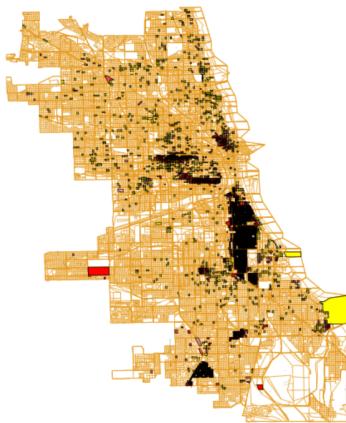


Detecting Community Structures in Social Networks by Graph Sparsification

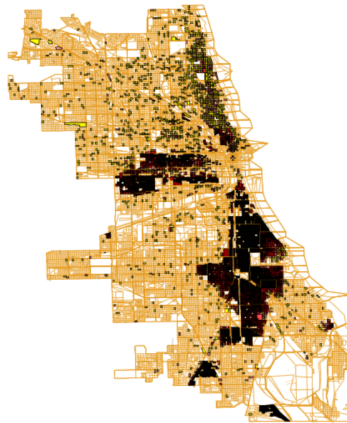
Partha Basuchowdhuri, **Satyaki Sikdar**, Sonu Shreshtha, Subhasis Majumder

Department of Computer Science and Engineering,
Heritage Institute of Technology, Kolkata, India





(a) *Chicago, 1940*



(b) *Chicago, 1960*

Figure: The tendency of people to live in racially homogeneous neighborhoods[1]. In yellow and orange blocks % of Afro-Americans ≤ 25 , in brown and black boxes % ≥ 75 .

Definition of a Community

For a given graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, find a cover $\mathbb{C} = \{C_1, C_2, \dots, C_k\}$ such that $\bigcup_i C_i = \mathcal{V}$.

- For disjoint communities, $\forall i, j$ we have $C_i \cap C_j = \emptyset$
- For overlapping communities, $\exists i, j$ where $C_i \cap C_j \neq \emptyset$

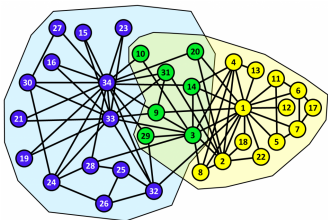


Figure: Zachary's Karate Club Network

$\mathbb{C} = \{C_1, C_2, C_3\}$, $C_1 = \text{yellow}$ nodes, $C_2 = \text{green}$, $C_3 = \text{blue}$ is a disjoint cover

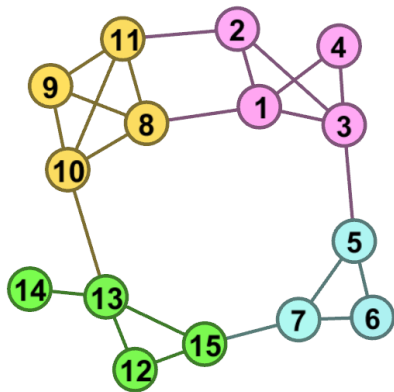
However, $\bar{\mathbb{C}} = \{\bar{C}_1, \bar{C}_2\}$, $\bar{C}_1 = \text{yellow}$ & green nodes and $\bar{C}_2 = \text{blue}$ & green nodes is an overlapping cover

For our problem, we concentrate on *disjoint* community detection



A Little Background: Edge Betweenness Centrality

$$c_B(e) = \sum_{\substack{s, t \in \mathcal{V} \\ s \neq t}} \frac{\sigma(s, t \mid e)}{\sigma(s, t)}$$



Top 6 edges

Edge	$c_B(e)$	Type
(10, 13)	0.3	inter
(3, 5)	0.23333	inter
(7, 15)	0.2079	inter
(1, 8)	0.1873	inter
(13, 15)	0.1746	intra
(5, 7)	0.1476	intra

Bottom 6 edges

Edge	$c_B(e)$	Type
(8, 11)	0.022	intra
(1, 2)	0.0269	intra
(9, 11)	0.031	intra
(8, 9)	0.0412	intra
(12, 15)	0.052	intra
(3, 4)	0.060	intra



The Girvan-Newman Algorithm

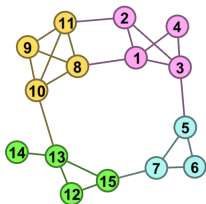
Proposed by Michelle Girvan and Mark Newman[2] in 2002

