

# An Algorithm to Find Overlapping Community Structure in Networks

Steve Gregory

Department of Computer Science  
University of Bristol, BS8 1UB, England  
steve@cs.bris.ac.uk

**Abstract.** Recent years have seen the development of many graph clustering algorithms, which can identify community structure in networks. The vast majority of these only find disjoint communities, but in many real-world networks communities overlap to some extent. We present a new algorithm for discovering overlapping communities in networks, by extending Girvan and Newman's well-known algorithm based on the *betweenness* centrality measure. Like the original algorithm, ours performs hierarchical clustering — partitioning a network into any desired number of clusters — but allows them to overlap. Experiments confirm good performance on randomly generated networks based on a known overlapping community structure, and interesting results have also been obtained on a range of real-world networks.

## 1 Introduction and Motivation

Many complex systems in the real world can be represented abstractly as networks (or graphs). Recently, with increasing availability of data about large networks and the need to understand them, the study of networks has become an important topic. A property that has been extensively studied is the existence of community structure in networks. A *cluster* (or *community* or *module*) is a subgraph such that the density of edges within it (*intracluster edges*) is greater than the density of edges between its vertices and those outside it (*intercluster edges*). A wide range of algorithms have been developed to discover communities in a network, including [4, 7, 11, 12, 13, 14].

Probably the best-known algorithm for finding community structure is Girvan and Newman's algorithm [7, 14], based on the betweenness centrality measure [6]. The *betweenness* (strictly, the *shortest-path betweenness*) of edge  $e$ ,  $c_B(e)$ , is defined as the number of shortest paths, between all pairs of vertices, that pass along  $e$ . A high betweenness means that the edge acts as a bottleneck between a large number of vertex pairs and suggests that the edge is an intercluster edge. Although the algorithm is quite slow and is no longer the most effective clustering algorithm, it does give relatively good results. The algorithm works as follows:

1. Calculate edge betweenness of all edges in network.
2. Find edge with highest edge betweenness and remove it.
3. Recalculate edge betweenness for all remaining edges.
4. Repeat from step 2 until no edges remain.

This is a hierarchical, divisive, clustering algorithm. Initially, the  $n$ -vertex network (if connected) is a single cluster. After one or more iterations, removing an edge (step 2) causes the network to split into two components (clusters). As further edges are removed, each component again splits, until  $n$  single-vertex clusters remain. The result is a *dendrogram*: a binary tree in which the distance of nodes from the root shows the order in which clusters were split. A cross-section of the dendrogram at any level represents a division of the network into any desired number of clusters.

In step 3, edge betweenness need not be recalculated for the whole network, but only for the component containing the edge removed in step 2, or for both components if removing the edge caused the component to split.

Most algorithms assume that communities are disjoint, placing each vertex in only one cluster. However, in the real world, communities often overlap. For example, in collaboration networks an author might work with researchers in many groups, in biological networks a protein might interact with many groups of proteins, and so on.

In this paper we present a new algorithm to find overlapping community structure in networks. It is a hierarchical, divisive algorithm, based on Girvan and Newman's but extended with a novel method of splitting vertices. We describe the design of the algorithm in Section 2. In Section 3 we present some results on both artificial (computer-generated) and real-world networks. Section 4 compares our algorithm with a few others that can detect overlapping communities. Conclusions appear in Section 5.

## 2 Finding Overlapping Clusters

In any divisive hierarchical clustering algorithm, clusters are repeatedly divided into smaller (normally disjoint) clusters that together contain the same items. To allow overlapping clusters, there needs to be some way of splitting (copying) an item so that it can be included in more than one cluster when the cluster divides.

In the context of network clustering, assuming it is based entirely on the network structure, it seems reasonable to assume that each vertex should be in the same cluster as at least one of its neighbours, unless it is in a singleton cluster or no cluster at all. Therefore, a vertex  $v$  should be split into at most  $d(v)$  copies, where  $d(v)$  is the degree of  $v$ . However, we still need to decide how many times a vertex should be split, and when a vertex should be split (e.g., at the beginning or when dividing a cluster).

Our algorithm extends Girvan and Newman's algorithm (the "GN algorithm") with a specific method of deciding *when* and *how* to split vertices. As in the original work, we only consider unipartite networks with undirected, unweighted edges. We name our new algorithm "CONGA" (Cluster-Overlap Newman Girvan Algorithm).

