high, indicating that the structure found by the greedy optimization algorithm is highly modular (at least according to the definition of modularity). Good modularity values typically lie between 0.3 – 0.7, while higher values are rare [2].

The main drawback of CPM is its somewhat rigid definition of communities. When a network is very dense, it can become super-critical in the sense of clique percolation, which leads to giant clique communities. At the other end of the spectrum, when the network is too sparse, the network is sub-critical and there are not enough connected cliques to find any communities. For example, in the metabolic network, CPM's coverage is largely due to one giant community containing almost all nodes, leading to a minuscule community quality. Removing this giant community increases the enrichment value, but only $\sim 5\%$ of nodes remain. This situation is not unchanged by increasing clique size. For the Y2H network, however, the problem is sparsity: there are not enough cliques to find structure.

2.3.2 Modularity optimization

To study how typical modularity [2, 18, 19] optimization methods perform, we choose the fast/greedy optimization method of Clauset, et al. [14]. Although this particular modularity algorithm is the most popular one, more accurate methods exist, based on simulated annealing, extremal optimization, and more. (See [3] for additional details.) However, the modularity values we found are often quite high (good modularity values typically lie between 0.3 – 0.7, while higher values are rare [2]), so the lack of accuracy in our comparison is less likely to be from failing to find partitions near the system's maximum modularity. The modularity values found for the test networks are shown in Table 1.

2.3.3 Infomap

The Infomap algorithm [12] is becoming accepted as one of the best and most accurate node partitioning methods [20]. It exploits deep results from information theory and uses a complex, multi-stage optimization scheme. In our application of this method, we used 100 restarts for the large networks (phone, amazon, etc.) and 1000 restarts for smaller networks. The final partition that minimized the map length was then used.

3 Properties of link communities

3.1 Link communities capture multiple memberships between nodes

While clustering links is a much more flexible approach than clustering nodes, one might wonder whether this method is flexible enough—after all, it does not appear to take into account links that appear in multiple contexts (overlapping links). In the main text, we briefly address the issue of multiple relations represented by a single link. Main text Fig. 1f shows that it is very natural that two nodes of

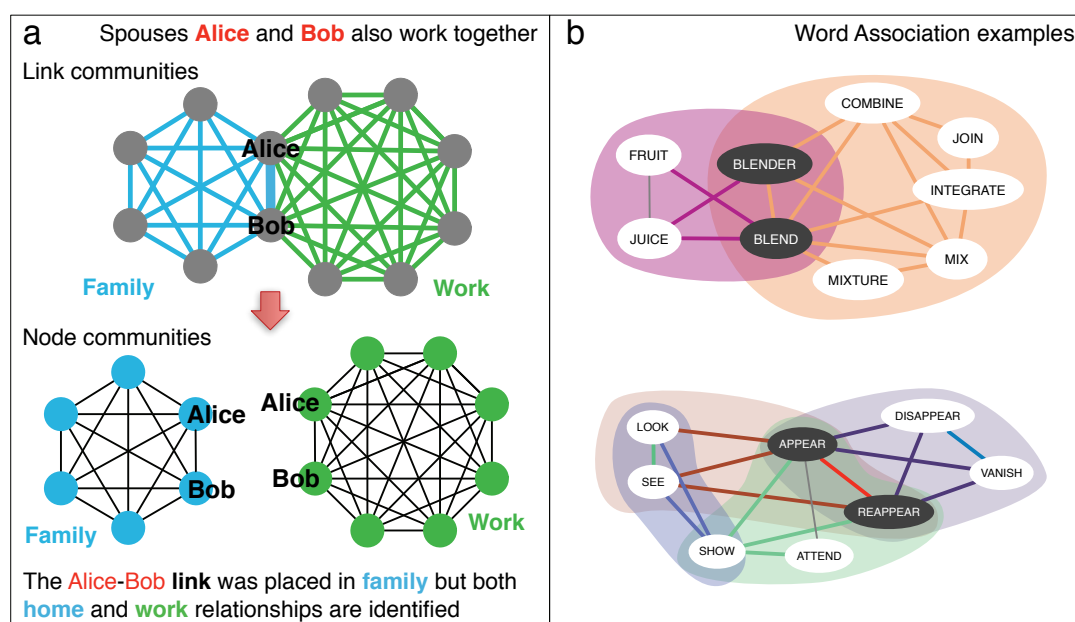


Figure 4: Overlapping links. In the link community framework, a link may be assigned to only one community. By deriving node communities, however, the problem of effectively discovering multiple relationships between nodes is effectively solved. Two nodes can belong to many communities together regardless of the membership of the link between them. Left: illustration of the situation. Right: real examples from word association network. In the upper example, *Blend* and *blender* belong to both ‘fruit juice’ community and ‘mix’ community. In the bottom example, the link between *appear* and *reappear* does not even belong to any of the other communities, but they belong to several communities together.

a given link can simultaneously belong to multiple communities even though the link itself belongs to only one community. Here, we let the examples in Fig. 4 provide further illumination of this point.

The simplistic cases in Fig. 4, however, do not address the complex community structure that arises in real life, where the multiple relationships may include more groups of many nodes and more than one link. Consider a high school with classes of about 30 students. These classes form clusters/communities and are likely to be located by the link community method. Now, students from these classes typically form a number of further communities: Some go to the same class to learn a foreign language, others play on the school’s basketball team, etc. Thus, there will be further overlapping communities in such a way that the members in these new communities are in touch with each other in two distinct ways: through going to the same regular class *and* through playing basketball together. Figure 5 show that the link communities do, in fact, extract these subtle relationships.

It is true that if a group is completely *subsumed* inside another group, and there are *no structural differences* distinguishing this group, such as different connectivity patterns, then link communities will not find the internal group. *No method* will find it, because it’s completely invisible (Fig. 5a). However, if the school’s social network is weighted based on the time students spend together, or if basketball players are slightly more likely to become friends with other basketball players than with students not on the team, or if the team has slightly different external connectivity, these will be identified (Fig. 5b). Notice that the link communities shown in Fig. 5b only separate the player-coach links. This is sufficient to completely identify the basketball team. Figure 5c shows a further example. We also identify these sub-communities in practice; note the ‘clever/wit’ community inside the ‘smart/intelligent’ community in main text Fig. 1f.

What about in practice? Are multiple relationships between nodes rare or abundant in link communities? To answer this, we study the network of communities, where each node is now a community in the original network, and the weights on each link are the number of shared members. The distribution of link weights s_{ov} in this network, studied by Palla *et al.* [11] (we use their notation), explicitly

shows how many nodes participate in the same communities together. (Whenever $s_{ov} > 1$ we have found multiple relationships between two or more nodes.) The broad distributions of s_{ov} in Fig. 6 (top row) show that link communities successfully capture multiple relationships in practice, for both sparse and dense networks. Examining the distribution of the number of community memberships per node m , also studied by Palla *et al.*, we see (Fig. 6 bottom row) that link communities capture a great deal of overlap. (See also Fig. 27.)

3.2 Link dendrograms, node hierarchy, and overlap

A link dendrogram can be very different from a node dendrogram. As an example, consider the graph shown in Fig. 7. Here we have constructed a simple network without overlap, but with two levels of node hierarchy, consisting of four very dense communities, loosely connected into pairs which are then more loosely connected. At the lower level of the link dendrogram, we find six communities, not the expected four. *The reason is that link clustering has correctly identified the two sets of cross-community links as structurally related groups.*

Several prominent methods for finding hierarchical organization exist [23, 22], however, none are able to handle overlap since hierarchical structure always assumes almost disjoint community partitions. For instance, see Fig. 8 for a case where simple overlap prevents node hierarchy from finding true hierarchical structure. Structurally, the red and orange node should be members of the full cliques to which they are connected, but node clustering assigns them to their own community. The situation is more severe than it appears since in a network with pervasive overlap, *all nodes* are in a situation similar to that of the orange and red node. Clique percolation finds overlapping community structure (cliques) in the example network very easily, while the hierarchical random graph model fails to find all of them. Figure 9 illustrates a similar situation.

3.3 Partition density

To support the relevance of the structure found at the optimum partition density, we examine the link communities of the metabolic and mobile phone networks, presented in Fig. 10. Here we show community coverage, the ratio of the number of links within the second largest to largest communities s_2/s_1 , and partition density D , as a function of the dendrogram cut threshold (Fig. 10a). That maxima in D coincide with $s_2/s_1 \rightarrow 1/2$ indicates that discovered link communities are well structured [11, 25]. Likewise, the community size distribution at the optimum D is heavy tailed for both networks (Fig. 10b). These properties suggest that the optimum D is related to a critical point where the link communities are neither fragmented nor gelated. These statistics for the remaining test corpus are shown in Figs. 11 and 12.

3.4 Link communities and fuzzy membership weights

Most fuzzy community methods require membership weights quantifying how strongly a node belongs to a particular community, such that the sum of every node's weights is 1. Link communities can be mapped into fuzzy community memberships simply by counting the number of link membership a node has. If node i with 8 total links has 5 links to community A and 3 links to community B then its membership weights are $w_{iA} = 5/8$ and $w_{iB} = 3/8$.

It is, however, often more natural to consider each node as a full member of its communities. A person's family would be disappointed if anyone proclaimed that he or she was only 1/5th of a member of it; in the metabolic network, it would also be strange to say that H_2O was only 1/200th a member of a given pathway.

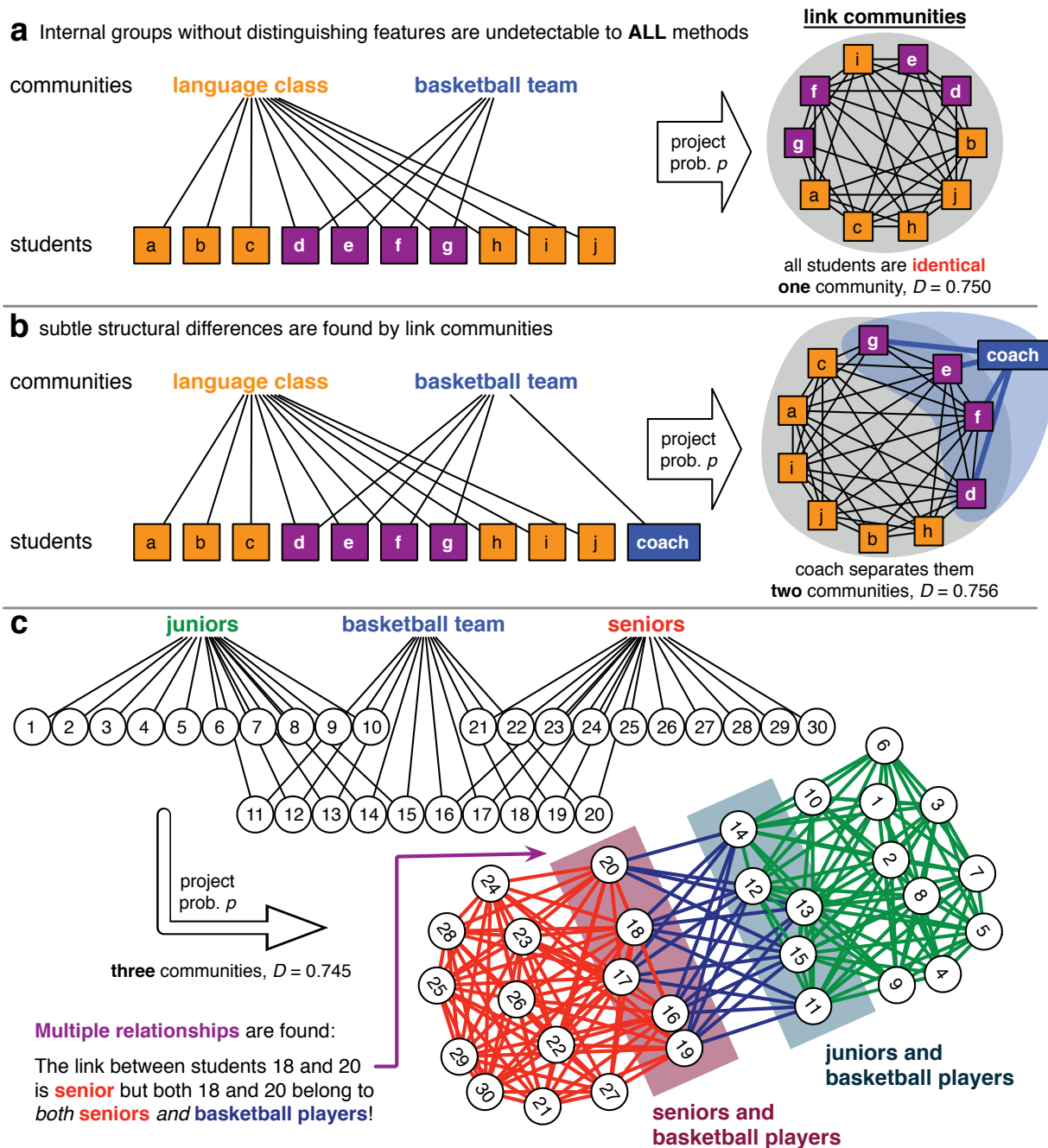


Figure 5: Some small, illustrative examples of the subtle structural changes that link communities detect, using the bipartite social model of [21] with $p = 0.8$, followed by our link communities algorithm. In (a) there are no distinguishing structural features to separate the “subsumed” basketball team from the language class. Detecting the team is impossible for all methods. In (b) however, a single change allows for 100% complete detection. The entire basketball team is successfully found, even though only the coach-team links are separated. It doesn’t take much to achieve the proper node communities. (c) A more extreme example. Class and team detection are again 100% accurate. Very subtle patterns are detectable (see, e.g., the word association communities in main text Fig. 1f and Figs. 3, 7, 14, 15).

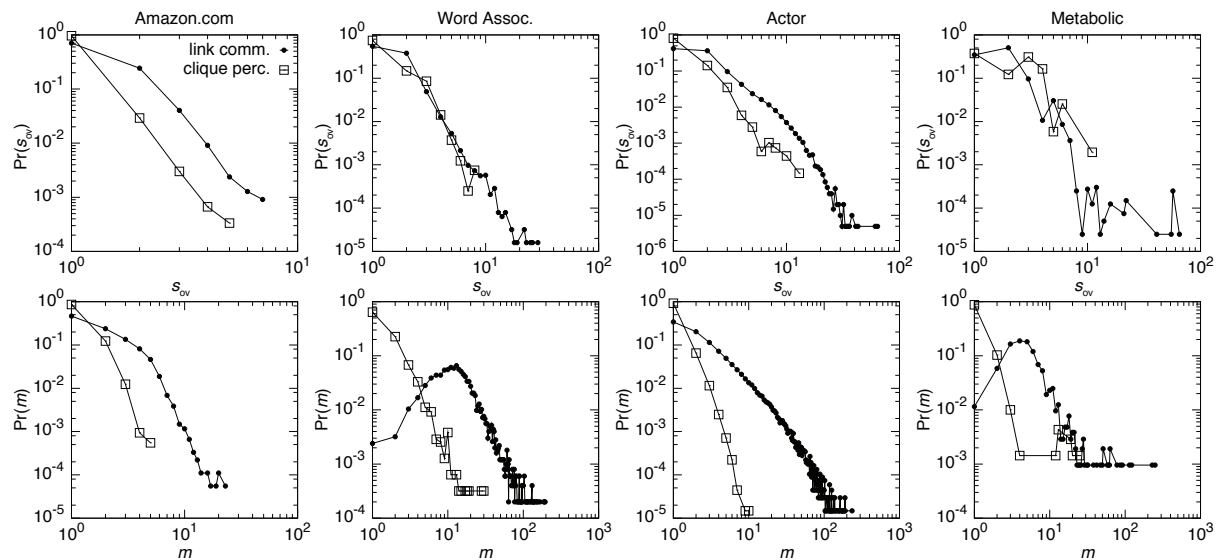


Figure 6: Membership and overlap statistics for link communities in sparse (Amazon.com, actor) and dense (word association, metabolic) networks. Shown are the distributions for overlap size s_{ov} (**top**) and membership number m (**bottom**), as introduced by Palla *et al.* [11]. Link communities were found at the maximum partition density D . We find that link communities extract more highly overlapping communities and a higher average number of overlapping memberships for the denser networks than the sparser ones. The distribution of s_{ov} corresponds to the distribution of weights in the community network. Statistics for clique percolation are shown for comparison (clique size k was chosen from existing literature precedents or else to maximize composite performance).

3.5 Filtering weighted networks

While the networks composing our test corpus are considered unweighted, it may happen that a researcher is presented with a weighted network. A common pre-processing step is filtering the network, deleting all edges below some defined weight threshold. This was done in [11], where the clique percolation method was applied to networks after removing links below some weight w_* . This approach may not be ideal, however, as useful information may be lost.

Since this technique is common, it is important to see how link communities are affected by such filtering. The word association network (Sec. 6.4.2) possesses such weights, and was filtered with $w_* = 0.025$ in Palla *et al.* [11] (using clique size $k = 4$). In Fig. 13 we show the composite performance for the tested methods on the original unfiltered word association network and the thresholded network. Several methods benefit a great deal, but the link communities remain the leader both overall and in community quality. This is strong evidence that link communities are better at dealing with dense networks than other methods, and at exploiting all available information.