The Key Ideas

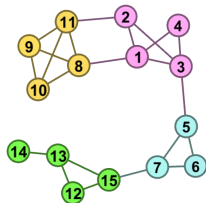
- Based on reachability of nodes - shortest paths
- Edges are selected on the basis of the edge betweenness centrality

The algorithm

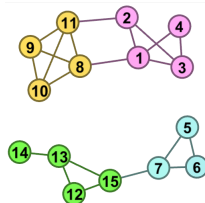
- 1 Compute centrality for all edges
- 2 Remove edge with largest centrality; ties can be broken randomly
- 3 Recalculate the centralities on the running graph
- 4 Iterate from step 2, stop when you get clusters of desirable quality



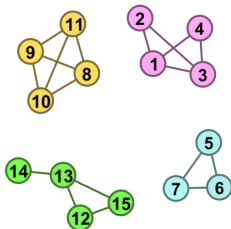
(a) Best edge: (10, 13)



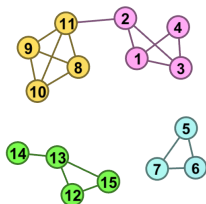
(b) Best edge: (3, 5)



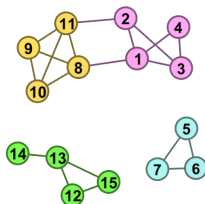
(c) Best edge: (7, 15)



(f) Final graph



(e) Best edge: (2, 11)



(d) Best edge: (1, 8)

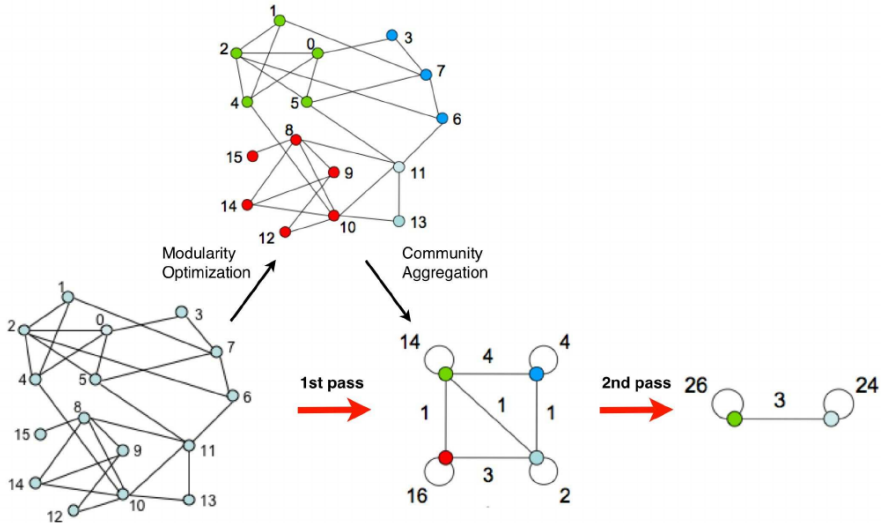


Louvain Method: A Greedy Approach

- Proposed by Blondel et al[3] in 2008
- Takes the greedy maximization approach
- Very fast in practice, it's the current state-of-the-art in disjoint community detection
- Performs hierarchical partitioning, stopping when there cannot be any further improvement in modularity
- Contracts the graph in each iteration thereby speeding up the process



Louvain Method in Action





Our Method

Input: An unweighted network $\mathcal{G}(\mathcal{V}, \mathcal{E})$

Output: A disjoint cover \mathbb{C}

- 1 Use Jaccard coefficient to turn \mathcal{G} into a weighted network $\mathcal{G}(\mathcal{V}, \mathcal{E}, \mathcal{W})$
- 2 Construct an t -spanner of $\mathcal{G}(\mathcal{V}, \mathcal{E}, \mathcal{W})$. Take the complement of \mathcal{G}_S , call it \mathcal{G}_{comm}
- 3 Use LINCOM to break \mathcal{G}_{comm} into small but pure fragments
- 4 Use the second phase of Louvain Method to piece all the small bits and pieces together to get \mathbb{C}