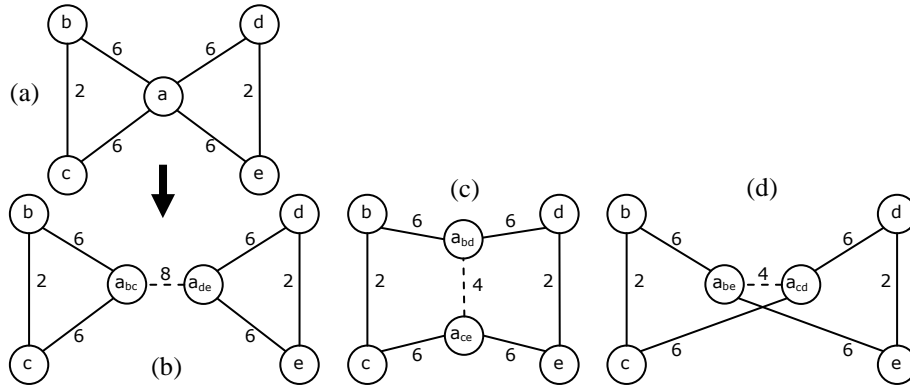
**Splitting Vertices.** In the GN algorithm, the basic operation is removing an edge. We introduce a second operation: splitting a vertex. If split, a vertex  $v$  always splits into *two* vertices  $v_1$  and  $v_2$ : edges with  $v$  as an endvertex are redirected to  $v_1$  or  $v_2$  such that  $v_1$  and  $v_2$  each has at least one edge. By splitting repeatedly, a vertex  $v$  can eventually split into at most  $d(v)$  vertices. Vertices are split incrementally during the clustering process. This binary splitting fits well into the GN algorithm because, like removing an edge, splitting a vertex may cause its cluster to split into two.

**Split Betweenness.** The key point of the CONGA algorithm is the notion of “split betweenness”. This provides a way to decide (1) *when* to split a vertex, instead of removing an edge, (2) *which* vertex to split, and (3) *how* to split it. Clearly,  $v$  should only be split into  $v_1$  and  $v_2$  if these two vertices belong to different clusters. We could verify this by counting the number of shortest paths that would pass between  $v_1$  and  $v_2$  if they were joined by an edge. Then, if there were more shortest paths on  $\{v_1, v_2\}$  than on any real edge, the vertex should be split; otherwise, an edge should be removed as usual. This is the basis of our method of splitting a vertex, which is as follows.

For any split of vertex  $v$  into  $v_1$  and  $v_2$ , we add a new “imaginary” edge between  $v_1$  and  $v_2$ . If  $u$  is a neighbour of  $v_1$  and  $w$  is a neighbour of  $v_2$ , all shortest paths that passed through  $v$  along edges  $\{u, v\}$ ,  $\{v, w\}$  now pass along  $\{u, v_1\}$ ,  $\{v_1, v_2\}$ ,  $\{v_2, w\}$ . The imaginary edge has zero cost: the lengths of paths traversing it are unchanged, and no new shortest paths are created: paths beginning from  $v$  do not traverse this edge. We then calculate the betweenness  $c_B(\{v_1, v_2\})$  of the imaginary edge. In general, there are  $2^{d(v)-1} - 1$  ways to split  $v$  into two. We call the split that maximizes  $c_B(\{v_1, v_2\})$  the *best split* of  $v$ , and the maximum value of  $c_B(\{v_1, v_2\})$  the *split betweenness* of  $v$ .

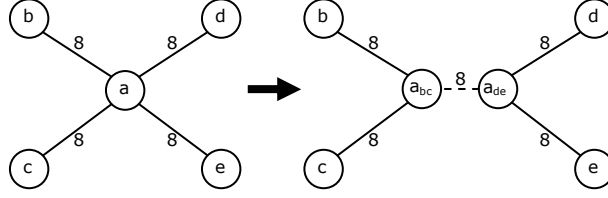
We modify the GN algorithm so that, at each step, it considers the split betweenness of every vertex as well as the edge betweenness of every edge. If the maximum split betweenness is greater than the maximum edge betweenness, the corresponding vertex is split, using its best split. (Note that imaginary edges are never actually added to the network, but are used only during the calculation of the split betweenness.)

Fig. 1(a) shows a network comprising two overlapping clusters:  $\{a, b, c\}$  and  $\{a, d, e\}$ . Labels on the edges show edge betweennesses (with shortest paths counted in both directions). Fig. 1(b) shows  $a$ ’s best split into  $a_{bc}$  and  $a_{de}$ , with the imaginary edge (betweenness 8) shown as a dashed line. Fig. 1(c-d) shows the other two possible splits of  $a$ . In these, the imaginary edge has a lower betweenness, 4, proving that the split of Fig. 1(b) is the best split and the split betweenness of  $a$  is 8. Because this is greater than any edge betweenness,  $a$  should indeed be split.