3.6 Examples of link community structure

This section contains additional examples of link communities in various networks, all intended to illustrate that link clustering finds meaningful and relevant structure.

3.6.1 Biological networks

Figure 14 shows the community structure around protein YML007W. There are three major communities, all three are related to the transcription process, identified as the mediator complex, NuA4 HAT complex, and SAGA complex [26, 27, 28], respectively. Note the overlapping membership of protein YHR099W, which is already known as a subunit of both the NuA4 complex and the SAGA complex [29, 30, 31]. Figure 15 shows three major communities around the protein YBL041W, which be-

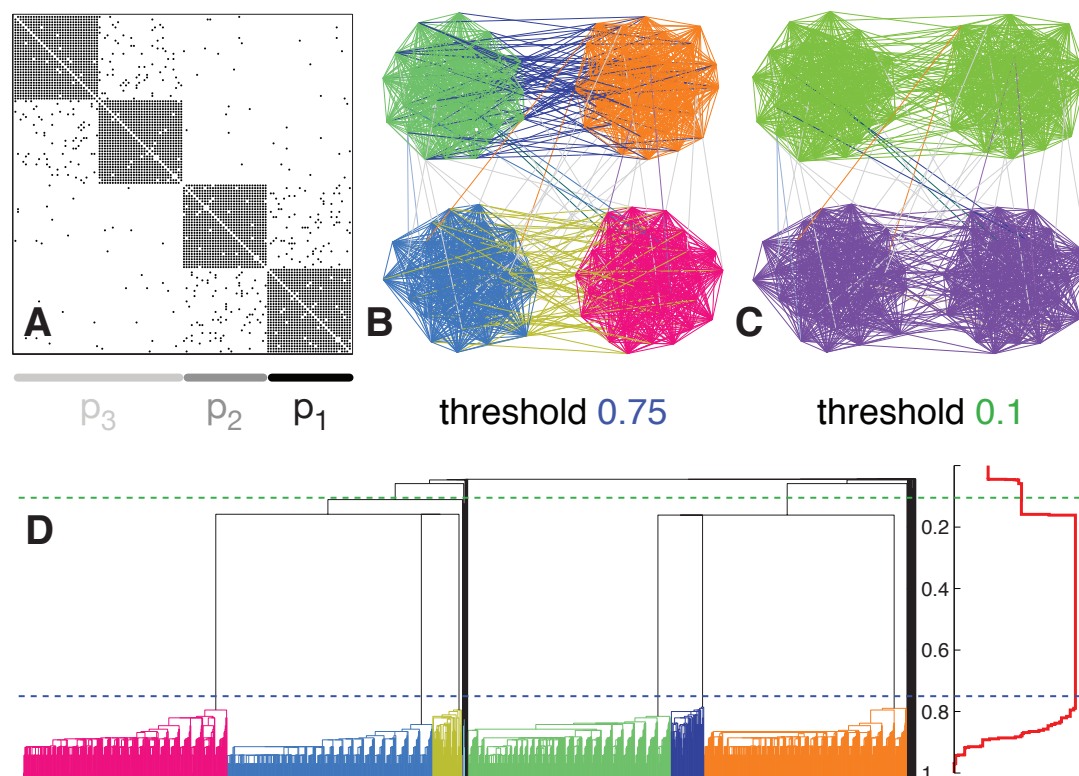


Figure 7: Building link dendrogram intuition. Shown is an example illustrating how hierarchy can be captured at multiple levels of the link dendrogram. **(A)** The 128×128 adjacency matrix for a network of four densely connected non-overlapping communities (each possible link exists with probability p_1), each connected to another community (p_2), and finally the two pairs are weakly connected (p_3). For this example, $p_i = \frac{1-\epsilon}{127-i}$, $\epsilon = 0.02$. The communities at a high **(B)** and low **(C)** threshold, and the full dendrogram **(D)** are shown. The chosen values of p_i lead to a very “stretched” dendrogram and partition density, as expected. While one expects to identify **four** communities at the higher threshold, **six** are actually found, since the inter-community edges are accurately identified by link clustering.

longs to the core of the proteasome complex [32]. We can directly observe that the proteasome consists of two parts: the core and the regulatory particle, and link clustering finds two corresponding communities plus a community connecting the two. As expected from the structure of the proteasome, the core is less exposed to other communities, while the regulatory particle has several connected communities. Likewise, Fig. 16 shows the community structure around Acetyl-CoA, illustrating several roles that Acetyl-CoA plays in the metabolic network.

In addition, we supply in Supplementary Table 1 the list of all communities found by link clustering along with its most relevant GO terms or pathway annotations. For the PPI networks, we use GO-TermFinder [33] version 0.82 to find enriched GO terms and estimate the p -values for each GO term. First, we find all GO terms with p -value less than 0.05, then we pick up only the most significant term for each aspect (biological process, cellular component, molecular function). These terms and p -values are listed along with the community members in Supplementary Table 1. This table shows that more than 80% of communities have at least one enriched GO-term with p -value lower than 0.0001 and more than 30% of communities have at least one enriched GO-term with p -value lower than 10^{-10} .

For the metabolic network, we first filter out communities where less than three members possess pathway annotations. Then, we calculate the enriched pathway annotations shared by the largest number of community members. We compile this information in Supplementary Table 2.

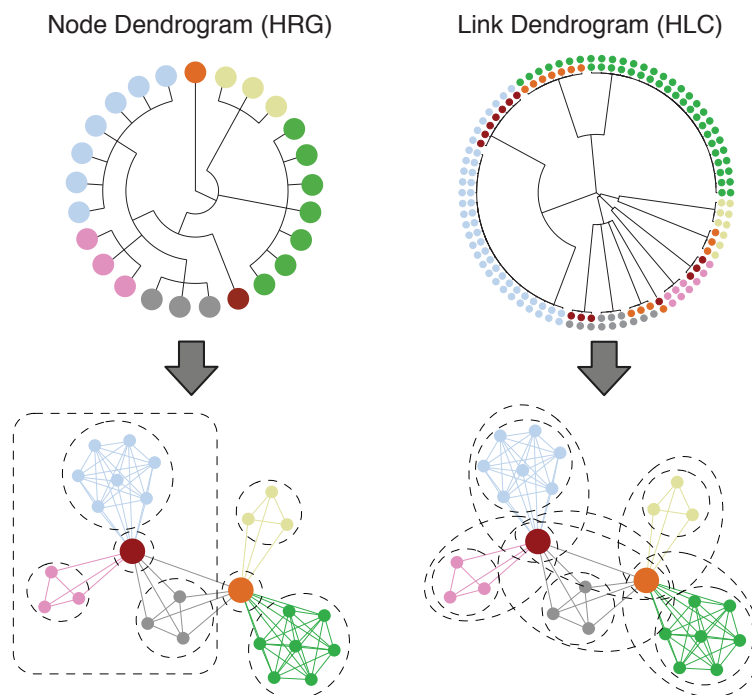


Figure 8: Comparison of a node dendrogram and link dendrogram in the presence of overlap. The node dendrogram is obtained by using the hierarchical random graph (HRG) method (consensus dendrogram) [22], and the link dendrogram is obtained from link clustering. Nodes are colored to distinguish each node or clique and dotted lines represent several hierarchies in the dendrogram. In the link dendrogram, two colored circles at each leaf represent the link between the nodes with the given colors. Note that HRG isolates the red, orange, and gray nodes in the dendrogram, even though they are central to the network and belong to the same clique: one cannot retrieve the full clique communities. In contrast, the link dendrogram captures every clique while at the same time constructing a reasonable hierarchical tree. Note that the links of the red node are placed in appropriate branches of the dendrogram according to their context. Also note the internal hierarchical structures found inside each clique.

3.6.2 Word association networks

We present more examples of link communities in the word association network in Fig. 17. We also attach the list of all link communities found by link clustering at the maximum D in the word association network as Supplementary Table 3.

4 Generalizations and extensions of link communities

4.1 Networks with weighted, directed, or signed links

The similarity between links can be easily extended to networks with weighted, directed, or signed links (without self-loops), since the Jaccard index generalizes to the Tanimoto coefficient [34]. Consider a vector $\mathbf{a}_i = (\tilde{A}_{i1}, \dots, \tilde{A}_{iN})$ with

$$\tilde{A}_{ij} = \frac{1}{k_i} \sum_{i' \in n(i)} w_{ii'} \delta_{ij} + w_{ij} \quad (4)$$

where w_{ij} is the weight on edge e_{ij} , $n(i) = \{j | w_{ij} > 0\}$ is the set of all neighbors of node i , $k_i = |n(i)|$, and $\delta_{ij} = 1$ if $i = j$ and zero otherwise. The similarity between edges e_{ik} and e_{jk} , analogous to Eq. (2), is now:

$$S(e_{ik}, e_{jk}) = \frac{\mathbf{a}_i \cdot \mathbf{a}_j}{|\mathbf{a}_i|^2 + |\mathbf{a}_j|^2 - \mathbf{a}_i \cdot \mathbf{a}_j} \quad (5)$$

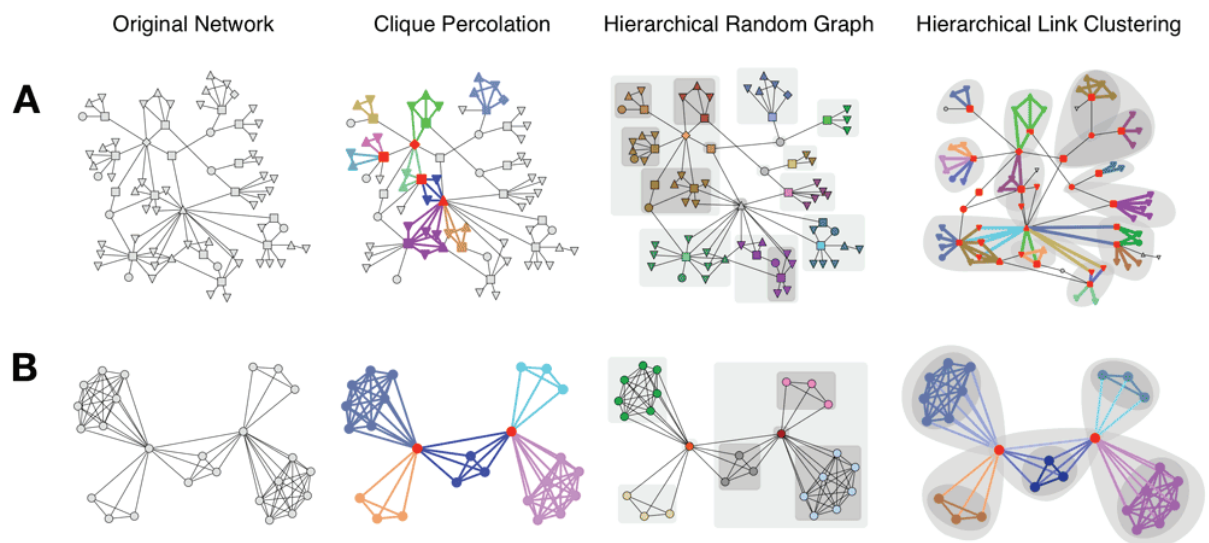


Figure 9: Comparison of methods on a network of UK grassland species interactions [24], which has evident hierarchical structure (A), and on a simple example network with overlapping communities (B). Colors and boxes indicate community structures while nested boxes illustrate hierarchical information. Red nodes possess multiple community memberships. The performance of existing methods depends heavily on the network's structural characteristics. CPM fails to detect the structure in sparse, hierarchical networks (A). The HRG model captures the hierarchical structure in (A) but neglects overlap, and forces the middle 5-clique in (B) to be arbitrarily spread across branches. In the case of hierarchical link clustering, both hierarchy and overlapping structures are well classified. Again, real social networks possess more overlap than in (B).

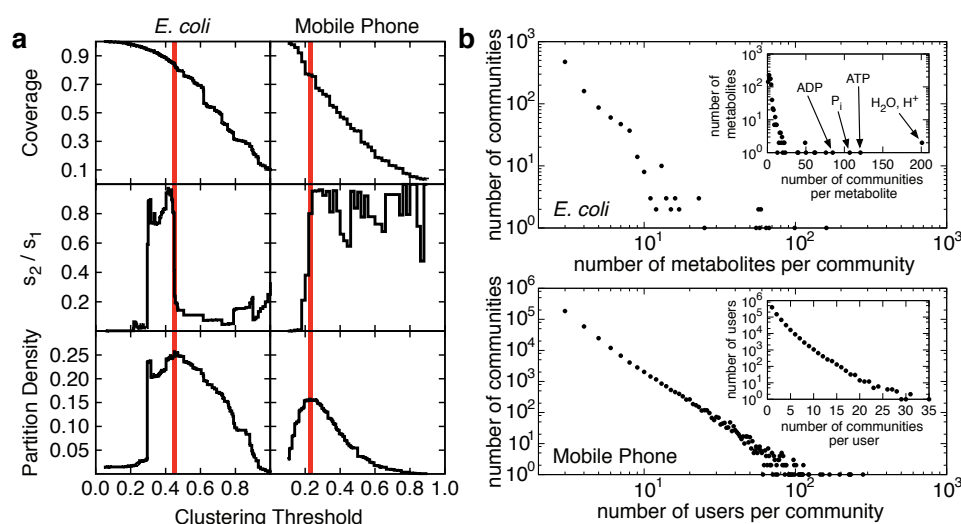


Figure 10: Statistics for the *E. coli* metabolic and mobile phone networks. (a) Community coverage, the ratio of the number of links in the two largest communities, and the partition density D , respectively. In both networks, peaks in D align with $s_2/s_1 \rightarrow 1/2$, implying that the maximum of D corresponds to the percolation transition point where community size exhibits a power-law distribution. (b) The distribution of community sizes and node memberships (insets). The distribution of community size shows a heavy tail. The number of memberships per node is reasonable for both networks: we do not observe phone users that belong to large numbers of communities and we correctly identify currency metabolites, such as water and ATP, that are prevalently used throughout metabolism. The appearance of currency metabolites in many metabolic reactions is naturally incorporated into link communities, whereas their presence hindered community identification in previous work.

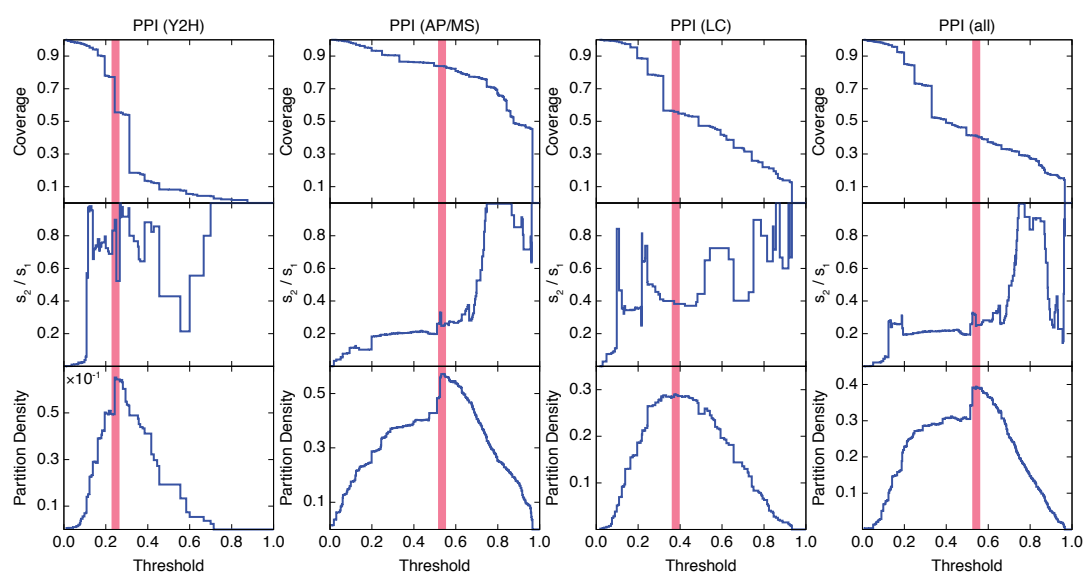


Figure 11: Several statistics for the **protein-protein interaction** networks, as a function of the link dendrogram cut threshold. Compare with Fig. 10a.