Jaccard Intro

Definition

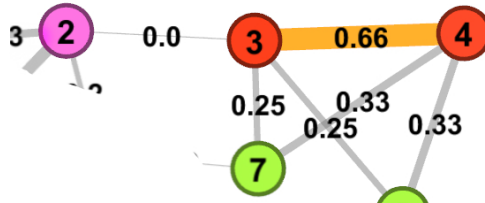
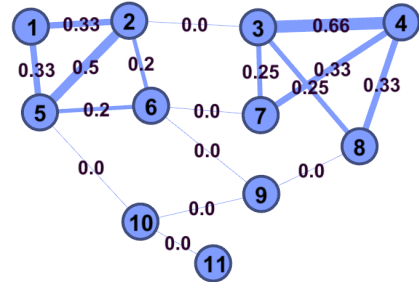
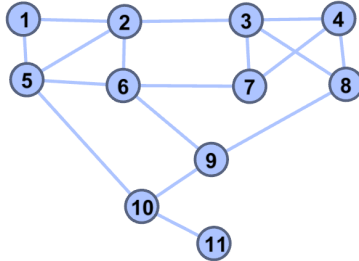
$$w_J(e(v_i, v_j)) = \frac{|\Gamma(v_i) \cap \Gamma(v_j)|}{|\Gamma(v_i) \cup \Gamma(v_j)|}$$

where $\Gamma(v_i)$ is the neighborhood of the node v_i
 $\therefore w_J \in [0, 1]$

- Jaccard works well in domains where local influence is important[4][5][6]
- The computation takes $O(m)$ time



Jaccard Example



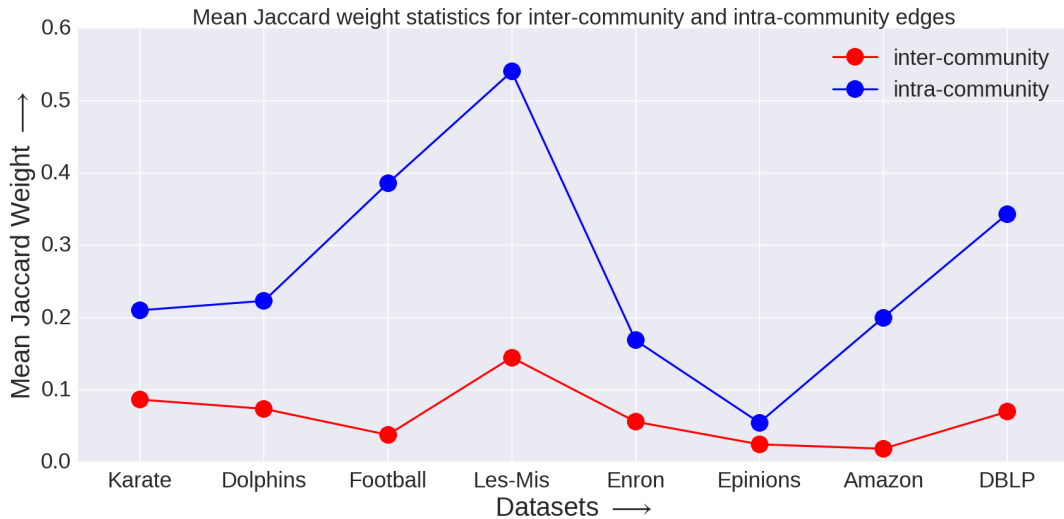


Table: Jaccard weight statistics for top 10% edges in terms of w_J .

Network	$ E $	intra-cluster edge count	top 10% edges in terms of w_J		
			Total edges	Intra-edge	Fraction
Karate	78	21	7	7	1
Dolphin	159	39	15	15	1
Football	613	179	61	61	1
Les-Mis	254	56	25	25	1
Enron	180,811	48,498	18,383	18,220	0.99113
Epinions	405,739	146,417	40,573	36,589	0.90180
Amazon	925,872	54,403	92,587	92,584	0.99996
DBLP	1,049,866	164,268	104,986	104,986	1



Spanner

- A (α, β) -spanner of a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{W})$ is a subgraph $\mathcal{G}_S = (\mathcal{V}, \mathcal{E}_S, \mathcal{W}_S)$, such that,

$$\delta_S(u, v) \leq \alpha \cdot \delta(u, v) + \beta \quad \forall \quad u, v \in V$$

- A t -spanner is a special case of (α, β) spanner where $\alpha = t$ and $\beta = 0$

Authors	Size	Running Time
Althöfer et al. [1993] [7]	$O(n^{1+\frac{1}{k}})$	$O(m(n^{1+\frac{1}{k}} + n \log n))$
Althöfer et al. [1993] [7]	$\frac{1}{2}n^{1+\frac{1}{k}}$	$O(mn^{1+\frac{1}{k}})$
Roditty et al. [2004] [8]	$\frac{1}{2}n^{1+\frac{1}{k}}$	$O(kn^{2+\frac{1}{k}})$
Roditty et al. [2005] [9]	$O(kn^{1+\frac{1}{k}})$	$O(km)$ (det.)
Baswana and Sen [2007] [10]	$O(kn^{1+\frac{1}{k}})$	$O(km)$ (rand.)