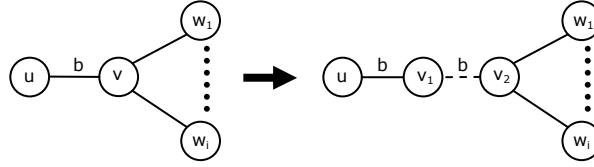
**Fig. 1.** (a) Network. (b) Best split of vertex  $a$ . (c), (d) Other splits of vertex  $a$

Fig. 2 shows a network which does not exhibit clustering. Here, any (2+2) split of  $a$  is a best split. The split betweenness of  $a$  is 8, which is the same as the betweenness of each edge. Therefore, by default, we remove any edge instead of splitting  $a$ .



**Fig. 2.** Best split of vertex  $a$ : split betweenness of  $a$  is 8

Our method will never split a vertex into  $v_1$  and  $v_2$  such that  $v_1$  has only one neighbour,  $u$ . This is because the betweenness of  $\{v_1, v_2\}$  would be the same as that of  $\{u, v_1\}$ , as shown in Fig. 3, so removing edge  $\{u, v\}$  would be preferred over splitting  $v$ . As a consequence of this, vertices with degree less than 4 are never split. In general, there are now only  $2^{d(v)-1} - d(v) - 1$  ways to split a vertex into two.



**Fig. 3.** Vertex will not split into vertices with degree 1

**Vertex Betweenness and Split Betweenness.** When splitting a vertex  $v$  into  $v_1$  and  $v_2$ , the split betweenness is the number of shortest paths passing between any neighbour of  $v_1$  and any neighbour of  $v_2$  via  $v_1$  and  $v_2$ . By definition, this is no greater than the number of shortest paths passing through  $v$ : the *vertex betweenness* of  $v$  [6]. It is simple to calculate vertex betweenness  $c_B(v)$  from edge betweenness  $c_B(e)$ :

$$c_B(v) = \frac{1}{2} \sum_{e \in \Gamma(v)} c_B(e) - (n-1) \quad (1)$$

where  $\Gamma(v)$  is the set of edges with  $v$  as an endvertex and  $n$  is the number of vertices in the component [5]. Therefore, as an optimization, we can use vertex betweenness as an upper bound on split betweenness: if the vertex betweenness of  $v$  is no greater than the maximum edge betweenness, there is no need to calculate  $v$ 's split betweenness.

**Calculating Split Betweenness.** To calculate the split betweenness, and best split, of a vertex  $v$ , we first compute the pair betweennesses of  $v$ . The *pair betweenness* of  $v$  for  $\{u, w\}$ , where  $u$  and  $w$  are neighbours of  $v$  and  $u \neq w$ , is the number of shortest paths that traverse both edges  $\{u, v\}$  and  $\{v, w\}$ . The vertex betweenness of  $v$  is the sum of all of its pair betweennesses.

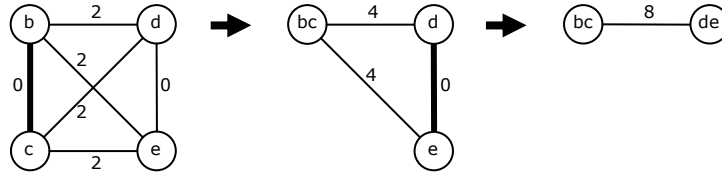
We can represent the pair betweennesses of  $v$ , degree  $k$ , by a  $k$ -clique in which each vertex is labelled by one of  $v$ 's neighbours and each edge  $\{u, w\}$  is labelled by the pair betweenness “score” of  $v$  for  $\{u, w\}$ . Then, to find the best split of  $v$ :

1. Choose edge  $\{u,w\}$  with minimum score.
2. Coalesce  $u$  and  $w$  to a single vertex,  $uw$ .
3. For each vertex  $x$  in the clique, replace edges  $\{u,x\}$ , score  $b_1$ , and  $\{w,x\}$ , score  $b_2$ , by a new edge  $\{uw,x\}$  with score  $b_1+b_2$ .
4. Repeat from step 1  $k-2$  times (in total).

The labels on the remaining two vertices show the split, and the score on the remaining edge is the split betweenness.

This algorithm is not guaranteed to find the best split. To do that, we would need to try *all* edges in step 1 of each iteration, which would require exponential time. Our “greedy” method is much more efficient and, in practice, usually finds the best split or a close approximation to it. Fig. 4 shows how it finds the best split of vertex  $a$  of Fig. 1.