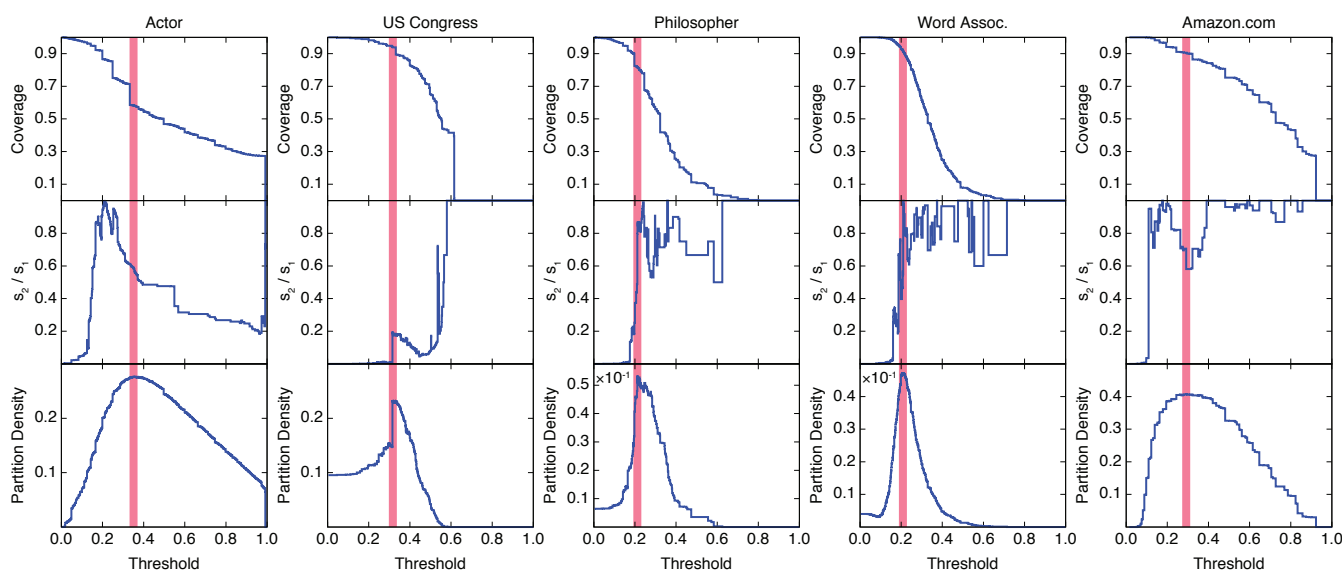


Figure 12: Several statistics for the remaining corpus networks. Compare with Fig. 11 and Fig. 10a.

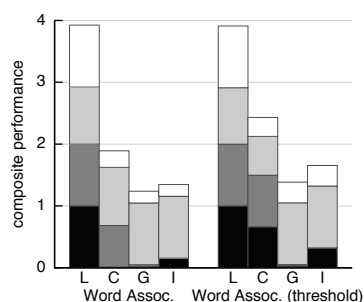


Figure 13: Thresholding or filtering a weighted network is not critical for link communities, whereas other methods benefit from this procedure. (Symbols and colors as per main text Fig. 2. Here we show the composite performance for the original word association network (**left**) and the same network after thresholding weak links (**right**). For the thresholding we use $w_* = 0.025$, the same value used in [11], as well as $k = 4$ for clique percolation. Clique percolation, particularly its community quality (black), greatly improves. We see that the link community procedure is robust to “noisy” links, unlike other approaches, and actually benefits from all available information.

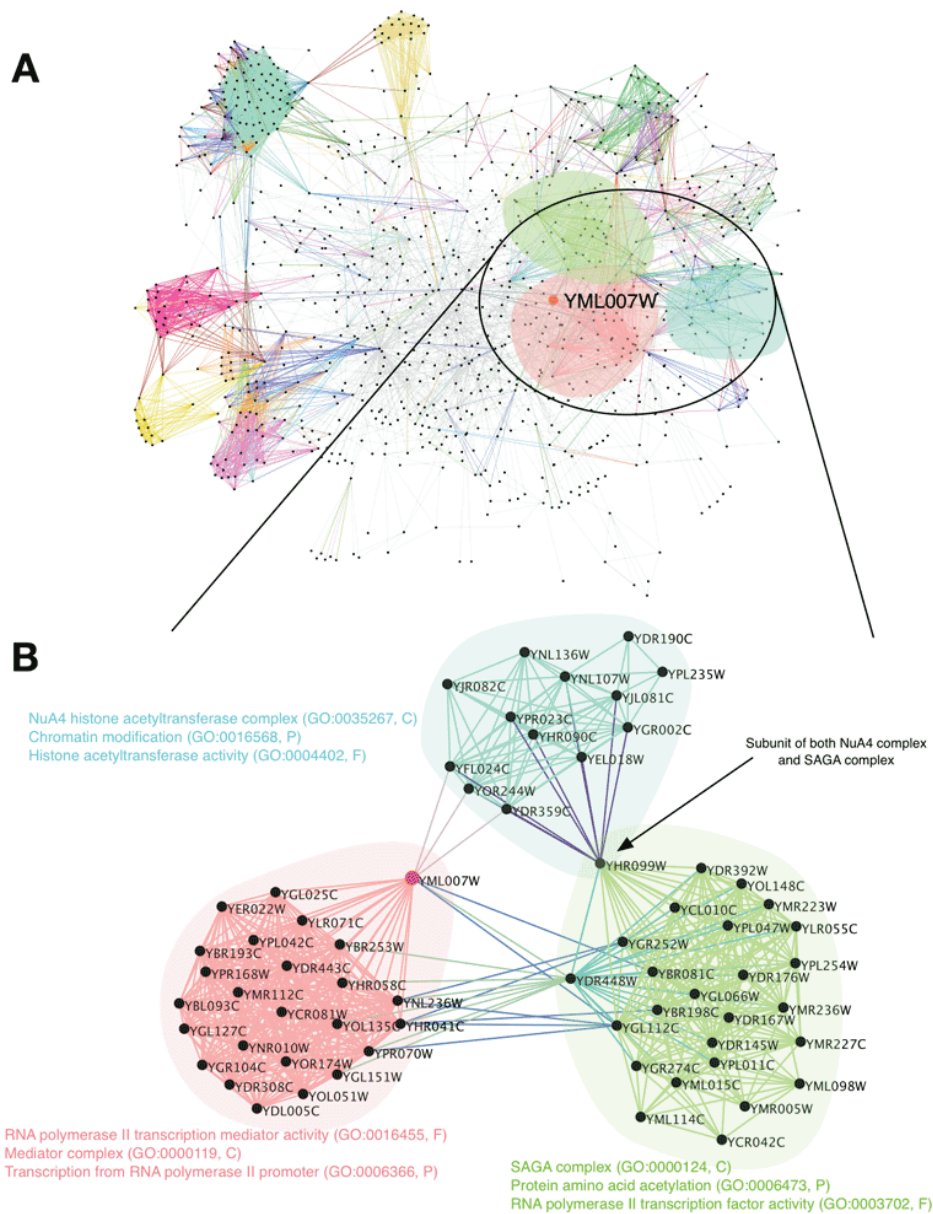
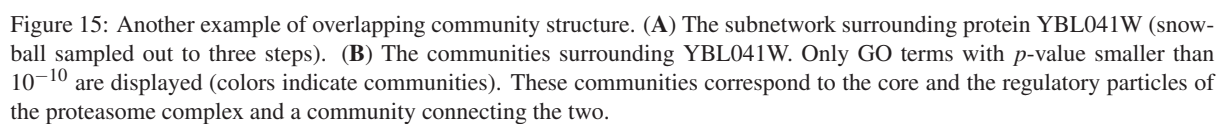


Figure 14: An example of overlapping community structure in the PPI compendium network. (A) The subnetwork surrounding protein YML007W (snowball sampled out to three steps). (B) The communities around YML007W. Only GO terms with p -value smaller than 10^{-10} are displayed (colors correspond to communities).



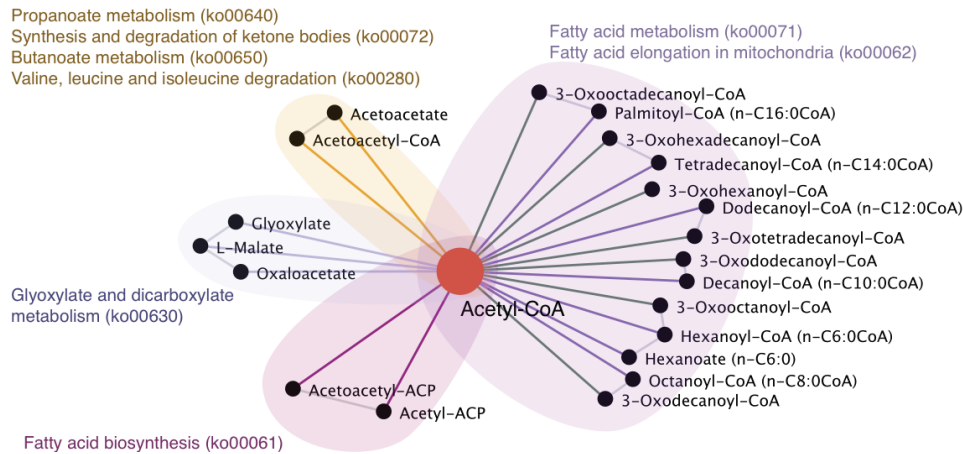


Figure 16: Overlapping community structure around Acetyl-CoA in the *E. coli* metabolic network. Acetyl-CoA plays several different and important roles in metabolism. Shown are only communities with homogeneity score equal to 1 (all compounds inside each community share at least one pathway annotation); all other links, including those that contribute to community structure, are omitted. Pathway annotations shared by all community members are displayed with corresponding colors. The two communities to the right of Acetyl-CoA are grouped since they share the same exact pathway annotations.

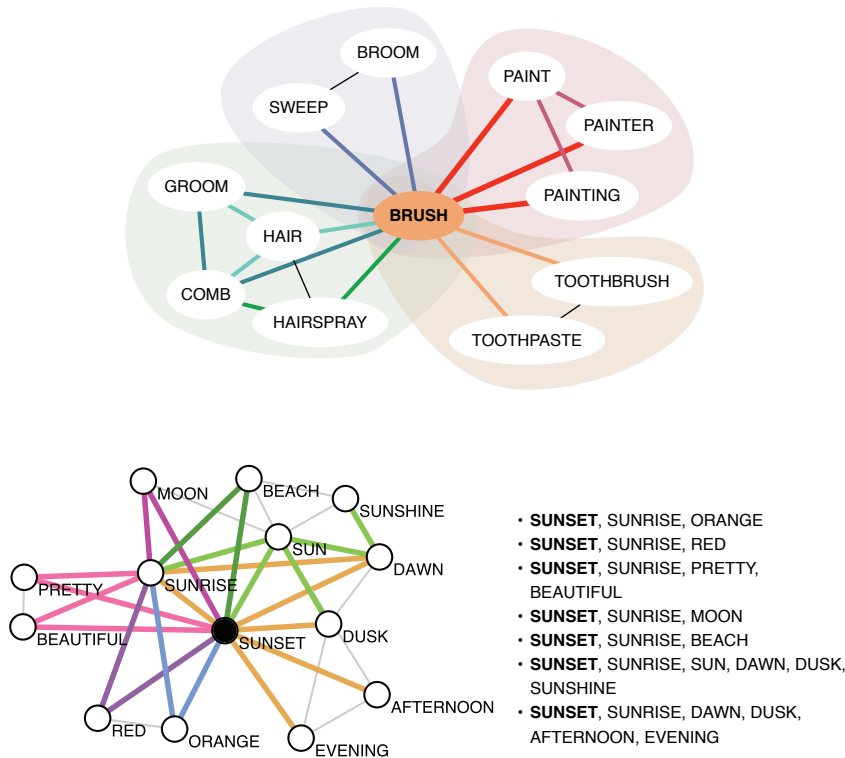


Figure 17: More link community examples in the word association network. Top: link communities successfully captures various meanings of the word BRUSH. Bottom: Link communities captures diverse associations of the word pair SUNRISE-SUNSET. The translated node communities are listed.

4.2 Multi-partite networks

A multi-partite network is a network in which the nodes can be divided into K disjoint sets and all links must terminate in two distinct sets. This creates additional constraints on the existence of certain edges which must be accounted for in both the link similarity and the partition density.

Link similarity: The similarity measures, Eqs. (2) and (5), depend only upon connectivity, and therefore automatically account for multi-partite structure. The one change necessary is incorporating the forbidden connections between the same kind of nodes, which can be achieved by using the set of neighbors instead of the inclusive neighbor set when calculating the similarity.

Partition density: We must modify the definition of partition density since a fully connected K -partite clique is much sparser than a clique in a unipartite network. In general, the K -partite partition density of a subset c can be written as

$$D_c^{(K)} = \frac{m_c + 1 - \sum_k n_c^{(k)}}{\sum_k \left(n_c^{(k)} \sum_{k' \neq k} n_c^{(k')} \right) - 2 \left[\left(\sum_k n_c^{(k)} \right) - 1 \right]}, \quad (6)$$

where the index k runs over the K node types and the notation $n_c^{(k)}$ refers to nodes of type k . The full partition density is achieved by summing over individual communities, $D^{(K)} = 2M^{-1} \sum_c m_c D_c^{(K)}$.

4.3 Local methods

Since our definition of similarity between links only uses local information, a local version [35, 36, 37] of link clustering can be trivially obtained. One can simply choose a starting *link*, compute its similarity S with all adjacent links, agglomerate the one with the largest S into the community, compute any new similarities between edges inside the community and bordering it, and repeat. A stopping criteria to determine when the community has been fully agglomerated is still necessary [36]. For instance, one can monitor the partition density as links are agglomerated, in order to establish a reasonable community boundary. Another, simpler approach is to fix the similarity threshold and agglomerate only links with similarity larger than that threshold. To find all the overlapping communities of a node one can simply begin the above methods with each of that starting node's links or start from one link, find its community (which may end up including another starting node link), then pick another unassigned link from the starting node, find that community, and repeat until all the starting node's links are contained within communities.

4.4 Partition density optimization

Since the partition density is a quality function of link community structures in networks, it is possible to find link communities by direct optimization. Begin by assigning links to communities at random, then use, e.g. simulated annealing. The fact that link communities are disjoint partitions enables us to apply many traditional optimization techniques to find overlapping communities.

5 Testing community methods

5.1 Methodology

Our goal is to provide a fair evaluation of all the community methods we test. Unfortunately, evaluation of community structure in real networks is akin to a “chicken and egg” problem: since we don't know what the actual communities are, we must use algorithms to try and discover them. But if we don't know the real communities, how can we determine if the found communities are any good?

While common in the biological sciences, where enrichment analysis or similarity analysis using annotations (e.g. GO terms) is the standard method to assess computational predictions about a group of proteins, quantitative validation using real-world networks has not been a common practice in community research. Even the most widely cited, state-of-the-art papers about community identification do not provide quantitative validation, but only provide qualitative arguments with one or two small networks that are small enough to draw and look at the structure [12, 38, 39, 11]. A recent survey paper [40] about community structure, although very extensive, does not contain even a single section regarding quantitative validation using real-world networks.

Some literature has answered the problem of validating community detection methods using model graphs (*benchmarks*) designed to generate a random, pre-programmed community structure as “ground truth”. However, since the community structure in these graphs reflects the *conceptual model* of communities held by their creator, there is no guarantee that the results can be extrapolated to real networks. Worse, this approach introduces serious biases towards the algorithms that conform with the same conceptual model as the benchmark graphs and are directly biased *against* other theories of community structure.

For instance, every existing benchmark graph has the underlying principle that a community should have more intra-community links than outgoing links, which is not true in networks with pervasive overlap. Furthermore, no existing benchmark graph takes into account the highly non-random abundance of triangles, one of the most important and fundamental characteristics of real world networks, and one of the earliest discoveries of the complex networks field [41]. The randomized nature of current benchmark graphs shows evident bias against algorithms such as clique percolation [11], which exploits these triangles (and cliques) and is based on a different community definition than modularity [2], which is the conceptual model behind current benchmark graphs.

To avoid requiring the hidden “ground truth” communities, we have focused on networks that possess descriptive *metadata*. This information does not directly contribute to the construction of the network, but it allows us to understand what the nodes in the network do, how similar they are to one another, and how many contexts or roles each node has. An example of a network and its metadata is presented in Fig. 22. Using these metadata to describe how similar nodes are within communities (community quality, see Sec. 5.2), we can compare and contrast the results of different methods, relating how much each method’s results tell us about the relevant (hidden) metadata.

5.2 Measures

There are some subtle aspects to consider when comparing disparate community algorithms. Some methods find excellent communities (high quality) but only for a very small fraction of the network (low coverage). Others find medium-quality communities but classify the majority of the network. Some methods find overlapping memberships, others do not. Since it is difficult and unfair to compare all methods along any one of these directions, we have introduced a simple *composite performance* measure to fairly account for these differences while also allowing a researcher to focus on the individual aspects.

We study four distinct aspects of the quality and coverage of the communities found—the quality measures are based on metadata and the measures of coverage focus on the amount of information extracted from the network.

Community Quality. Many of the networks studied here possess metadata that attaches a small set of *annotations* or *tags* to each node. For example, in the Amazon.com network, each product is categorized into several subjects (see Figs. 18, 22); each actor’s career in the Actor collaboration network can be described by a set of plot keywords; each protein in the Protein-Protein Interaction networks is given a set of GO-terms, which describe the biological process that the protein participates in. Assuming that these metadata form a description of the node, beyond the network

itself, we can reasonably state that “similar” nodes share more metadata than dissimilar nodes. To quantify this, we compute, e.g., the *enrichment of node pair similarity*:

$$\text{Enrichment} = \frac{\left\langle \mu(i, j) \right\rangle_{\substack{\text{all } i, j \text{ within} \\ \text{same community}}}}{\left\langle \mu(i, j) \right\rangle_{\text{all pairs } i, j}}, \quad (7)$$

where $\mu(i, j)$ is a metadata-based similarity between node i and j whose exact definition depends on the particular network (each similarity is discussed in detail in Sec. 6). In other words, enrichment is the average metadata similarity between all pairs of nodes that share a community, divided by the average metadata similarity between all pairs of nodes³. The denominator serves as a baseline similarity and larger values of enrichment show that the communities are “tighter,” according to the metadata. Note that it is important to compare all pairs of nodes, not just links, since links themselves are often enriched beyond average, depending on the properties of the metadata. See Fig. 18, top left.