Algorithm 1: Greedy Construction of t -spanner
(GREEDY(G, t))

Input : Undirected, weighted graph $G(V, E, W)$, (t)

Output: t -spanner edge set E_S

```
1 begin
2    $E_S \leftarrow \emptyset$ 
3   foreach  $e(s, t) \in E$ , in non-decreasing order of
     weight do
4     if  $d_{(V, E_S)}(s, t) > (2k-1)w(e(s, t))$  then
5        $E_S \leftarrow E_S \cup \{e(s, t)\}$ 
6   return  $E_S$ 
```

Algorithm 2: Spanner based Community Detection
(SPAN(G, k))

Input : Undirected, weighted graph $G(V, E, W)$, (k)

Output: $G_{comm} = (V, E_{comm})$

```

1 begin
2    $E_S \leftarrow \emptyset, C_0 \leftarrow \{\{v\} | v \in V\}, E' \leftarrow E, E_{comm} \leftarrow E$ 
3    $\mathcal{E}'_i \leftarrow \mathcal{E}_i \leftarrow \emptyset, \mathcal{R}_0 \leftarrow V$ 
4   for  $i \leftarrow 1$  to  $k-1$  do
5      $\mathcal{R}_i \leftarrow$  randomly pick  $n^{1-\frac{i}{k}}$  clusters from  $C_{i-1}$ 
6      $\mathcal{E}_i \leftarrow \emptyset, C_i \leftarrow \mathcal{R}_i$ 
7     forall the  $e(u, v) \in \mathcal{E}_{i-1}$  do
8       if  $u, v \in c, c \in \mathcal{R}_i$  then
9          $\mathcal{E}_i \leftarrow \mathcal{E}_i \cup \{e(u, v)\}$ 
10    foreach  $v \in V'$  AND  $c \notin \mathcal{R}_i$  do
11      foreach  $e(v, v'), v' \in c | c \in C_{i-1}$  do
12        if  $c \in \mathcal{R}_i$  then
13           $v_{min} \leftarrow \arg \min_{v' \in \Gamma(v), c} d(v, v')$ 
14        else
15           $v_{min}(c) \leftarrow \arg \min_{v' \in \Gamma(v), c} d(v, v')$ 
16    forall the  $v_{min}, v_{min}(c) \in \Gamma(v), c | c \in C_{i-1}$ 
    do
       $\triangleright v_{min}$  chosen from sampled nodes only

```

```

17
18   if  $v_{min}$  exists for  $v$  then
19      $E_S \leftarrow E_S \cup \{e(v, v_{min})\}$ 
20      $E_{comm} \leftarrow E_{comm} \setminus \{e(v, v_{min})\}$ 
21      $\mathcal{E}_i \leftarrow \mathcal{E}_i \cup \{e(v, v_{min})\}$ 
22      $E' \leftarrow E' \setminus \{E'(v, c_{min})\} |$ 
23        $v_{min} \in c_{min}$ 
24     if  $v_{min}(c) < v_{min}$  then
25        $E_S \leftarrow E_S \cup \{e(v, v_{min}(c))\}$ 
26        $E_{comm} \leftarrow E_{comm} \setminus$ 
27          $\{e(v, v_{min}(c))\}$ 
28        $E' \leftarrow E' \setminus \{E'(v, c) | v_{min}(c) \in c$ 
29     else
30        $E_S \leftarrow E_S \cup \{e(v, v_{min}(c))\}$ 
31        $E_{comm} \leftarrow E_{comm} \setminus \{e(v, v_{min}(c))\}$ 
32        $E' \leftarrow E' \setminus \{E'(v, c) | v_{min}(c) \in c$ 
33   forall the  $e(v, u) | u, v \in c', c' \in E'$  do
34      $E' \leftarrow E' \setminus \{e(v, u)\}$ 
35   forall the  $v \in V'$  do
36     forall the  $c \in C_{k-1}$  do
37        $E_S \leftarrow E_S \cup \{e(v, v_{min}(c))\}$ 
38        $E_{comm} \leftarrow E_{comm} \setminus \{e(v, v_{min}(c))\}$ 
39        $E' \leftarrow E' \setminus \{E'(v, c) | v_{min}(c) \in c$ 
40   return  $G_{comm}$ 