1. There are  $k-2 = 2$  phases; the edge chosen in step 1 of each phase is highlighted.



**Fig. 4.** Finding the best split of vertex  $a$  of Fig. 1

**Calculating Pair Betweennesses.** Pair betweennesses are computed while calculating edge betweenness, by a straightforward modification of the GN algorithm. The GN algorithm increments the betweenness of edge  $\{i,j\}$  for all shortest paths beginning at each vertex  $s$ . CONGA does this *and* increments the pair betweennesses of  $i$  for all pairs  $\{j,k\}$  such that  $k$  is a neighbour of  $i$  on a path between  $s$  and  $i$ .

There is some overhead, in both time and space, in computing pair betweennesses during the betweenness calculation. In most cases this information is not used because we can often determine, from the vertex betweenness, that a vertex should not be split. Therefore, our betweenness calculation is split into two phases, as shown below.

**The CONGA Algorithm.** Our complete algorithm is as follows:

1. Calculate edge betweenness of all edges in network.
2. Calculate vertex betweenness of vertices, from edge betweennesses, using Eq. (1).
3. Find candidate set of vertices: those whose vertex betweenness is greater than the maximum edge betweenness.
4. If candidate set is non-empty, calculate pair betweennesses of candidate vertices, and then calculate split betweenness of candidate vertices, using Eq. (1).
5. Remove edge with maximum edge betweenness or split vertex with maximum split betweenness (if greater).
6. Recalculate edge betweenness for all remaining edges in same component(s) as removed edge or split vertex.
7. Repeat from step 2 until no edges remain.

**Complexity and Efficiency.** The GN algorithm has a worst-case time complexity of  $O(m^2n)$ , where  $m$  is the number of edges and  $n$  is the number of vertices. In CONGA, each vertex splits into an average of up to  $2m/n$  vertices, so the number of vertices after splitting is  $O(m)$ ; the number of iterations is still  $O(m)$  and the number of edges is unchanged. This makes the time complexity  $O(m^3)$  in the worst case.

In practice, the speed depends on the number of vertices that are split. If more are split, more iterations are needed, the network becomes larger, and step 4 needs to be performed more frequently. Conversely, vertex splitting can cause the network to decompose into separate components more readily, which reduces the execution time.

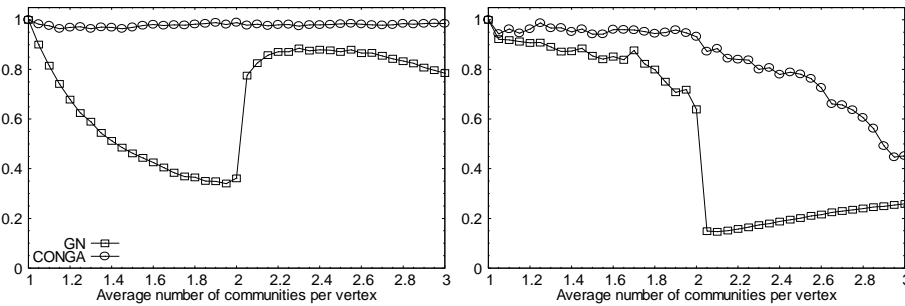
### 3 Results

We have tested CONGA on computer-generated networks based on a known, possibly overlapping, community structure. Each network contains  $n$  vertices divided into  $c$  equally-sized communities, each containing  $nr/c$  vertices. Vertices are randomly and evenly distributed between communities such that each vertex is a member of  $r$  ( $\geq 1$ ) communities on average. Edges are randomly placed between pairs of vertices with probability  $p_{in}$  if the vertices belong to the same community and  $p_{out}$  otherwise.

We assess how well an algorithm can recover the community structure from a network by using the algorithm to compute  $c$  clusters and comparing the result with the  $c$  known communities. We calculate two values (all averaged over 10 graphs):

- *recall*: the fraction of vertex pairs belonging to the same community that are also in the same cluster.
- *precision*: the fraction of vertex pairs in the same cluster that also belong to the same community.

First (Fig. 5), we generated networks of 256 vertices divided into 32 communities, set  $p_{out} = 0$  (i.e., no intercommunity edges) and  $p_{in} = 0.5$ , and increased the amount of overlap from  $r = 1$  (i.e., no overlap) to  $r = 3$ . The number of edges (and hence the average degree) increases roughly quadratically with  $r$ , because the average community size is proportional to  $r$  and each vertex is a member of  $r$  communities. So the average degree is 4 for  $r = 1$  but increases to approximately 15 for  $r = 2$  and 32 for  $r = 3$ .