This approach is very similar to that used in [42] to quantify the relevance of interactions.

Overlap Quality. For each node i in the network, we extract from the metadata a scalar quantity (call this the overlap metadata) that we expect to be closely related to the number of true communities that node i participates in. For example, in the word association network, each community corresponds to a set of words that share the same general topic. The more definitions a word has, the more topics the word is expected to belong to. In the metabolic network, the number of reaction pathways that a metabolite participates in corresponds to the number of communities (contexts or roles) of the metabolite.

To rigorously quantify the amount of information gained by community algorithms, we use *mutual information* to relate the number of memberships and the overlap metadata. This quantity tells us how much information about the true overlap of a node is gained by knowing or learning the number of communities that a particular method has assigned to the node. Mutual information works well since detected relationships need not be linear or obey a predisposed functional form. By running multiple algorithms and computing this mutual information, we can see which methods let us know the most about the overlap metadata. Note that even non-overlapping methods may learn information about the overlap metadata, since some nodes may be placed within zero communities. See Fig. 18, bottom left.

Community Coverage. To measure community coverage, we simply count the fraction of nodes that belong to at least one community of three or more nodes. A size of three was chosen since it is the smallest *nontrivial* community. This measure provides a sense of how much of the network is analyzed. See Fig. 18, top right.

Overlap Coverage. Two algorithms may both completely classify a network, giving complete coverage, but one method may extract more information by finding many more densely overlapping communities than the other. It is therefore important to consider overlap coverage as well as community coverage. To do so, we count the average number of memberships in nontrivial communities that nodes are given. For non-overlapping community methods, both coverage measures are identical. This measure shows how much information is extracted from that portion of the network that the particular algorithm was able to analyze. See Fig. 18, bottom right.

³For very large networks or very large communities, we may not be able to test every possible pair of nodes. In this case, if the network is more than around 1M nodes, we compute the baseline from 10^7 randomly chosen pairs of nodes. Likewise, for communities of more than 1000 nodes, we chose 10^5 random pairs to compute the numerator in Eq. (7).

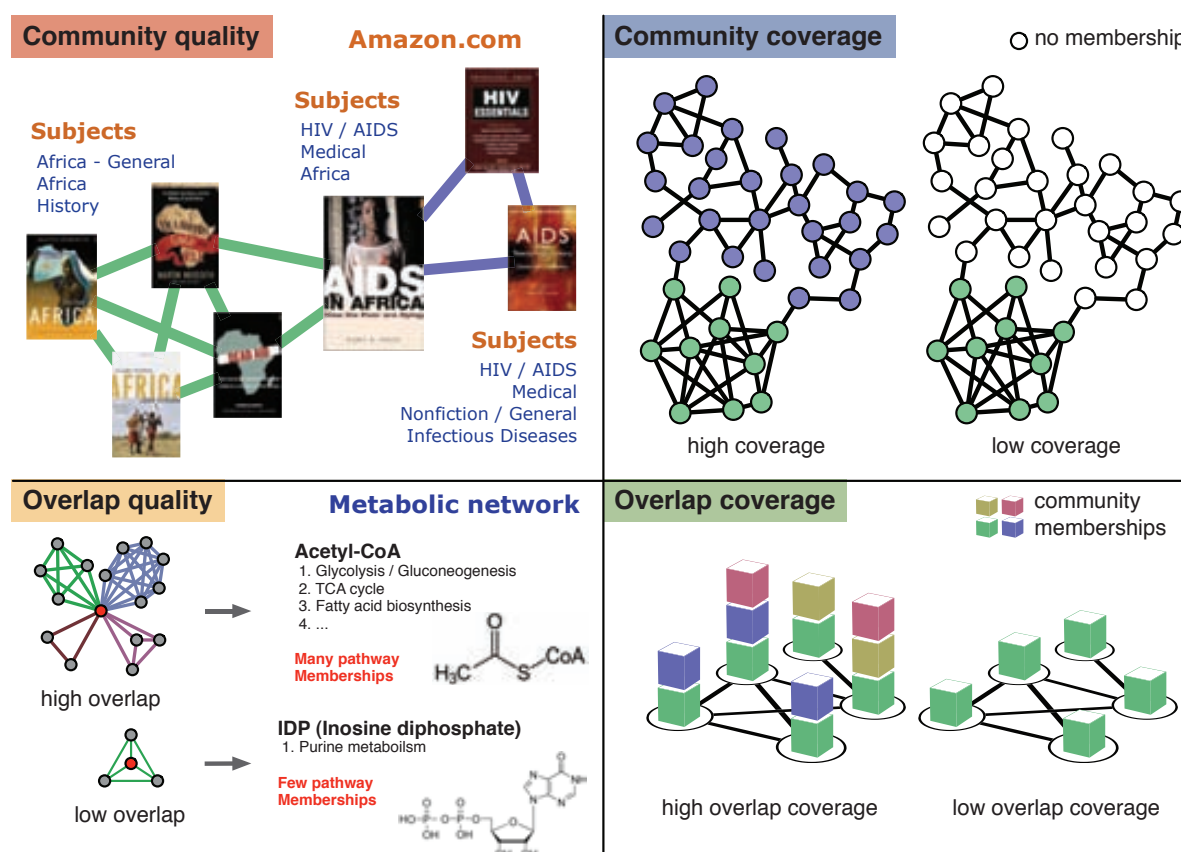


Figure 18: The elements of composite performance. **(top left)** *Community quality* measures the similarity between nodes within each community compared to a null model, based on metadata. **(bottom left)** *Overlap quality* compares the amount of overlap found for each node with a measure of real-world overlap, based on metadata. **(top right)** *Community coverage* is simply the fraction of nodes categorized by the algorithm. **(bottom right)** Two methods may have the same community coverage but one may extract many more overlapping memberships and will yield more information about the network. Thus we introduce *Overlap coverage*, the average number of memberships per node. This is equivalent to community coverage for non-overlapping methods.

Note that the evaluation of the community and overlap quality include neither trivial communities nor singleton nodes, since their absence is considered by the coverage measures.

For many networks, these measures do not necessarily fall between 0 and 1. For example, in the Amazon.com product network and the word association network, link communities find enrichments 80–100 times higher than the global baseline. Therefore, we renormalize all community and overlap quality values such that the maximum value is 1 for the best performing method⁴. This allows us to directly compare performance across networks whose metadata similarities may cover vastly different ranges of values. Likewise, overlap coverage is often greater than 1 for overlapping methods; these values are likewise rescaled. Community coverage is also renormalized, although there is typically always one algorithm that yields complete coverage and the values are already constrained to [0, 1].

We are now left with four measures quantifying the performance of each algorithm. In order to provide a clean, simple representation of each algorithm's performance, we show a stacked bar chart summing all four measures. Since each measure is normalized to have values between 0 and 1, so that the best method for each measure has a value of 1, the maximum composite performance will be 4. Note that this composite performance measure weighs each of the four aspects equally, while providing a

⁴If a method happens to yield a negative value for a particular measure, all the methods are subsequently scaled such that the minimum value is 0.

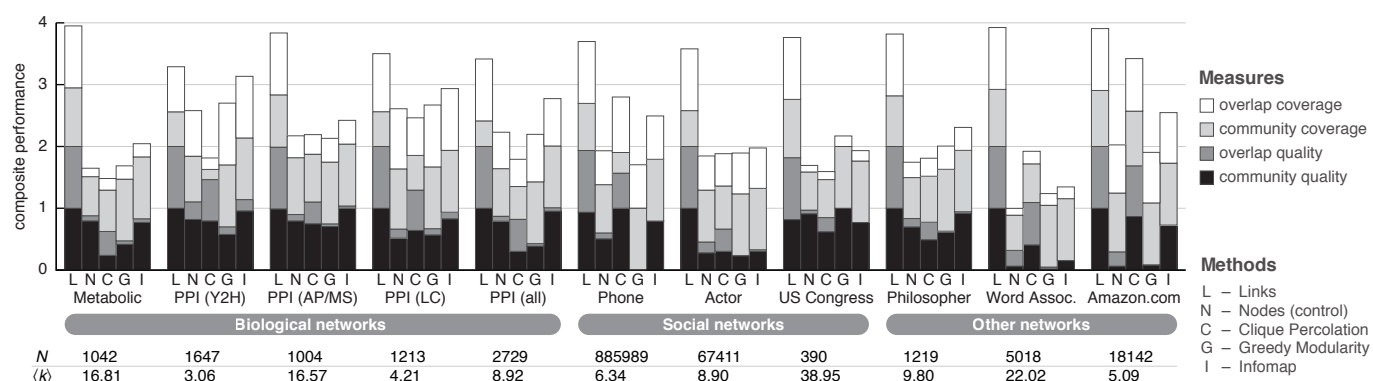


Figure 19: Data-driven evaluation of community algorithms over a large corpus of real networks. (Compare with main text Fig. 2, which lacks the node clustering control algorithm.) Each column represents an algorithm's **composite performance**, measuring community/overlap accuracy and sensitivity. Also shown for each network is the number of nodes N and the average degree $\langle k \rangle$. Link communities achieve the best performance in every network.

simple and easily understood bar chart that nevertheless allows the researcher to evaluate the individual merits of each performance criterion. We find this stacked representation simpler to understand than multiple bar charts while still presenting sufficient information to be fair to all aspects of the problem. Results are shown in Fig. 19 (compare with main text Fig. 2).

6 Network datasets

6.1 Overview

Here we discuss the network datasets used throughout this work, including properties of their metadata, how they were collected, and how the metadata was used to compute the composite performance. Table 2 summarizes all the networks used in this study.

We have chosen eleven networks to test (one is the union of three other networks). This test set contains some of the most relevant networks in recent network research: protein-protein interaction networks for *S. cerevisiae* [42], the metabolic network reconstruction of *E. coli* [43], and a large, dynamic social network derived from mobile phone telecommunication records [44, 45, 46, 47]. A variety of other networks were also chosen to serve as diverse test topologies, representative of the diverse datasets used in complex networks research, and to enable the comprehensive validation procedure of Sec. 5, due to their rich metadata. Table 2 includes brief descriptions of this associated metadata.

6.2 Biological networks

6.2.1 Protein-protein interaction

We analyzed the protein-protein interaction (PPI) network of *S. cerevisiae*, the most studied PPI network.

Construction We use a recently published dataset of PPI networks compiled into three genome-scale networks: yeast two-hybrid (Y2H), affinity purification followed by mass spectrometry (AP/MS), and literature curated (LC) [42]. We also use the union of these three networks (PPI (all)). We use only the largest component of each network.

Metadata We use the Gene Ontology (GO) terms as metadata for the PPI network. The GO project is “a major bioinformatics initiative with the aim of standardizing the representation of gene and gene product attributes across species and databases.” [54] And it provides controlled vocabulary (GO terms) which describes certain aspects of protein characteristics (function, location, etc).

We choose GO terms as the most reasonable metadata for PPI networks, since they are the most elaborate protein annotations available, provide structured information along with statistical information for each term, and there are established methods to calculate the functional similarity between proteins.

Community quality We adopt the same measure as the paper that published the datasets [42]. First, a p -value that two proteins share similar GO terms by chance is calculated using GO biological process terms and the total ancestry measure [55]. The similarity between two proteins $\mu(i, j)$ is defined as either one (if $p < 10^{-3}$) or zero (if $p \geq 10^{-3}$). Then, the enrichment of functionally similar pairs is calculated using Eq. (7):

Overlap quality We use the total number of GO terms as a proxy for the amount of overlap, since it is likely that a protein with many GO terms functions in more diverse contexts. We compute the mutual information between the number of GO terms and the number of discovered memberships as overlap quality.

6.2.2 Metabolic

We use a metabolic network reconstruction of *E. coli* K-12 MG1655 strain (iAF1260), one of the most elaborate metabolic network reconstructions currently available [43].

Construction From the metabolic network reconstruction iAF1260, we retain only cellular reactions, ignore information regarding the compartments (cytoplasm and periplasm), and project the network into metabolite space (two metabolites are connected if they share a reaction). For instance, if an enzyme catalyzes the reaction where metabolites A and B are transformed into C and D , the resulting network would contain a clique of A, B, C , and D .

Metadata We use the pathway annotations from KEGG database [56], which is one of the most widely used metabolic network databases. Each metabolite has zero or more metabolic pathway annotations. For instance, Acetyl-CoA is annotated with 38 pathways including Glycolysis, citrate cycle, and fatty acid biosynthesis.

Community quality To measure the similarity between a pair of metabolites a and b , we calculate the Jaccard index between their pathway sets, i.e. $\mu(a, b) = |P_a \cap P_b| / |P_a \cup P_b|$, where P_m is the set of pathways that contain metabolite m . With this similarity, the community quality is then calculated using Eq. (7).

Overlap quality The number of pathways represents the number of contexts that a given metabolite participates in. We measure the mutual information between the number of pathways and the number of community memberships found by the algorithms.

6.3 Social networks

6.3.1 Mobile phone

This dataset catalogs approximately 8 million users, all calls among these users, and the locations of users when they initiate a phone call (the tower from which the call originated). Self-reported demographic information such as age and gender is also available for some users.

Construction We generate the social network by constraining the location to a 350 km by 80 km region and two nodes in the region are connected only if they each call the other person at least once during a 30-week period. We assign to each user a single location, that of the tower they most frequently used. The final network contains approximately 2.8 million links.

Community quality Unlike most other networks, we do not possess tags for each node, but instead the nodes are embedded spatially, using each phone user's most likely location. To compute the similarity between nodes, we use the euclidean distance between their most likely locations, hypothesizing that social contact is more frequent for users that are geographically related. Since nodes with higher similarity have *smaller* distance, we do not use Eq. (7), but instead:

$$\text{Community quality} = 1 - \frac{\langle d(i, j) \rangle_{\substack{\text{all } i, j \text{ within} \\ \text{same community}}}}{\langle d(i, j) \rangle_{\text{all pairs } i, j}}, \quad (8)$$

where $d(i, j)$ is the euclidean distance between the most likely locations of nodes i and j .

Overlap quality To quantify how much information was discovered about the amount of overlap, we use the total number of phone calls each user made during the observation window. This operates under the assumption that frequent phone users may fulfill broader roles in their social networks.

6.3.2 Actor

For this network, we use the Internet Movie Database (IMDb) to find working collaborations between film actors. We focus on actors who star in at least one movie during the years 2000 and 2009, and at least two movies during their entire career. Television shows, video games, and other performances were not used.

Construction The raw IMDb files were downloaded from <http://us.imdb.com/interfaces> on 2009-12-08. From this data, we construct a bipartite network of movies and actors. We remove films and actors who do not satisfy the above criteria and then project the bipartite network onto the actors, creating a network where two actors i and j are linked with a weight w_{ij} if they co-star in w_{ij} films. Finally, we remove projected links with weights $w < 2$ and keep only the largest connected component. By ensuring that the actors have appeared together in at least two films, we increase the likelihood that they developed a working relationship.

Community quality Associated with each film is a set of plot keywords. We can roughly summarize each actor's career during 2000–2009 by taking the union of all the keywords of the movies that actor appeared in. Since many keywords are very finely grained, we consider only those that label at least 100 films (over the entire IMDb dataset). The Jaccard index between these sets is then used as the node-node similarity in Eq. (7) to compute the “keyword enrichment” of each community algorithm.

Overlap quality One option for overlap metadata is to use the *seniority* of the actor, defined as the year of his or her first film role (not necessarily during 2000–2009). We expect actors with longer careers to be professionally capable of participating in more collaborative groups. The mutual information between the number of communities an actor belongs to and the first year of his or her career is then used to quantify this relationship.

6.3.3 US Congress

The network of legislative collaborations between US congressional representatives (not senators) during the 108th US congress (2003–2005).

Construction Using the dataset of [49, 50]⁵, we construct a bipartite network B of representatives and the legislative bills they (co)-sponsored. Many bills are co-sponsored by the majority of repre-

⁵Downloaded from <http://jhffowler.ucsd.edu/cosponsorship.htm>.