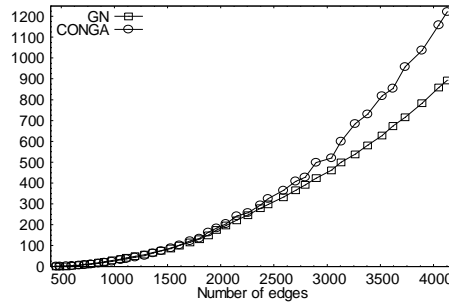


**Fig. 5.** recall (left), precision (right):  $n=256$ ,  $c=32$ ,  $p_{in}=0.5$ ,  $p_{out}=0$ , various  $r$

For the GN algorithm, as  $r$  increases, recall declines steadily because the (non-overlapping) clusters are smaller than the communities; precision is quite high, though certainly not perfect, in this range. Suddenly, at around  $r = 2$ , recall increases and precision decreases, as most vertices are placed in a single cluster. In contrast, CONGA behaves very well up to about  $r = 2$  and then deteriorates gradually.

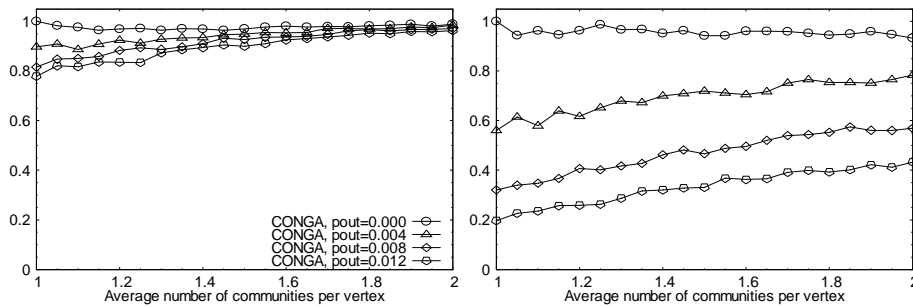
We have repeated this experiment for various values of  $c$  and  $p_{in}$ . The curves always show a similar shape, though the value of  $r$  at which precision drops varies.

Fig. 6 shows the execution times of our implementation of both GN and CONGA algorithms, implemented in Java and run on a 2.4GHz Pentium 4 processor. Execution times increase with the number of edges, which varies from 453 ( $r = 1$ ) to 4128 ( $r = 3$ ). CONGA is expected to be slower than the GN algorithm. However, in the cases where it gives good results (up to about  $r = 2$ , or 2000 edges, for these examples) its execution time is similar to, or even less than, that of the GN algorithm.



**Fig. 6.** Execution time (seconds) for  $n=256$ ,  $c=32$ ,  $p_{in}=0.5$ ,  $p_{out}=0$ , various  $r$

Fig. 7 shows the effect on CONGA of combining overlapping communities with intercommunity edges. The precision reduces as  $p_{out}$  increases, because the extra edges mislead the algorithm into placing some vertices into too many clusters. For example, for  $r = 1$  and  $p_{out} = 0.012$ , each vertex has an average of three intercommunity edges and four intracommunity edges.



**Fig. 7.** recall (left), precision (right):  $n=256$ ,  $c=32$ ,  $p_{in}=0.5$ , various  $r$ ,  $p_{out}$

To evaluate the algorithm on real-world networks, there is no correct solution with which to compare, so the quality of a clustering must be assessed in a different way.

This is usually done by measuring the relative number of intracluster and intercluster edges, for example, by the *modularity* measure [13, 14]. However, modularity is defined only for disjoint clusters and it is not clear how to extend the definition.

When clusters may overlap, it is generally possible to avoid an intercluster edge by including one or both endvertices in both of the clusters that it joins, but this is not necessarily desirable. A good clustering should minimize both the number of intercluster edges *and* the amount of overlap. To measure this we use the *average degree* measure, previously used for overlapping clusters in [3]. We define the *vertex average degree* (*vad*) of a set of clusters  $S$ , where each cluster is a set of vertices, as:

$$vad(S) = \frac{2 \sum_{C \in S} |E(C)|}{\sum_{C \in S} |C|} \quad (2)$$

where  $E(C)$  is the set of edges with both endvertices in  $C$ .

Like modularity, the absolute value of *vad* has no real meaning, but can potentially be used to compare algorithms and to decide how many clusters to divide a network into. Another useful measure is the *overlap* of a set of clusters  $S$ : the sum of the cluster sizes divided by the size of the union of all clusters. We have run the CONGA and GN algorithms on several real-world examples, listed in Table 1.