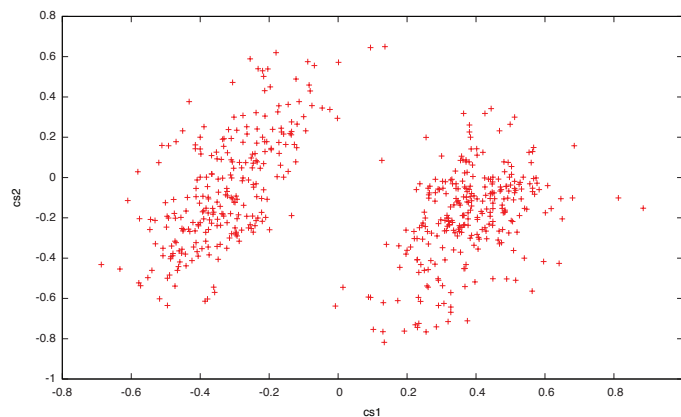


Figure 20: Scatter plot of the Common Space Scores for the 108th US Congress (House and Senate). The ideological and political breakdown is visible in the clustering of the points, which closely follow party lines (republicans and democrats).

sentatives and there were many bills introduced (7765 total), so projecting this bipartite network onto the representatives results in a very dense, nearly complete graph. To avoid this, we filter out edges to capture only the tightest working relationships. To do this, we apply two filtering criteria. First, we remove all introduced bills that contain more than 10 (co)-sponsors total. This network is then projected onto the representatives to form network G_1 . Meanwhile, we also project the unfiltered B onto the representatives and then delete all links with weights less than 75, forming network G_2 . The final network G that we feed to the community detection algorithms is then the intersection of G_1 and G_2 , i.e., each link in G must exist in both G_1 and G_2 . This network is still fairly dense but was disconnected, so we focus on only the giant connected component. This is why there are only 390 representatives.

Community quality Associated with each representative are two values between -1 and 1 known as the *common space score* [51, 52]. These values form a two-dimensional space where distances capture political and ideological similarity (Fig. 20). The first dimension generally represents liberal/conservative bias while the second is related to women's rights and abortion issues. We simply compute the euclidean distance between pairs of points as the node-node similarity measure, and compute the overall “enrichment” of an algorithm's communities using Eq. (8).

Overlap quality For the overlap metadata we use the seniority of each congressional representative, measured as the number of elected terms that person has served. We roughly expect that longer-serving representatives will more easily participate in multiple collaborations than those who are newly elected. The mutual information between the number of community memberships and the number of elected terms is then used to quantify this relationship.

6.4 Other networks

6.4.1 Philosopher

Network of famous philosophers and their philosophical influences, as recorded by users of the english-language Wikipedia⁶.

Construction The raw data consists of the file `enwiki-latest-pages-articles.xml` containing all articles in Wikipedia per 2009-12-02, 22:35:45, which was obtained from the site's download section⁷. Wikipedia maintains a list of all philosophers, sorted by name⁸. This set of names

⁶<http://en.wikipedia.org>

⁷<http://download.wikimedia.org/enwiki/latest/>


⁸http://en.wikipedia.org/wiki/Lists_of_philosophers

Retrieved from "http://en.wikipedia.org/wiki/Alfred_North_Whitehead"

Categories: 20th-century philosophers | American philosophers | English mathematicians | Academics of Imperial College London | Academics of University College London | Harvard University faculty | Ontologists | Old Shirburnians | English philosophers | People from Ramsgate | Logicians | Metaphysicians | Philosophers of science | Western mystics | 1861 births | 1947 deaths

- This page was last modified on 13 December 2009 at 15:12.

Figure 21: The network of **philosopher**'s and their philosophical influences, as captured by Wikipedia. Here we show the *infobox* for mathematician and philosopher A. N. Whitehead (**right**), and the categories that his page is grouped into (**top**), many of which represent his chosen profession. The bottom of the infobox lists the other philosophers who influenced his work and the philosophers who were later influenced by him. The page also has a collection of hyperlinks to other wikipedia pages, which we use to quantify the similarity between pairs of philosophers.

Alfred North Whitehead	
	
Full name	Alfred North Whitehead
Born	February 15, 1861
Died	December 30, 1947 (aged 86)
Era	19th century philosophy 20th century philosophy
Region	Western Philosophy
School	Process Philosophy
Main interests	Metaphysics, Mathematics
Notable ideas	Process Philosophy
Influenced by	Kant, Bergson, Plato, James, Dewey
Influenced	Gilles Deleuze, Philip Clayton, Charles Hartshorne, Latour, Maurice Merleau-Ponty, Bertrand Russell, Wolfgang Smith, Isabelle Stengers, Mordecai Kaplan, William Irwin Thompson

forms the nodes of the philosopher network; an example is shown in Fig. 21. Internal Wikipedia hyperlinks between philosophers form the network links⁹.

Community quality Associated with each philosopher's webpage is the set of all (internal) Wikipedia hyperlinks. Besides links to other philosophers, used to build the network, each page has many hyperlinks to philosophical concepts, philosophical schools of thought, time periods, geographical areas, and so on. We expect more similar philosophers to have more Wikipedia pages in common, so we use the Jaccard index between these sets as the node-node similarity measure in Eq. (7).

Overlap quality Each philosopher is placed into a number of categories (see Fig. 21 top). We expect that philosophers that belong to more categories will participate in more communities, due to their broader interests, etc., though the relationship is not necessarily linear. The mutual information between the number of community memberships and the number of categories is then used to quantify this relationship.

6.4.2 Word association

This network is constructed from existing datasets about free association of word pairs [57]. This dataset is not only interesting as is, but also acts as a nice testbed for community identification: Since nodes are plain english words, we can qualitatively evaluate how reasonable each community is just by looking at the members of a community. This network is quite dense and possesses pervasive overlap.

Construction The dataset was created at the University of South Florida and University of Kansas [57]. They presented 5,019 stimulus words to more than 6,000 participants and asked them to write the first word that came to mind. For instance, if you hear the word *cheddar*, you will almost certainly think about the word *cheese*. They gathered all of these word pairs and assigned a weight that

⁹ Another choice of links between philosophers would have been the set of links listed under *Influenced by* and *Influenced* in the philosopher 'infobox' (see Fig. 21). However, most of the articles describing lesser known philosophers do not have infoboxes, so in order to work with the largest possible dataset, we chose to use all internal hyperlinks.



Figure 22: Example of the network and available metadata for the **Amazon.com** product co-purchases network. Here we show a particular book, some of the books it is often bought with, the set of subjects it is classified into by Amazon.com, and the set of popular “tags” Amazon.com users have chosen to describe or annotate the book’s content. We can use shared tags to quantify how similar pairs of books are, and the more subjects a book has, the more communities it might be expected to belong to. Other combinations of metadata are certainly possible. Other networks have similar quantities.

represents how frequently two given words are associated. This data itself is a weighted, directed network between words. We reduce this network into an undirected, unweighted network by ignoring weight and direction (cf. Palla *et al.* [11]).

Metadata We use the WordNet database for the metadata [53], assigning a set of meanings/definitions or *senses* to each word (known as *synsets*). Since this database was specifically built for semantic analysis, each detailed meaning of a word has a unique ID, which enables quantitative analysis.

Community quality We define a pair of words to be similar when they share at least one meaning ID, i.e. $\mu(i, j) = 1$ if i and j share at least one meaning, 0 otherwise. Then the community quality is defined using Eq. (7).

Overlap quality We calculate the mutual information between the number of meanings for each word and the number of non-trivial community memberships for the node.

This network was previously studied using clique percolation in [11]. They used clique size $k = 4$ but first removed all edges with weights less than $w_* = 0.025$. Here we consider the unweighted, unfiltered network and so instead use $k = 5$, which gives much higher quality k -clique communities and improved composite performance. In Sec. 3.5 we discuss this filtering, and show results for $k = 4$ with and without weight thresholding (Fig. 13).

6.4.3 Amazon.com products

Products that are frequently purchased at the same time by customers at Amazon.com. The Amazon Web Service (<http://aws.amazon.com/>) provides a tool to programmatically access information about any given product sold on their website. For a particular product, we retrieve the top five most frequently co-purchased products, the set of tags or annotations that users have applied to describe the product, and the list of subjects the product is sold under. The former is used to construct the network while the latter two are used for metadata. See Fig. 22 for an example product.

Construction Using Amazon.com's XML web service, on 2009-12-24 we performed a breadth-first search (BFS) crawl (or snowball sample) of co-purchased products by repeatedly retrieving encountered products' top five co-purchases (along with relevant metadata), starting from the number one bestselling book at the time, *THE HELP* by Kathryn Stockett. This crawl continued out to depth $d = 12$. At the final layer of the BFS snowball, many nodes may point to unexplored products at the next step. These unexplored products are removed from the network, since we do not know their connectivity, resulting in a final network of $N = 18142$ nodes. This network is interesting not only because of the rich metadata that is available but also because this snowball sampling technique does not completely capture the network yet is a common approach when sampling dynamic web data. Likewise, since Amazon.com only returns the top five most co-purchased products, the network's degree distribution is not accurate (we treat the final network as being undirected). This provides an interesting test to see how reliant or customized a community method is to the broader degree distributions that are commonly encountered.

Community quality Each product is associated with a set of keywords or annotations known as tags. These tags were applied by users of the website and describe the product, e.g., the plot or characters of a book. The Jaccard index between the sets of tags was used as the node-node similarity in Eq. (7) to compute the overall "tag enrichment" for each algorithm.

Overlap quality Similar to user tags, each product is associated with a set of subjects categorizing it. We expect that products with more subjects will belong to more communities due to the broader nature of the product, as well as user purchasing interests. Thus we use the number of subjects as the overlap metadata and compute the mutual information between the number of communities and number of subjects. This tells us how much we have "learned" about the subjects a product belongs to merely by learning the number of communities the algorithm has placed the product into.

Reversing this metadata choice (using subjects for community quality and number of tags for overlap quality) does not qualitatively alter our composite performance results, indicating that our test procedure is not reliant on particular metadata.

network	description	N	$\langle k \rangle$	metadata	
				community	overlap
PPI (Y2H)	PPI network of <i>S. cerevisiae</i> obtained by yeast two-hybrid (Y2H) experiment [42]	1647	3.06	Set of each protein's known functions (GO terms) ^a	The number of GO terms
PPI (AP/MS)	Affinity purification mass spectrometry (AP/MS) experiment	1004	16.57	GO terms	GO terms
PPI (LC)	Literature curated (LC)	1213	4.21	GO terms	GO terms
PPI (all)	Union of Y2H, AP/MS, and LC PPI networks	2729	8.92	GO terms	GO-terms
Metabolic	Metabolic network (metabolites connected by reactions) of <i>E. coli</i>	1042	16.81	Set of each metabolite's pathway annotations (KEGG) ^b	The number of KEGG pathway annotations
Phone	Social contacts between mobile phone users [45, 46, 47]	885989	6.34	Each user's most likely geographic location	Call activity (number of phone calls)
Actor	Film actors that appear in the same movies during 2000–2009 [48]	67411	8.90	Set of plot keywords for all of the actor's films	Length of career (year of first role)
US Congress	Congressmen who co-sponsor bills during the 108th US Congress [49, 50]	390	38.95	Political ideology, from the common space score [51, 52]	Seniority (number of congresses served)
Philosopher	Philosophers and their philosophical influences, from the English Wikipedia ^c	1219	9.80	Set of (wikipedia) hyperlinks exiting in the philosopher's page	Number of wikipedia subject categories
Word Assoc.	English words that are often mentally associated [53]	5018	22.02	Set of each word's senses, as documented by WordNet ^d	Number of senses
Amazon.com	Products that users frequently buy together	18142	5.09 ^e	Set of each product's user tags (annotations)	Number of product categories

^aGO terms are “structured, controlled vocabularies (ontologies) that describe gene products in terms of their associated biological processes, cellular components and molecular functions in a species-independent manner.” See http://wiki.geneontology.org/index.php/GO_FAQ

^bKEGG database provide metabolic pathway annotations for metabolites. See <http://www.genome.jp/kegg/>

^cThese influences are treated independently from the global wikipedia hyperlink structure and are particularly easy to extract for philosopher biographies.

^dSee <http://wordnet.princeton.edu/wordnet/man/wngloss.7WN.html>

^eAmazon.com's XML Service only returns the five most co-purchased products, though considering the network as undirected will boost some node degrees. This artificial constraint makes the network to have very narrow degree distribution, and serves as a unique test set.

Table 2: A brief description of the networks used in the paper. Shown are the number of nodes N , the average degree $\langle k \rangle$, and brief descriptions of the metadata available to study node similarity and the expected amount of overlap. Full details in Sec. 6.

7 Validating hierarchical organization

The main text shows that link communities present an excellent way to reconcile the apparently disparate notions of hierarchy and overlap, something which has not been accomplished before. As illustrated in main text Fig. 1 and Fig. 8, it is impossible to find a node hierarchy that captures any pervasively overlapping community structure, even in a very simple case. In this sense, the current approach contrasts with all other hierarchical community methods, because our approach—link communities—is a straightforward way to unify hierarchy and overlap.

In most of the examples used in the main text, we pick out a scale (determined by the maximal partition density D), resulting in the set of ‘best’ communities to study. However, we believe that the choice of a best level of communities is often made because the tools to analyze hierarchy are not as advanced as the tools for communities and that the full structure is currently more difficult to deal with, and not because the best level is the only level worth exploring.

Here, we elaborate on the part of the main text showing that the *best* level of a hierarchy is not the *only* level worth exploring. This is true in many domains: For example, faculty, staff, and students at a university may organize at multiple scales, from schools (school of science, school of business, etc.), down to the departmental level (physics department, chemistry department, etc.) and then further down to research groups and small-scale collaborations. The most modular structure may form at, say, the departmental level, but the structures of both smaller research groups and larger school-wide organizations are still relevant.

In the main text, the evidence for this point is contained in main text Fig. 4. Below, we present additional evidence for the presence of meaningful, multi-scale structure represented in the link dendrogram, as well as results for the full network corpus. A small number of networks possess metadata about the hierarchy itself, so we also provide alternative evidence for the existence of such structure in those networks.

No previous methods have captured pervasively overlapping structures across multiple system levels; the combination of pervasive overlap and meaningful community structure on multiple levels of the dendrogram is the multi-scale complexity to which we refer throughout the text.

7.1 Examples of hierarchical structure

Before we begin a quantitative analysis, it is useful to qualitatively inspect samples of the detected hierarchical organization. Here we choose the word association network to illustrate the multi-scale hierarchical structures; in other networks, it is more difficult to appreciate the meanings of communities and their hierarchical organization since we are less familiar with the node labels.

We use two approaches to decipher complex, hierarchical structure. One is tracking how a single link forms larger and larger super-communities (bottom-up) and the other is drilling down into the sub-communities of a large community (top-down). As shown in Fig. 23, both perspectives clearly (but qualitatively) illustrate the success of the link dendrogram in capturing the network’s meaningful communities at multiple levels.

Figure 24 presents a further example of the spatial hierarchy of link communities within the mobile phone network, expanding on that shown in main text Fig. 4.

7.2 When is hierarchical structure meaningful?

We begin by noting that finding a hierarchical tree does not necessarily imply the discovery of meaningful structure; one can always build a random tree, for example. The hierarchical tree is only meaningful when the encoded structure is relevant to the system being studied.

To show that the link dendrogram contains meaningful structure at multiple levels, we now investigate the following:

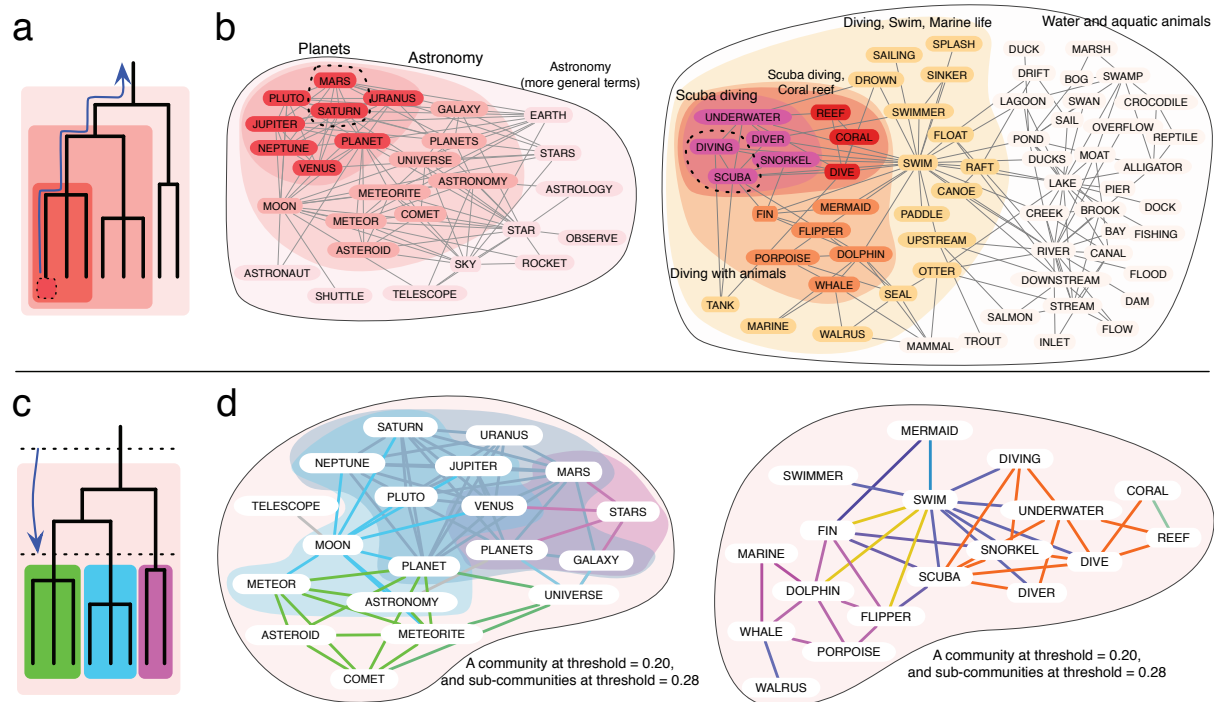


Figure 23: Examples of hierarchical structure in the word association network. The word association network is a nice example for this purpose, since it is easy to appreciate the meanings and contexts of the individual words and communities. (a) Here we pick a link and follow how the link merges with others as we climb the hierarchical tree. (b) We start from the link MARS–SATURN on the left, and the link SCUBA–DIVING on the right. As we move towards the root of the hierarchical tree, the link MARS–SATURN forms a ‘planet’ community, an ‘astronomy’ community, and then a more general ‘astronomy’ community. The link SCUBA–DIVING results in richer hierarchical structure: the link’s community becomes more and more general until we reach a large community of water-related words. (c) Here we delve into the hierarchical structure from a high level community into its sub-communities at a lower level. (d) We pick a sub-community from the example in (b) at threshold 0.20. We then identify its sub-communities at threshold 0.28. These sub-communities are represented by links with different colors. The sub-communities split into meaningful groups of similar words. Note that many links are not shown here because we are only drawing the link communities from these branches.