**Table 1.** Algorithm’s results on real-world networks

Name	Ref.	Vertices	Edges	Runtime (s)
Karate club	[18]	34	78	0.2
Dolphins	[9]	62	159	0.5
College football	[7]	115	613	7.8
Network science	[12]	379	914	12.5
Blogs	[17]	3982	6803	30411
Words	[10]	1000	3471	6767

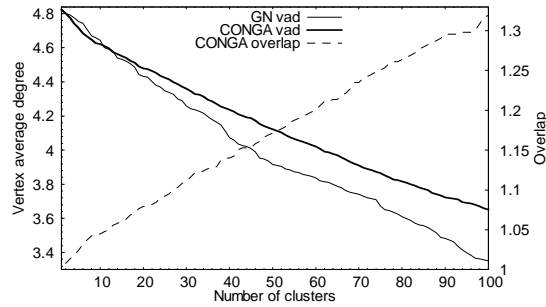
“Karate club” [18], discussed in [7], represents a social network based on two disjoint real-world communities. The communities are not reflected clearly in the network structure: there are eight intercommunity edges. GN finds an almost perfect two-cluster solution, misclassifying one vertex. CONGA finds a different solution with a small overlap, 1.03. The *vad* is 4.45 for CONGA and 4.0 for GN, suggesting that the overlapping clustering is a good one (albeit incorrect).

“Dolphins” [9], discussed in [14], is a social network of dolphins, also based on two disjoint communities. Here there are only six intercommunity edges. GN finds the two communities correctly and CONGA finds the same division but with two vertices from the larger community included in both clusters: the overlap is 1.03. The *vad* is 4.91 for CONGA and 4.94 for GN.

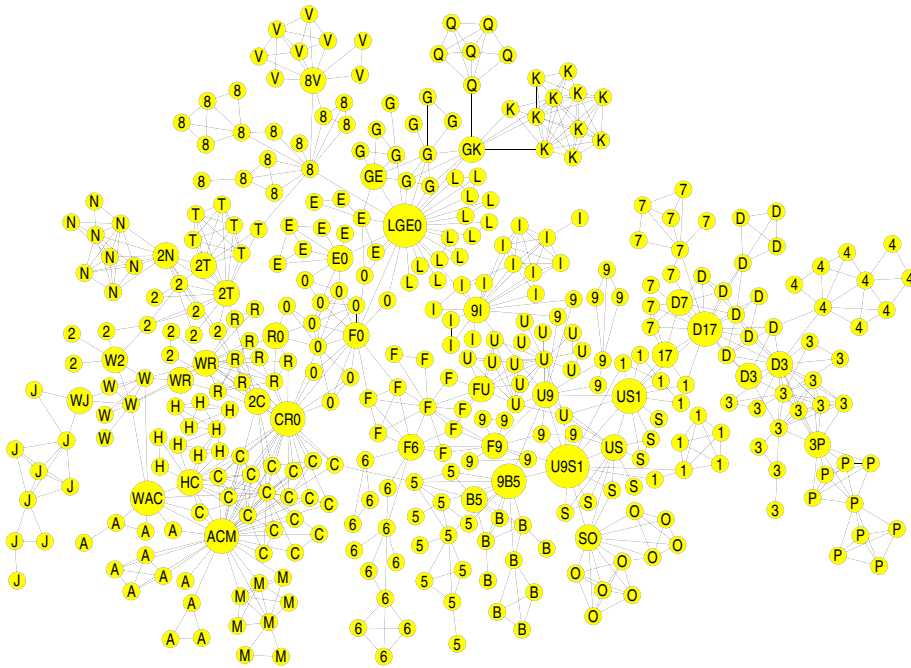
“College football” [7] is a network based on games between teams that belong to 15 disjoint real-world communities. This network has many intercommunity edges. Neither algorithm finds a perfect 15-cluster solution; the one found by CONGA has a lower *vad* (5.87 vs. 7.18) and a large overlap: 1.75.



“Network science” [12] is a collaboration network of coauthorships. For such networks it may be impossible to determine the number of real-world communities, and it seems reasonable to assume they might overlap. Fig. 8 shows the *vad* of solutions found by both GN and CONGA for up to 100 clusters, and the overlap of the CONGA clusterings. CONGA’s solution has a higher *vad* than GN’s for 14 or more clusters. CONGA’s solution for 33 clusters is illustrated in Fig. 9: each cluster is identified by a letter or digit and each vertex is labelled with the cluster(s) to which it belongs.



**Fig. 8.** Network science collaboration network: *vad* and overlap for GN and CONGA



**Fig. 9.** Network science collaboration network divided into 33 overlapping clusters

“Blogs” [17] is a network of blogs on the MSN (now known as Windows Live™ Spaces) platform. An edge links people who have left two or more comments on each

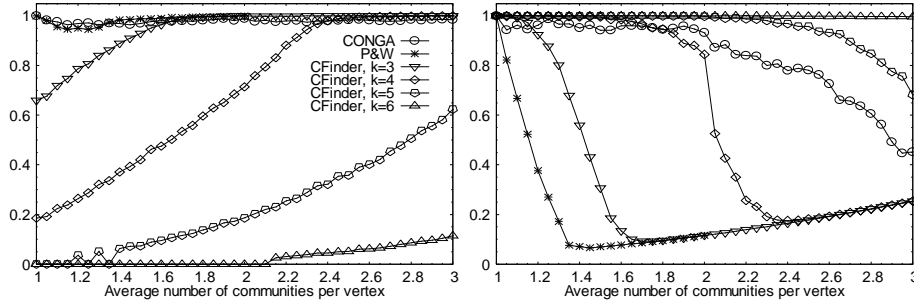
other’s blog, and so are deemed to be acquainted. CONGA’s solution has a consistently higher *vad* than GN’s, especially for more than 90 clusters. The overlap increases with the number of clusters but levels off, reaching a maximum of 1.39.

“Words” is a non-social network: a contiguous 1000-vertex subgraph of a word association network from [15], converted from an original directed, weighted version [10]. CONGA successfully groups words with similar meaning. For example, dividing it into 400 clusters, the word “form” appears in four: {contract, document, form, order, paper, signature, write}, {blank, entry, fill, form, up}, {compose, create, form, make}, {form, mold, shape}. Again, the *vad* for CONGA’s solution is consistently higher than GN’s; the overlap increases and tails off, reaching a maximum of 2.23.

## 4 Related Work

Pinney and Westhead [16] have also proposed extending the GN algorithm with the ability to split vertices between clusters. The decision of whether to split a vertex or remove an edge is based entirely on edge betweenness and vertex betweenness. The highest-betweenness edge is removed only if its two endvertices have similar (within 10%) betweenness; otherwise the highest-betweenness vertex is split. When splitting a vertex, it is removed and “copied into the child subnetworks”; we assume that this means that  $v$  splits into  $d(v)$  vertices, each linked to one of  $v$ ’s neighbours.

We have implemented this algorithm with the above assumption. Fig. 10 (the “P&W” curve) shows its results on the random networks of Section 3. Although the recall is good, the precision drops rapidly at a much smaller overlap than CONGA’s. This appears to be because of the prolific way in which vertices are copied, and is worse for networks with a higher average degree.



**Fig. 10.** recall (left), precision (right):  $n=256$ ,  $c=32$ ,  $p_{in}=0.5$ ,  $p_{out}=0$ , various  $r$

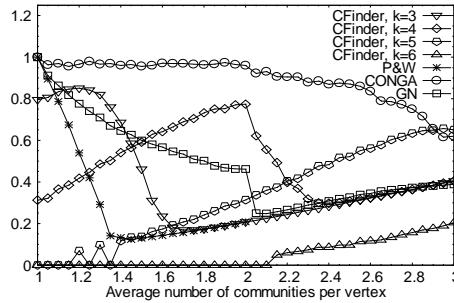
The clique percolation algorithm of Palla *et al.* [15], implemented in CFinder [1], finds overlapping clusters in a different way. Instead of dividing a network into its most loosely connected parts, it identifies the most densely connected parts. The parameter is not the number of clusters to be found but their density,  $k$ . A cluster is defined as the set of  $k$ -cliques that can all be reached from each other via a sequence of adjacent  $k$ -cliques; two  $k$ -cliques are *adjacent* if they share  $k-1$  vertices. Each vertex may be in many clusters, or even none: e.g., degree-1 vertices are always ignored.

We have run this algorithm (CFinder v1.21) on the random networks of Section 3. Fig. 10 (the “CFinder” curves) shows results for values of  $k$  from 3 to 6. The recall is low, especially for large  $k$ , but increases as the overlap, and hence average degree, of the network increases. The precision drops as overlap increases. For  $k \geq 5$  precision is better than CONGA’s, but recall is very low. This algorithm is useful if we only wish to find the densest parts of communities, and is faster than CONGA.

Baumes *et al.* [2, 3] present a collection of algorithms to find overlapping clusters. One algorithm iteratively improves a candidate cluster by adding vertices to and removing vertices from it while its density improves. Another removes vertices from a network until it breaks into disjoint components, forming initial clusters, and then replaces each removed vertex into one or more of the clusters, which might overlap.

Li *et al.* [8] form overlapping clusters using both the structure of the network *and* the content of vertices and edges. The first phase of their algorithm finds densely connected “community cores”, similarly to the method of [15]. In the second phase, clusters are formed from cores by adding further triangles and edges whose content (assessed using keywords) is similar to that of the core.

To conclude this section, we plot in Fig. 11 the F-measure (harmonic mean of recall and precision) for all algorithms tested. CONGA has the best score for  $r$  up to 2.9.