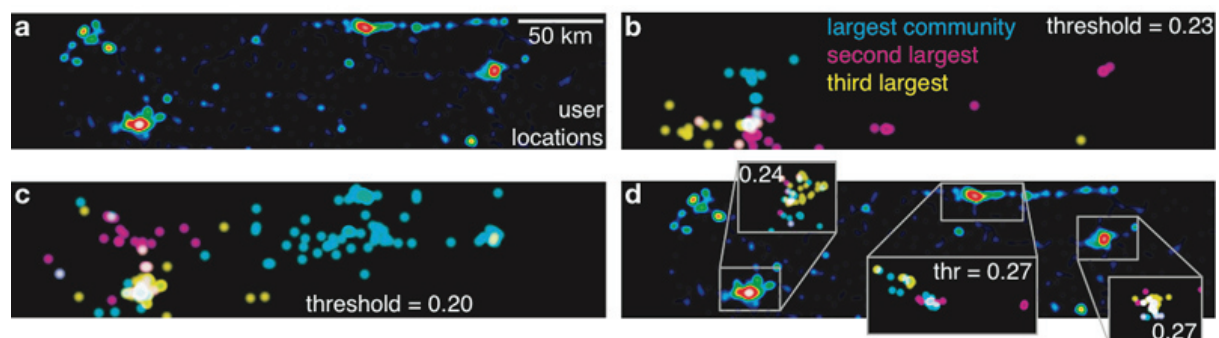


Figure 24: A spatial hierarchy of link communities amongst mobile phone users. (a) A heatmap showing the most likely geographic locations of all users in the network, several cities are present. (b) The three largest communities at the link dendrogram threshold with maximum partition density cluster around a single city. (c) Cutting the dendrogram at a lower threshold reveals regional but still spatially correlated communities. (d) At thresholds above that shown in (b) we see smaller, intra-city communities. Compare with main text Fig. 4.

- (i). **Structural changes across the dendrogram.** We show that dendrogram structure is ‘dynamic’ in the sense that when we cut the dendrogram at different thresholds, the community structure changes significantly. This means that there is not one optimal structure frozen into the dendrogram across a wide range of thresholds.
- (ii). **Meaningful communities.** We have already established the partition density D as a measure of the *structural* quality of a given partition of the dendrogram. At the optimal value of D , our algorithm finds high quality communities (see main text Fig. 2). As discussed in Sec. 3.3, the partition density D may take on a variety of shapes as a function of the dendrogram cut. The fact that D is sharply peaked does not necessarily imply that multiple, meaningful levels of community structure do not exist. This is both because a large amount of very different structure may be captured in a very narrow band of the dendrogram and because the partition density is an averaged quantity such that there may be many high quality communities alongside less dense groups.

While structural quality is important—in particular to community detection algorithms—the network structure *a priori* does not reveal information about how ‘meaningful’ the structure is. In order to quantitatively show that structures at multiple scales are ‘meaningful’ we use metadata to study community quality (see Sec. 5) as a function of the link dendrogram cut threshold.

The remainder of Sec. 7 is devoted to exploring these two aspects in further detail.

7.3 Dynamic dendrogram structure

To begin, we now explore the rate of change of the overlapping community structures encoded in the link dendrograms. One possible concern is that the number of mergers could potentially drop over a range of the dendrogram, resulting in large gaps where the structure is fixed (e.g., Fig. 7d). In this section, we present evidence that the dendrogram structures for networks in our test corpus are indeed dynamic over a large range of thresholds.

7.3.1 Branching probability

One straightforward way to illustrate the dynamic nature of the link dendrogram is to compute the *branching probability*, the fraction of communities at some threshold t that subsequently split into multiple communities slightly farther down the dendrogram, at threshold $t + \Delta t$. Low branching probability means that few communities are changing in that level of the dendrogram; conversely, the dendrogram’s structure is rapidly changing when the branching probability is high. As shown in Fig. 25, all networks in our test corpus possess significant and steady branching probabilities over a wide range of thresholds.

Here we use $\Delta t = 0.06$, but we have tested the dependence of the branching probability on Δt in Fig. 26 and find high probabilities over a wide range of values.

7.3.2 Distributions of community sizes and node memberships

In addition to the branching probabilities, we also examine the distribution of community sizes (nodes per community) and memberships (communities per node) at multiple cuts of the link dendrogram. These distributions tell us the scales of the detected communities for each threshold, and how those communities overlap.

In Fig. 27, we show these distributions at three different levels of each network’s link dendrogram. We observe that many networks possess broad distributions of community sizes, indicating that a variety of size scales are encoded at each level of the dendrogram. The broad membership distributions simultaneously indicate that the amount of overlap remains significant at those same levels. These results mean that the structures encoded in the link dendrogram do not suddenly collapse but vary smoothly

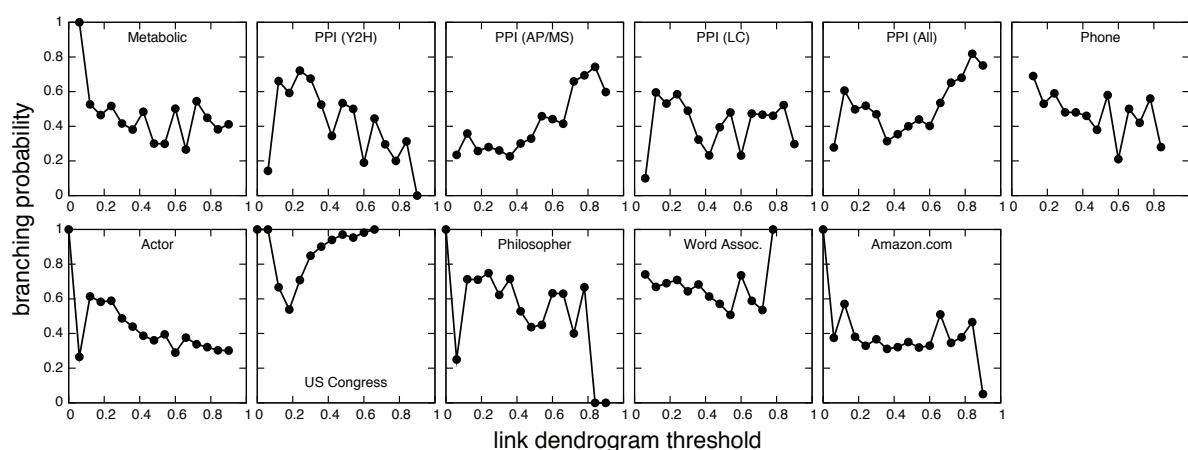


Figure 25: Branching probabilities for the link dendrograms of the networks studied in our test corpus. In all networks, the branching probability is high over a large range of thresholds, indicating that the structures encoded by the dendrograms are constantly changing.

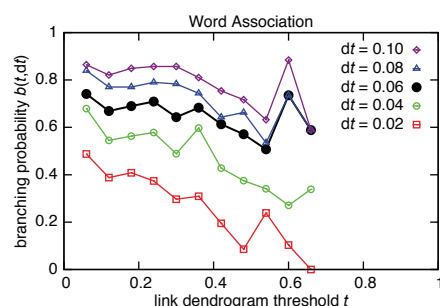


Figure 26: Studying the dependence of the branching probability $b(t, \Delta t)$ on the threshold window Δt . Since $b \rightarrow 0$ when $\Delta t \rightarrow 0$ and $b \rightarrow 1$ when $\Delta t \rightarrow 1$, we must demonstrate that there is a range where Δt is small but b is still large. To do so, we plot b versus t for several small values of Δt . We see that even for the lowest value, b is substantial for a wide range of t . Here we show the word association network, but this fact is generic over the test corpus. (We start the curves at $t = 0.7$ because the dense word association network does not begin clustering until $t \approx 0.8$, see Fig. 29.)

as a function of dendrogram threshold. We also observe that in some networks the community scales change while the amount of overlap remains steady (particularly the phone network), whereas in other networks the distributions of sizes vary less but the amount of overlap changes drastically (particularly the metabolic and PPI (all) networks). In conjunction with the branching probability, these properties highlight how the link dendrogram can reveal multiple aspects of the network's levels of hierarchical community structure.

7.4 Revealing meaningful communities at multiple scales

Now that we have shown that very different scale structures are contained throughout the link dendrograms, we must also demonstrate that these structures are meaningful.

7.4.1 Community quality as a function of cut-level

As we move from the leaves of the dendrogram (where each link is isolated) towards the root (where all links are merged into a giant community) communities must grow in size. Due to the construction of the community quality measures (see Sec. 5 for details about specific types of metadata), the community quality is likely to drop whenever two communities are joined—since a larger community is likely to be more diverse. For example, while ‘Physics’ and ‘Chemistry’ may be subsumed under the heading ‘Natural Science’, each field on its own is more homogeneous than the merger of the two.

Thus, it is likely that, relative to the optimal communities, the community quality will decay as the dendrogram cut approaches the root of the dendrogram. For this reason, meaningful communities are

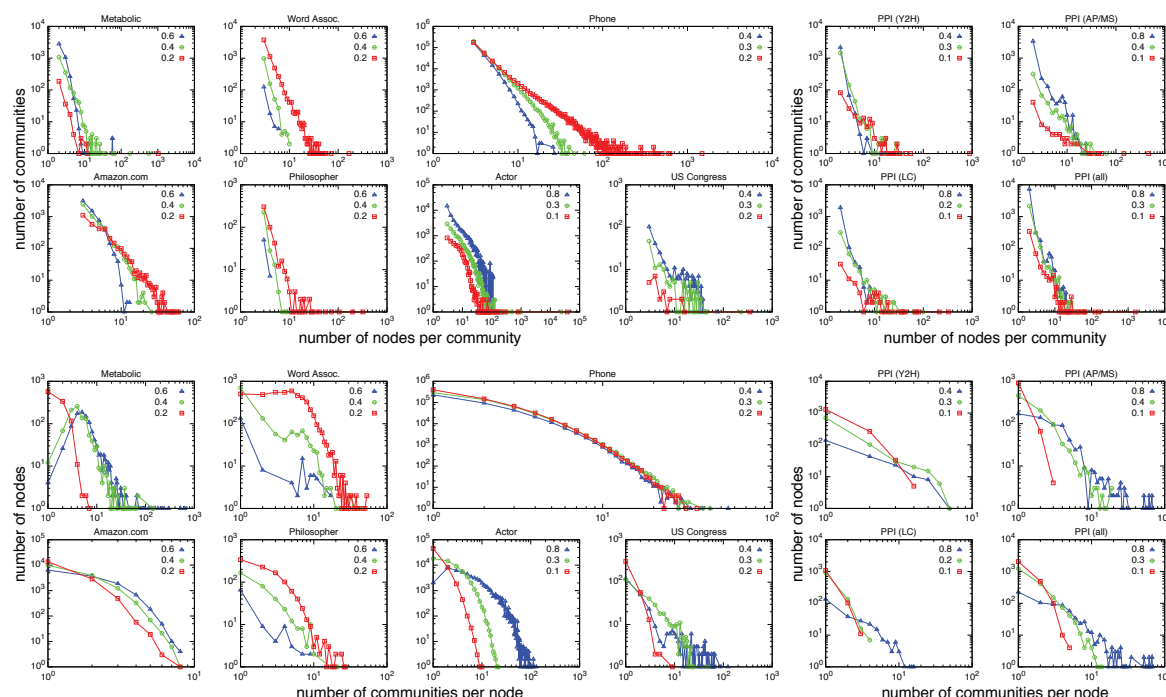


Figure 27: Overlapping community structure is very different when cutting the link dendrograms at different thresholds. Shown are the distributions of community sizes and memberships for the networks in our test corpus, each at three different link dendrogram thresholds. A broad, heavy-tailed distribution of community sizes arises at high thresholds in most networks and then persists over a wide range of the link dendrogram, indicating that the link dendrogram does not suddenly collapse but changes smoothly over much of its range. Meanwhile, the distributions of community memberships per node remain broad over the same region of the dendrogram (this effect is particularly striking in the phone network), indicating that overlapping structure is maintained throughout the dendrograms in nearly all networks. These results show that the community structures contained in the link dendrograms cover a wide range of scales while maintaining significant overlap.

expressed as a *slow* decay of community quality, compared to a properly randomized control dendrogram. We now show that all link dendrograms for our test corpus exhibit such slow decay, compared with the following control.

Randomized control dendrogram We wish to test whether the hierarchical structure is valid beyond some threshold t_* , e.g., that with maximum partition density. To do this, we introduce the following control: first, compute the similarities $S(e_{ik}, e_{jk})$ for all connected edge-pairs (e_{ik}, e_{jk}) , as normal. Then perform our standard single-linkage hierarchical clustering, merging all edge-pairs in descending order of S while $S \geq t_*$, fixing the community structure at $t = t_*$.

Below t_* , randomly shuffle similarities amongst the remaining edge-pairs with $S < t_*$, then proceed with the merging process as before. This randomization only alters merging order, and ensures that the *rate* of edge-pair merging is preserved, since the same similarities are clustered. This strictly controls not only the merging rate, but also the similarity distributions and the high-quality community structure found at t_* . This procedure ensures that the dendrogram is properly randomized while other salient features are conserved. See Fig. 28.

If there is significant, meaningful structure for $t < t_*$, we expect the actual community's quality Q to decay slower than the randomized control quality Q_{rand} . As shown in Fig. 29, this is the exact behavior we find across the entire network corpus¹⁰.

¹⁰Notice in the Actor network we see that the very large link communities appear *worse* than the control. The IMDb data is known to strongly split at very large scales, according to language groups [58]. Since our quality measure is based on plot keywords and not languages, the dendrogram may capture the true, large scale structure but this is not reflected in the metadata.

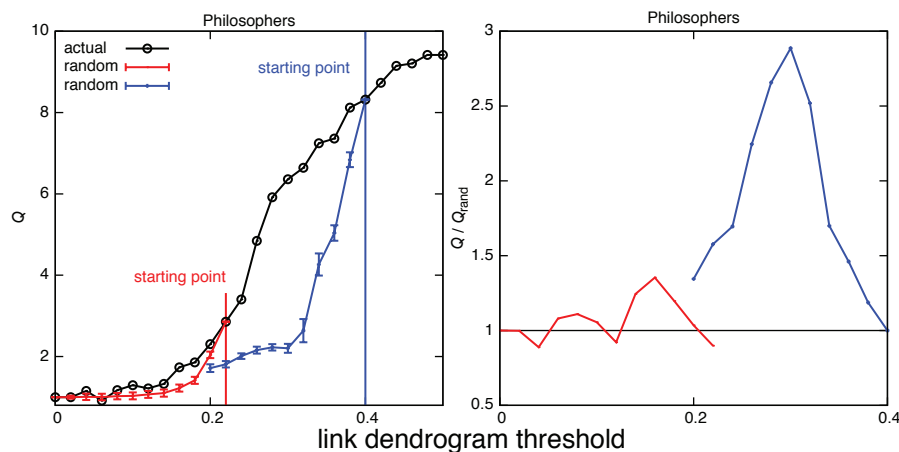


Figure 28: An illustration of the link dendrogram control, using the Philosopher network. We wish to test whether the hierarchical structure is valid beyond some threshold t_* . To do so, we first compute the edge-pair similarities of all “cluster-able” edges. We then cluster edges according to their similarity (as normal) until we have reached t_* . Afterwards, we then cluster the remaining edge-pairs at random. This control is much stronger than, e.g., clustering random pairs of edges, since the exact same edge-pairs are being clustered together, only the ordering of the clustering is changed. If there is significant, meaningful structure for $t < t_*$, we should expect the actual community’s quality Q to decay slower than the control’s quality Q_{rand} . In this example, we choose two values of t_* (vertical lines) and show that the philosopher network’s communities possess significant structure beyond $t_* = 0.4$, but little structure beyond 0.22.

7.4.2 Hierarchical metadata

Finally, the Amazon.com and PPI networks in our test corpus possess multi-level metadata. For these networks, we can construct a direct test of whether there are meaningful communities at different levels of the link dendrogram. For instance, a book in the Amazon.com network has category information at multiple levels of granularity, see Fig. 30 (top) for an example. The PPI networks also contain hierarchical information: GO terms (see Sec. 6.2.1) are organized hierarchically, forming a directed acyclic graph; the MIPS functional catalog also provides a hierarchical categorization of each protein.

From these hierarchical metadata, we now extract two sets of metadata: *coarse* and *fine*. If our method is able to find meaningful structures at multiple scales, we expect that the community quality based on the fine metadata will have high values at cuts near the leaves of the dendrogram, and the community quality based on coarse metadata will have high values for lower thresholds (higher than those using the fine metadata). That is, coarse-grained communities at the lower threshold will conform well with the coarse metadata while detailed, fine communities at higher thresholds will conform well with the fine metadata, as illustrated in Fig. 30.

For the Amazon.com network, we use the available subject categories given for each book, stored as lists, each of which are ordered by level of granularity (one list for THE BOOK THIEF is shown at the top of Fig. 30). Broad categories such as ‘General’ are removed. The coarse metadata for each book is then the set of first elements of that book’s category lists, and the fine metadata are the last elements.

In the PPI network, we use the MIPS functional catalog annotations since they provide a clearly defined set of hierarchical metadata: For instance, *metabolism* is labeled ‘01’, *amino acid metabolism* is ‘01.01’, *assimilation of ammonia* is ‘01.01.03’, and so on. Each level is separated by a period, and each level is represented by two digits. The coarse metadata is obtained by reducing every annotation to its first hierarchical level. For instance, if a protein has an annotation ‘01.01.03’, we can represent it by ‘01’. These metadata constitute the coarse metadata for the protein. The fine metadata is obtained by removing all metadata that have two or less levels of information, and reducing longer metadata to three levels. For example, ‘01.01’ or ‘01’ will be removed from the annotation, and ‘01.01.01.01.01’ becomes ‘01.01.01’. We choose the third level as the fine metadata because there are only a few proteins

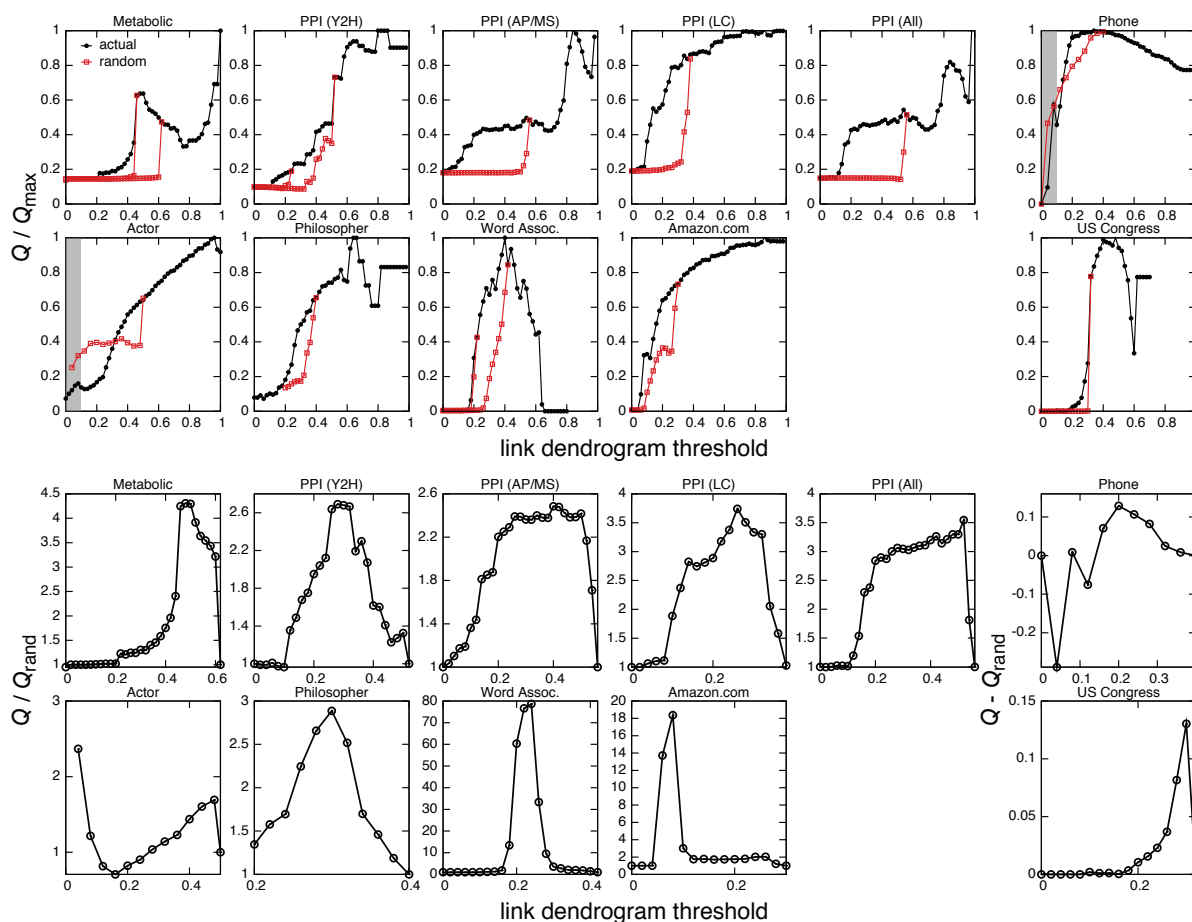


Figure 29: **(top)** The community quality Q (see Sec. 5.2) as a function of dendrogram threshold for the corpus networks. We see that most networks possess very slow decay of quality across a wide range of the dendrogram. This is particularly true for PPI (AP/MS), PPI (LC), PPI (All), Phone, word association, and Amazon.com networks. The control, shown in red, indicates that all networks possess meaningful hierarchical structure beyond the examined threshold. (For metabolic, PPI (Y2H), and the word association networks, we test multiple thresholds.) Notice in the Actor network we see that the very large link communities appear *worse* than the control. The IMDb data is known to strongly split at very large scales, according to language groups [58]. Since our quality measure is based on plot keywords and not languages, the dendrogram may capture the true, large scale structure but this is not reflected in the metadata. We plot Q/Q_{\max} , normalizing the enrichments (dispersions in the case of the Phone and US Congress networks) by their maximal value. For the large Phone and Actor networks, we sample communities to speed up the calculation of the quality of the null partitions. This may introduce a small positive bias in the shaded regions. **(bottom)** The relative quality Q/Q_{rand} (the ratio of the two curves), highlighting the validity of each link dendrogram's hierarchy. For the Phone and US Congress networks we instead plot $Q - Q_{\text{rand}}$ as the difference is more meaningful than the ratio for dispersive measures.

that have finer levels of annotations, and thus these finer levels are too noisy.

With these two sets of metadata, we calculate community quality and coverage for the different networks. Figures 31 and 32 clearly show the difference between coarse and fine metadata. In every case, the coarse metadata remains relatively more important at lower thresholds (near the root of the dendrogram) and the fine metadata becomes less important. This confirms our hypothesis shown in Fig. 30 and indicates that the structures throughout the link dendrogram correspond well to the hierarchical metadata.

Finally, it is interesting to note that the highly clustered AP/MS network shows a distinct pattern in the link dendrogram compared to the LC network. By calculating 'normalized performance,' the normalized sum of community quality Q/Q_{\max} and coverage, we see that the dense AP/MS protein co-complex clusters give that network a clear optimum at higher thresholds (~ 0.6) than the LC network,

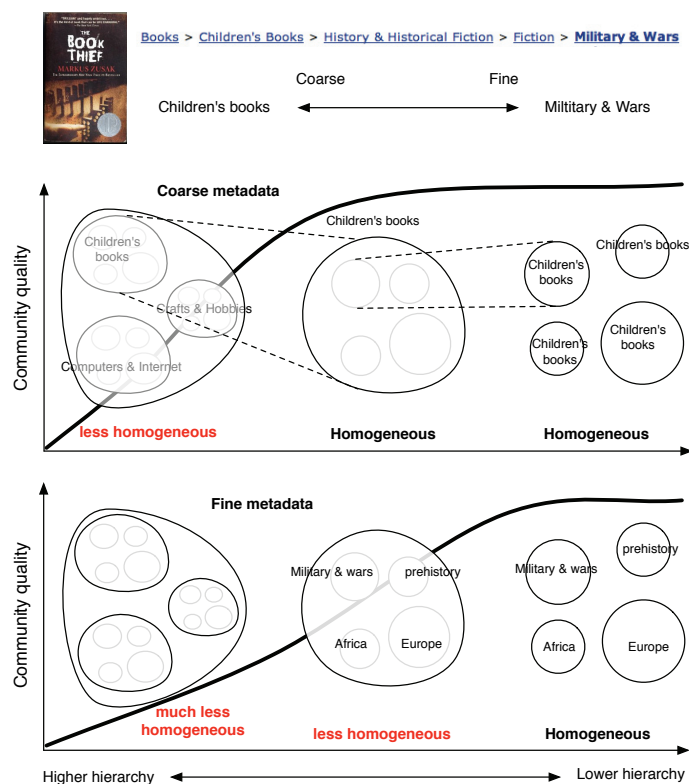


Figure 30: A cartoon explaining multi-scale metadata. **(top)** Nodes in some of our test networks have metadata that are organized hierarchically. We can use these data to study the hierarchical organization of the communities we detect. This schematic figure illustrates the case where community structure at multiple levels is successfully revealed. **(middle)** If we use coarse metadata to evaluate the community quality, it will remain high until we reach the point where the scale of communities is larger than the scale described by the coarse metadata. **(bottom)** Meanwhile, if we use fine metadata, the quality will remain high until the point where the scale of communities is larger than the scale described by the fine metadata. That is, a clear distinction between the two curves of community quality versus threshold will emerge: one with coarse metadata and the other with fine metadata. The difference will vanish if one fails to capture the hierarchical structure between the two scales that are described by coarse and fine metadata. See Figs. 32 and 31 for results.

which peaks at ~ 0.2 . Meanwhile, the PPI (all) network, which contains all other PPI networks, shows *two* distinct peaks in performance, one corresponding to the AP/MS structure and one corresponding to LC. Thus the link dendrogram for the PPI (all) network captures AP/MS-specific structure at one level and LC-specific structure at another. The sparse Y2H network does not exhibit as much community structure as LC and AP/MS, and thus has little impact on the community structures of PPI (all), compared with the other constituent networks.

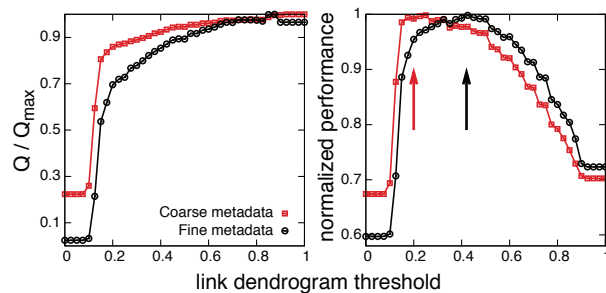


Figure 31: Hierarchically organized product category meta-data for the Amazon.com network confirms the validity of the discovered link dendrogram. **(left)** Community quality remains high for the coarse metadata for longer than the fine metadata, although both decay quite slowly. Note that controlling for the global baseline enrichment by normalizing with $(Q - Q_{\min}) / (Q_{\max} - Q_{\min})$ does not change this effect. **(right)** Normalized performance, the normalized sum of community quality and coverage, reveals that the fine metadata peaks earlier (threshold ~ 0.4) than the coarse metadata (threshold ~ 0.2), indicating that the community partitions at multiple levels of the link dendrogram are meaningful according to the hierarchical metadata.

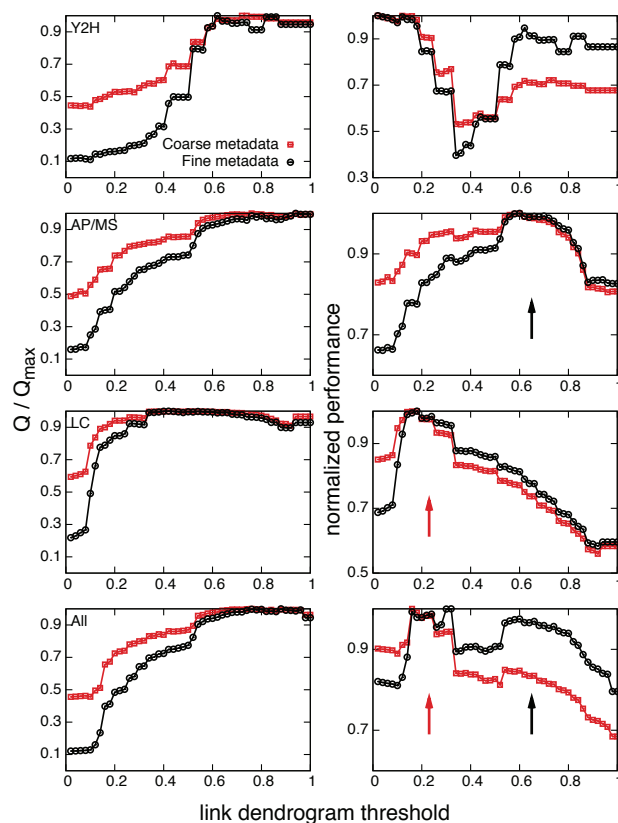


Figure 32: Hierarchical metadata confirms that distinct structures are visible throughout the link dendrograms of the PPI networks. Here we compute community quality (left column) and normalized performance, the normalized sum of quality and coverage (right column) for all four networks. As with the Amazon.com network, the quality decays more rapidly for the fine metadata than for the coarse (see Fig. 30), indicating that each link dendrogram's structures correspond well with the networks' existing metadata. Using normalized performance, the highly clustered AP/MS network shows a distinct pattern in the link dendrogram compared to the LC network. The dense AP/MS protein co-complex clusters give that network a clear optimum at higher thresholds (~ 0.6 , black arrow) than the LC network, which peaks at ~ 0.2 (red arrow). The PPI (all) network, which contains AP/MS and LC, shows *two* distinct peaks in performance, one corresponding to the AP/MS structure and one corresponding to LC. Thus the link dendrogram for the PPI (all) network captures AP/MS-specific structure at one level and LC-specific structure at another.

A Tables of measures

Here we list the raw (unnormalized) values for the four calculated measures, the networks and the algorithms that were shown in main text Fig. 2 and Fig. 19. For clique percolation we have chosen the value of k that gives the best overall composite score (see Appendix A.2), unless there is an existing precedent in the literature. Note that this weighs coverage and quality equally, but an experimenter may wish to prioritize coverage for quality, or vice versa.