**Fig. 11.** F-measure for random networks:  $n=256$ ,  $c=32$ ,  $p_{in}=0.5$ ,  $p_{out}=0$ , various  $r$

## 5 Conclusions

We have presented an algorithm that seems to be effective in discovering overlapping communities in networks. Good results have been obtained for a range of random networks with overlap of more than 2, which is large relative to the number of communities: if a network has only 32 communities, an overlap of 3 means that each vertex is in the same community as  $1/4$  of the whole network. As the number of communities is increased, the algorithm can cope with a larger overlap. The algorithm is not fast, but its speed is comparable with that of the GN algorithm from which it is derived.

Future work includes trying to improve the algorithm further and applying similar ideas to faster clustering algorithms than the GN algorithm. It is also worth evaluating how well the *vad* (vertex average degree) measures the quality of an overlapping clustering. Finally, it would be interesting to investigate the overlapping nature of real-world networks, a subject that has received little attention (but see [15]). For example,

Fig. 8 might suggest that the collaboration network (Fig. 9) can be broken down into a small number (up to 10-14) of disjoint clusters, maybe corresponding to research groups, but to decompose it further requires clusters to overlap.

Further information related to this paper, including the networks analysed and more results, can be found at <http://www.cs.bris.ac.uk/~steve/networks/>.

**Acknowledgements.** Many thanks are due to Peter Flach for his expert comments on several drafts of this paper.

## References

1. Adamcsek, B., Palla, G., Farkas, I., Derényi, I., and Vicsek, T. CFinder: locating cliques and overlapping modules in biological networks. *Bioinformatics* 22, 1021-1023. 2006.
2. Baumes, J., Goldberg, M., Krishnamoorthy, M., Magdon-Ismael, M., and Preston, N. Finding communities by clustering a graph into overlapping subgraphs. *Proc. IADIS Applied Computing 2005*, 97-104. 2005.
3. Baumes, J., Goldberg, M., and Magdon-Ismael, M. Efficient identification of overlapping communities. *Intelligence and Security Informatics* (LNCS 3495), 27-36. Springer. 2005.
4. Brandes, U., Gaertler, M., and Wagner, D. Experiments on graph clustering algorithms. In *Algorithms - ESA 2003* (LNCS 2832), 568-579. Springer. 2003.
5. Brandes, U. and Erlebach, T. (eds.). *Network Analysis: Methodological Foundations* (LNCS 3418). Springer. 2005.
6. Freeman, L.C. A set of measures of centrality based on betweenness. *Sociometry* 40, 35-41. 1977.
7. Girvan, M. and Newman, M.E.J. Community structure in social and biological networks. *Proc. Natl. Acad. Sci. USA* 99, 7821-7826. 2002.
8. Li, X., Liu, B., and Yu, P.S. Discovering overlapping communities of named entities. *Knowledge Discovery in Databases: PKDD 2006* (LNCS 4213), 593-600. Springer. 2006.
9. Lusseau, D., Schneider, K., Boisseau, O.J., Haase, P., Slooten, E., and Dawson, S.M. The bottlenose dolphin community of Doubtful Sound features a large proportion of long-lasting associations. *Behavioral Ecology and Sociobiology* 54, 396-405. 2003.
10. Nelson, D.L., McEvoy, C.L., and Schreiber, T.A. The University of South Florida word association, rhyme and word fragment norms. <http://w3.usf.edu/FreeAssociation/>. 1998.
11. Newman, M.E.J. Fast algorithm for detecting community structure in networks, *Phys. Rev. E* 69, 066133. 2004.
12. Newman, M.E.J. Finding community structure in networks using the eigenvectors of matrices, *Phys. Rev. E* 74, 036104. 2006.
13. Newman, M.E.J. Modularity and community structure in networks, *Proc. Natl. Acad. Sci. USA* 103, 8577-8582. 2006.
14. Newman, M.E.J., and Girvan, M. Finding and evaluating community structure in networks, *Phys. Rev. E* 69, 026113. 2004.
15. Palla, G., Derényi, I., Farkas, I., and Vicsek, T. Uncovering the overlapping community structure of complex networks in nature and society. *Nature* 435, 814-818. 2005.
16. Pinney, J.W. and Westhead, D.R. Betweenness-based decomposition methods for social and biological networks. In Barber, S., Baxter, P.D., Mardia, K.V., and Walls, R.E. (eds.), *Interdisciplinary Statistics and Bioinformatics*, 87-90. Leeds University Press. 2006.
17. Xie, N. Social network analysis of blogs. MSc dissertation. University of Bristol. 2006.
18. Zachary, W.W. An information flow model for conflict and fission in small groups. *Journal of Anthropological Research* 33, 452-473. 1977.