A.1 Overall methods

```
#1### metabolic #####
# comm. quality comm. coverage overlap quality overlap coverage
L 4.77233482 0.95009597 0.41809907 4.65642994
N 3.79668962 0.63339731 0.03529230 0.63339731
C 1.14228117 0.66890595 0.16156926 0.87523992
G 1.99859878 0.99808061 0.02294501 0.99808061
I 3.69066921 1.00000000 0.02419719 1.00000000
# overall winner: L

#2### PPI (Y2H) #####
# comm. quality comm. coverage overlap quality overlap coverage
L 2.35223830 0.55555556 0.08653618 0.72374013
N 1.92594816 0.73102611 0.02482696 0.73102611
C 1.87169405 0.16393443 0.05782404 0.18397086
G 1.36153975 0.99149970 0.01060934 0.99149970
I 2.25216779 0.98785671 0.01590075 0.98785671
# overall winner: L

#3### PPI (AP/MS) #####
# comm. quality comm. coverage overlap quality overlap coverage
L 2.73145231 0.83864542 0.38406704 2.57669323
N 2.19482167 0.91135458 0.03996271 0.91135458
C 2.07359450 0.76792829 0.13443620 0.81673307
G 1.94785560 0.99103586 0.01632051 0.99103586
I 2.75864056 0.99203187 0.01480369 0.99203187
# overall winner: L

#4### PPI (LC) #####
# comm. quality comm. coverage overlap quality overlap coverage
L 4.52990197 0.55812036 0.17366791 0.93075021
N 2.34665560 0.96537510 0.02541182 0.96537510
C 2.91313090 0.55647156 0.11309138 0.60428689
G 2.58449740 0.99175598 0.01737294 0.99175598
I 3.76173052 0.99175598 0.01857447 0.99175598
# overall winner: L

#5### PPI (All) #####
# comm. quality comm. coverage overlap quality overlap coverage
L 3.51593751 0.41260535 0.19629188 1.29754489
N 2.78187442 0.76511543 0.01589616 0.76511543
C 1.07221941 0.52876512 0.10141995 0.57053866
G 1.36531606 0.99523635 0.00797972 0.99523635
I 3.35047694 0.99340418 0.01124843 0.99340418
# overall winner: L

#6### phone #####
# comm. quality comm. coverage overlap quality overlap coverage
L 0.75761102 0.76180404 0.13760916 1.42556059
N 0.33284757 0.78113498 0.01301070 0.78113498
C 0.82114799 0.33514186 0.07811141 1.27819838
G -0.17040443 0.99970880 0.00029690 0.99970880
I 0.61369550 0.99967268 0.00031225 0.99967268
# overall winner: L

#7### actor #####
# comm. quality comm. coverage overlap quality overlap coverage
L 6.65974811 0.57986085 0.04076675 1.51764549
N 1.87424645 0.83947724 0.00706428 0.83947724
C 2.03239313 0.69482725 0.01468963 0.79485544
G 1.56548814 1.00000000 0.00000000 1.00000000
I 2.01709951 0.99273116 0.00106879 0.99273116
# overall winner: L

#8### congress #####
# comm. quality comm. coverage overlap quality overlap coverage
L 0.34647780 0.94358974 0.68222751 5.89743590
N 0.38692427 0.61794872 0.03855387 0.61794872
C 0.26286049 0.61282051 0.15720036 0.77435897
G 0.42350813 1.00000000 0.00000000 1.00000000
I 0.32595601 0.99487179 0.00000000 0.99487179
# overall winner: L

#9### philosopher #####
# comm. quality comm. coverage overlap quality overlap coverage
L 2.40739272 0.81788351 0.45773225 2.66119770
N 1.68405991 0.66037736 0.06243616 0.66037736
C 1.18575668 0.74405250 0.12858942 0.77276456
```

```

G 1.47736530 0.99835931 0.00791000 0.99835931
I 2.20936130 0.99015587 0.01235264 0.99015587
# overall winner: L

#10### word assoc. #####
# comm. quality comm. coverage overlap quality overlap coverage
L 83.16274063 0.92447190 0.09459306 5.23455560
N 5.31083477 0.56954962 0.02424692 0.56954962
C 33.94752060 0.62554803 0.06495803 1.05579912
G 1.69216772 0.99820646 0.00275916 0.99820646
I 12.98083645 1.00000000 -0.00000000 1.00000000
# overall winner: L

#11### amazon #####
# comm. quality comm. coverage overlap quality overlap coverage
L 102.81247272 0.90629479 0.01281968 1.22103406
N 6.71393780 0.95022599 0.00296223 0.95022599
C 89.11665793 0.88562452 0.01051039 1.03836402
G 8.95745118 1.00000000 0.00000000 1.00000000
I 75.04521188 1.00000000 0.00000000 1.00000000
# overall winner: L

```

A.2 Clique Percolation

When applying clique percolation we picked the value of clique size k that gave the best overall (normalized) composite score. Here we list the raw values for multiple k (shown as cp3, cp4, etc.). The overall winner lists the chosen value of k used in the main text and in Appendix A. If there is an existing precedent for which value of k to use, such as with the mobile phone data [16], we follow the original work.

It is important to note that choosing the k to maximize the composite performance score weighs coverage and quality equally, whereas a researcher may wish to sacrifice coverage for quality. Higher values of k tend to find very high quality communities; it is up to the researcher's discretion if such a choice is appropriate to his or her particular application.

```

#1### metabolic #####
# cp3 cp4 cp5 cp6 cp7 cp8 cp9
comm. quality 1.05405749 1.08502531 1.14092849 1.14228117 1.24034244 1.31172522 1.29127564
comm. coverage 0.99328215 0.97696737 0.88291747 0.66890595 0.46065259 0.32053743 0.17850288
over. quality 0.02817960 0.05334991 0.10411422 0.16156926 0.14877168 0.17817169 0.15137653
over. coverage 0.99328215 1.01919386 0.97312860 0.87523992 0.54798464 0.44913628 0.19673704
# overall winner: cp6

#2### PPI (Y2H) #####
# cp3 cp4 cp5
comm. quality 1.87169405 8.85655602 0.00000000
comm. coverage 0.16393443 0.01700061 0.00000000
over. quality 0.05782404 0.01462019 0.00000000
over. coverage 0.18397086 0.01700061 0.00000000
# overall winner: cp3

#3### PPI (AP/MS) #####
# cp3 cp4 cp5
comm. quality 1.84479625 2.07359450 1.98196162
comm. coverage 0.87250996 0.76792829 0.67430279
over. quality 0.08580265 0.13443620 0.11253619
over. coverage 0.91235060 0.81673307 0.70019920
# overall winner: cp4

#4### PPI (LC) #####
# cp3 cp4 cp5
comm. quality 2.91313090 3.59607890 3.98440063
comm. coverage 0.55647156 0.30502885 0.19043693
over. quality 0.11309138 0.09132684 0.06699107
over. coverage 0.60428689 0.32976092 0.20857378
# overall winner: cp3

#5### PPI (All) #####
# cp3 cp4 cp5
comm. quality 1.07221941 2.43214030 2.36350203
comm. coverage 0.52876512 0.35507512 0.28325394
over. quality 0.10141995 0.09320351 0.07553584
over. coverage 0.57053866 0.38402345 0.29754489
# overall winner: cp3

#6### phone #####
# cp4
comm. quality 0.82114799
comm. coverage 0.33514186
over. quality 0.07811141
over. coverage 1.27819838
# overall winner: cp4

#7### actor #####
# cp3 cp4 cp5

```

```

comm. quality    2.03239313    2.16441181    2.65452628
comm. coverage   0.69482725    0.49747074    0.36856003
over. quality    0.01468963    0.01664214    0.01300756
over. coverage   0.79485544    0.60769014    0.45345715
# overall winner: cp3

#8### congress #####
#
comm. quality    cp3      cp4      cp5      cp6      cp7      cp8      cp9      cp10     cp11     cp12
comm. quality    0.00062736 0.00447642 0.00983444 0.00973442 0.02589494 0.05050994 0.26286049 0.29205756 0.31483181 0.31804022
comm. coverage   0.94358974 0.88717949 0.83846154 0.78461538 0.71794872 0.66410256 0.61282051 0.55384615 0.48717949 0.40512821
over. quality    0.04074278 0.02453166 0.03971995 0.07602430 0.07531725 0.06842514 0.15720036 0.04223622 0.05535796 0.0
over. coverage   0.94358974 0.88717949 0.84871795 0.81538462 0.75641026 0.68717949 0.77435897 1.11574074 1.04736842 1.0
# overall winner: cp9

#9### philosopher #####
#
comm. quality    cp3      cp4      cp5      cp6      cp7
comm. quality    1.18575668 1.42170714 1.76620515 3.78843554 4.61831912
comm. coverage   0.74405250 0.49056604 0.26579163 0.11812961 0.05004102
over. quality    0.12858942 0.19299861 0.20831356 0.16526745 0.11187828
over. coverage   0.77276456 0.55865463 0.34536505 0.18375718 0.07957342
# overall winner: cp3

#10### word assoc. #####
#
comm. quality    cp3      cp4      cp5      cp6      cp7      cp8      cp9
comm. quality    1.00181181 1.16419232 33.94752060 61.96886046 63.47129301 101.26453476 47.66309121
comm. coverage   0.99860502 0.93941809 0.62554803 0.25308888 0.05759267 0.01215624 0.00378637
over. quality    0.00787718 0.06339547 0.06495803 0.03193241 0.01055262 0.00234616 0.00188072
over. coverage   1.03786369 1.41072140 1.05579912 0.34436030 0.06695895 0.01335193 0.00378637
# overall winner: cp4 (note: we use k=5 because the k=4 comm. quality was too low. These results are unfiltered)

#11### amazon #####
#
comm. quality    cp3      cp4      cp5      cp6
comm. quality    89.11665793 123.14041107 132.90590155 138.69567284
comm. coverage   0.88562452 0.60577665 0.30729798 0.07871238
over. quality    0.01051039 0.01587309 0.01210945 0.00709472
over. coverage   1.03836402 0.66563775 0.32526734 0.08020064
# overall winner: cp3

```

References

- [1] Jaccard, P. étude comparative de la distribution florale dans une portion des alpes et des jura. *Bulletin del la Société Vaudoise des Sciences Naturelles* **37**, 547–579 (1901).
- [2] Newman, M. E. J. & Girvan, M. Finding and evaluating community structure in networks. *Physical Review E* **69**, 026113 (2004).
- [3] Danon, L., Diaz-Guilera, A., Duch, J. & Arenas, A. Comparing community structure identification. *Journal of Statistical Mechanics: Theory and Experiment* **2005**, P09008 (2005).
- [4] Shen, H., Cheng, X., Cai, K. & Hu, M.-B. Detect overlapping and hierarchical community structure in networks. *Physica A* **388**, 1706–1712 (2009).
- [5] Nicosia, V., Mangioni, G., Carchiolo, V. & Malgeri, M. Extending the definition of modularity to directed graphs with overlapping communities. *J Stat Mech-Theory E* P03024 (2009).
- [6] Reichardt, J. & Bornholdt, S. Detecting fuzzy community structures in complex networks with a potts model. *Phys. Rev. Lett.* **93**, 218701 (2004).
- [7] Li, D. *et al.* Synchronization interfaces and overlapping communities in complex networks. *Phys. Rev. Lett.* **101**, 168701 (2008).
- [8] Lancichinetti, A., Fortunato, S. & Kertesz, J. Detecting the overlapping and hierarchical community structure in complex networks. *New Journal of Physics* **11**, 033015 (2009).
- [9] Knuth, D. E. *The Stanford GraphBase: A Platform for Combinatorial Computing* (Addison-Wesley, Reading, MA, 1993).
- [10] Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. N. & Barabási, A.-L. Hierarchical organization of modularity in metabolic networks. *Science* **297**, 1551–1555 (2002).
- [11] Palla, G., Derény, I., Farkas, I. & Vicsek, T. Uncovering the overlapping community structure of complex networks in nature and society. *Nature* **435**, 814 (2005).
- [12] Rosvall, M. & Bergstrom, C. T. Maps of random walks on complex networks reveal community structure. *Proceedings of the National Academy of Sciences* **105**, 1118–1123 (2008).
- [13] Newman, M. E. J. Fast algorithm for detecting community structure in networks. *Physical Review E* **69**, 066133 (2004).

- [14] Clauset, A., Newman, M. E. J. & Moore, C. Finding community structure in very large networks. *Phys. Rev. E* **70**, 066111 (2004).
- [15] Palla, G., Farkas, I. J., Pollner, P., Derenyi, I. & Vicsek, T. Directed network modules. *New Journal of Physics* **9**, 186 (2007).
- [16] Palla, G., Barabási, A.-L. & Vicsek, T. Quantifying social group evolution. *Nature* **446**, 664 (2007).
- [17] Kumpula, J. M., Kivelä, M., Kaski, K. & Saramäki, J. Sequential algorithm for fast clique percolation. *Phys. Rev. E* **78**, 026109 (2008).
- [18] Newman, M. E. J. Detecting community structure in networks. *The European Physical Journal B* **38**, 321–330 (2004).
- [19] Girvan, M. & Newman, M. E. J. Community structure in social and biological networks. *Proceedings of the National Academy of Sciences* **99**, 7821–7826 (2002).
- [20] Lancichinetti, A. & Fortunato, S. Community detection algorithms: a comparative analysis. *Phys. Rev. E* **80**, 056117 (2009).
- [21] Newman, M. E. J. & Park, J. Why social networks are different from other types of networks. *Physical Review E* **68**, 036122 (2003).
- [22] Clauset, A., Moore, C. & Newman, M. E. J. Hierarchical structure and the prediction of missing links in networks. *Nature* **453**, 98 (2008).
- [23] Sales-Pardo, M., Guimera, R., Moreira, A. & Amaral, L. Extracting the hierarchical organization of complex systems. *PNAS* **104**, 15224–15229 (2007).
- [24] Martinez, N. D., Hawkins, B. A. & adn B. P. Feifarek, H. A. D. *Ecology* **80**, 1044–1055 (1999).
- [25] Derényi, I., Palla, G. & Vicsek, T. Clique percolation in random networks. *Phys. Rev. Lett.* **94**, 160202 (2005).
- [26] Doyon, Y., Selleck, W., Lane, W. S., Tan, S. & Côté, J. Structural and functional conservation of the nua4 histone acetyltransferase complex from yeast to humans. *Mol. Cell. Biol.* **24**, 1884 (2004).
- [27] Dotson, M. R. *et al.* Structural organization of yeast and mammalian mediator complexes. *Proceedings of the National Academy of Sciences of the United States of America* **97**, 14307–14310 (2000).
- [28] Wu, P.-Y. J., Ruhlmann, C., Winston, F. & Schultz, P. Molecular architecture of the s. cerevisiae saga complex. *Molecular Cell* **15**, 199 – 208 (2004).
- [29] Saleh, A. *et al.* Tra1p Is a Component of the Yeast AdaSpt Transcriptional Regulatory Complexes. *J. Biol. Chem.* **273**, 26559–26565 (1998).
- [30] Brown, C. E. *et al.* Recruitment of HAT Complexes by Direct Activator Interactions with the ATM-Related Tra1 Subunit. *Science* **292**, 2333–2337 (2001).
- [31] Bhaumik, S. R., Raha, T., Aiello, D. P. & Green, M. R. In vivo target of a transcriptional activator revealed by fluorescence resonance energy transfer. *Genes & Development* **18**, 333–343 (2004).
- [32] Baumeister, W., Walz, J., Zühl, F. & Seemüller, E. The proteasome: Paradigm of a self-compartmentalizing protease. *Cell* **92**, 367 – 380 (1998).
- [33] 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).
- [34] Tanimoto, T. T. An elementary mathematical theory of classification and prediction. Tech. Rep., IBM Internal Report (1958).
- [35] Bagrow, J. P. & Boltt, E. M. A local method for detecting communities. *Phys. Rev. E* **72**, 046108 (2005).
- [36] Bagrow, J. P. Evaluating local community methods in networks. *J. Stat. Mech.* **2008**, P05001 (2008).
- [37] Clauset, A. Finding local community structure in networks. *Physical Review E* **72**, 026132 (2005).
- [38] Newman, M. E. J. & Leicht, E. A. Mixture models and exploratory analysis in networks. *Proceed-*

- ings of the *National Academy of Sciences* **104**, 9564–9569 (2007).
- [39] Guimerà, R. & Amaral, L. A. N. Functional cartography of complex metabolic networks. *Nature* **433**, 895–900 (2005).
- [40] Fortunato, S. & Castellano, C. *Community Structure in Graphs* (Springer, 2009).
- [41] Watts, D. J. & Strogatz, S. H. Collective dynamics of ‘small-world’ networks. *Nature* **393**, 440 (1998).
- [42] Yu, H. *et al.* High-Quality Binary Protein Interaction Map of the Yeast Interactome Network. *Science* **322**, 104–110 (2008).
- [43] Feist, A. M. *et al.* A genome-scale metabolic reconstruction for *Escherichia coli* K-12 MG1655 that accounts for 1260 orfs and thermodynamic information. *Molecular Systems Biology* **3**, 1 (2007).
- [44] Onnela, J.-P. *et al.* Structure and tie strengths in mobile communication networks. *PNAS* **104**, 7332 (2007).
- [45] Onnela, J.-P. *et al.* Analysis of a large-scale weighted network of one-to-one human communication. *New Journal of Physics* **9**, 179 (2007).
- [46] Palla, G., Barabási, A. & Vicsek, T. Quantifying social group evolution. *Nature* **446**, 664–667 (2007).
- [47] Gonzalez, M. C., Hidalgo, C. A. & Barabási, A.-L. Understanding individual human mobility patterns. *Nature* **453**, 479 (2008).
- [48] IMDb. <http://www.imdb.com> (2009).
- [49] Fowler, J. Connecting the congress: A study of cosponsorship networks. *Political Analysis* **14**, 456–487 (2006).
- [50] Fowler, J. Legislative cosponsorship networks in the U.S. house and senate. *Social Networks* **28**, 454–465 (2006).
- [51] Poole, K. T. Recovering a basic space from a set of issue scales. *American Journal of Political Science* **42**, 954–993 (1998).
- [52] Poole, K. T. *Spatial Models of Parliamentary Voting* (Cambridge University Press, New York, 2005).
- [53] Fellbaum, C. *WordNet: An Electronical Lexical Database* (The MIT Press, Cambridge, MA, 1998).
- [54] Gene Ontology Consortium. *Nucleic Acids Res.* **36**, D440 (2008).
- [55] Yu, H., Jansen, R., Stolovitzky, G. & Gerstein, M. Total ancestry measure: quantifying the similarity in tree-like classification, with genomic applications. *Bioinformatics* **23**, 2163–2173 (2007).
- [56] Kanehisa, M. & Goto, S. *Nucleic Acids Res.* **28**, 27–30 (2000).
- [57] 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/>) (1998).
- [58] Leskovec, J., Lang, K. J., Dasgupta, A. & Mahoney, M. W. Community structure in large networks: Natural cluster sizes and the absence of large well-defined clusters. *CoRR* **abs/0810.1355** (2008).