```

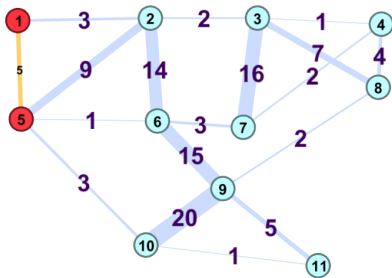



Figure: Original network $n = 11, m = 18$
 $\delta(1, 5) = 5$

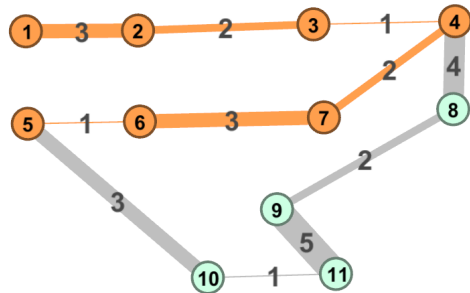


Figure: A 3-spanner of the network
 $n = 11, m = 11$ $\delta_s(1, 5) = 12$

Since $\delta_s(1, 5) < t \cdot \delta(1, 5)$, the edge $(1, 5)$ is discarded
The other edges are discarded in a similar fashion.

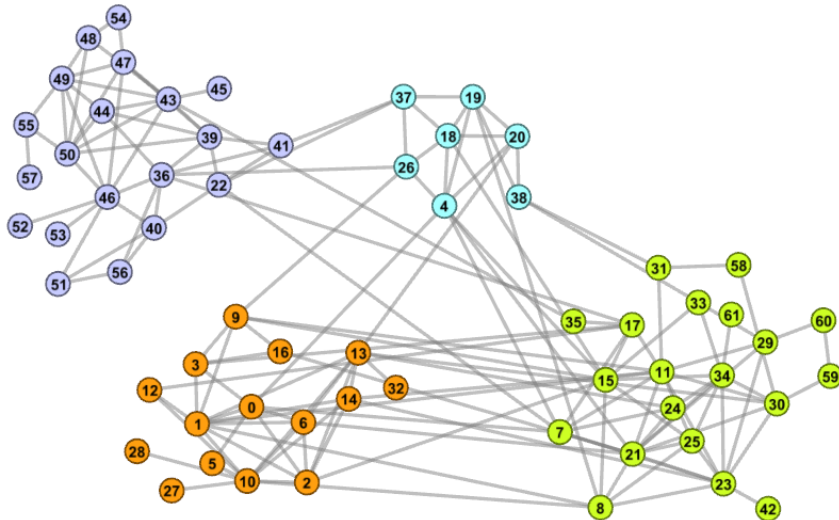


Figure: Dolphin network. $n = 62$, $m = 159$

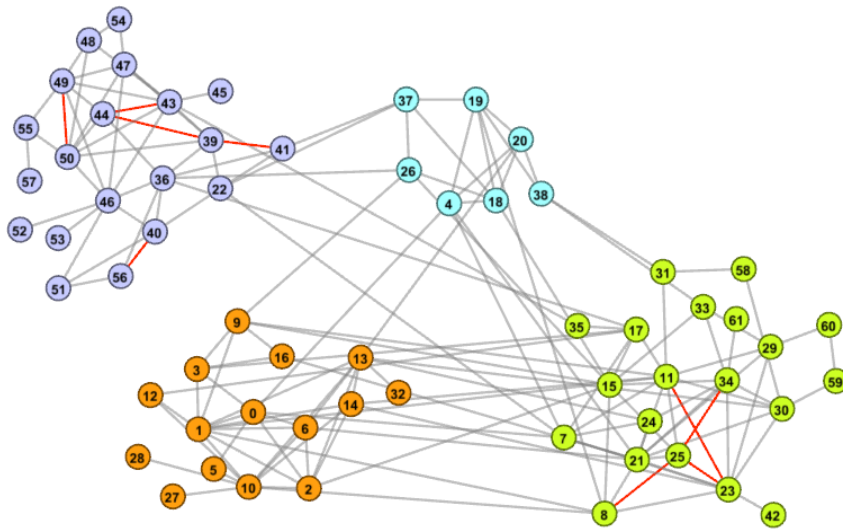


Figure: 3-spanner. $n = 62$, $m = 150$

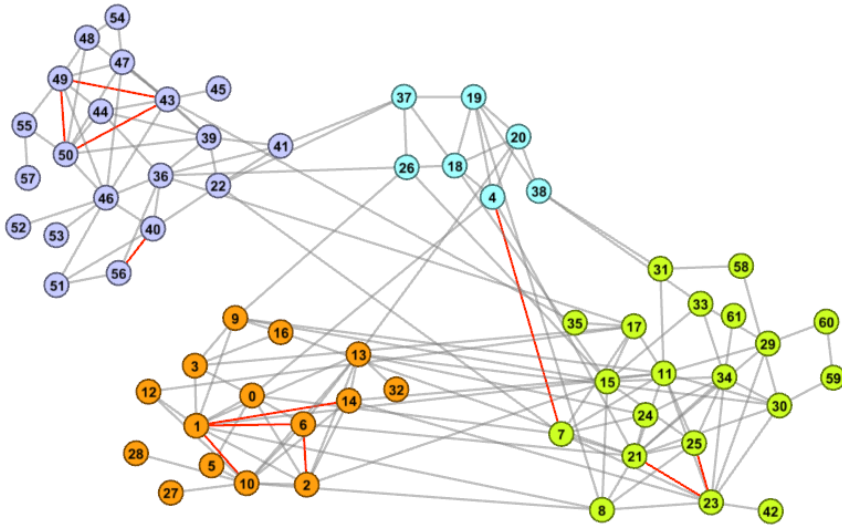


Figure: 5-spanner. $n = 62$, $m = 148$

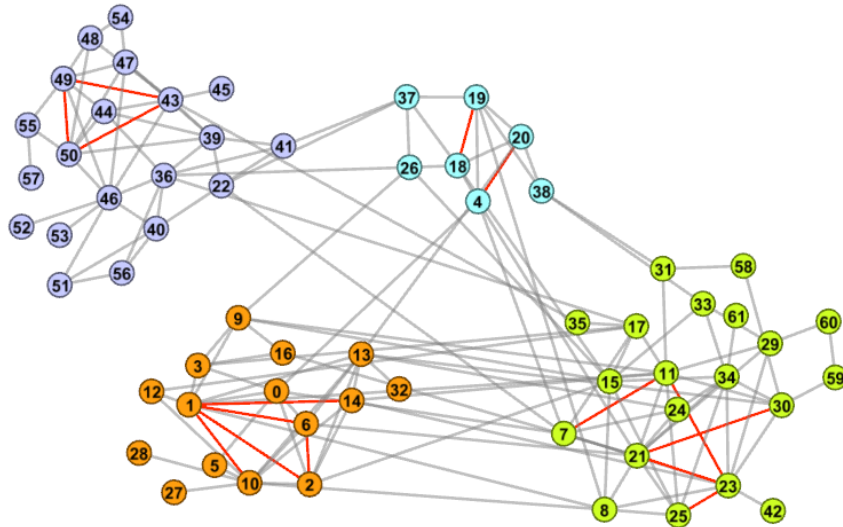


Figure: 7-spanner. $n = 62$, $m = 144$

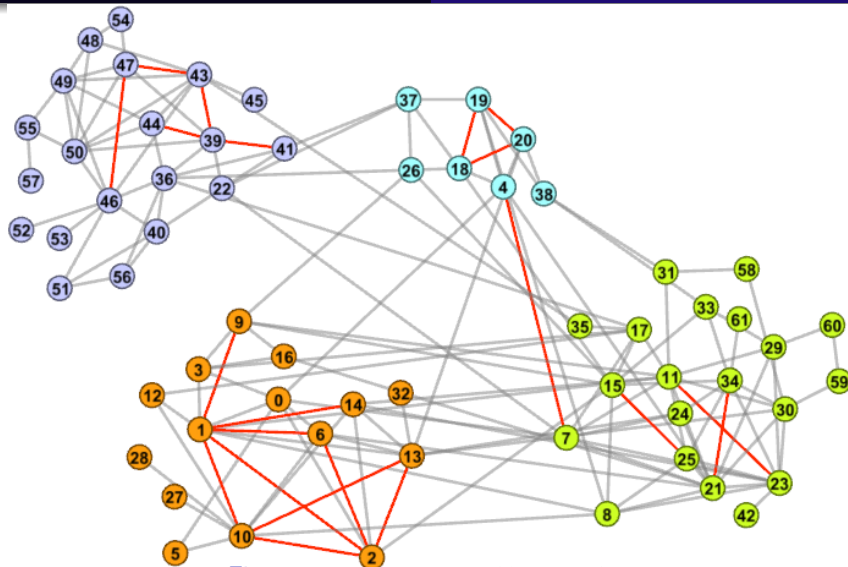
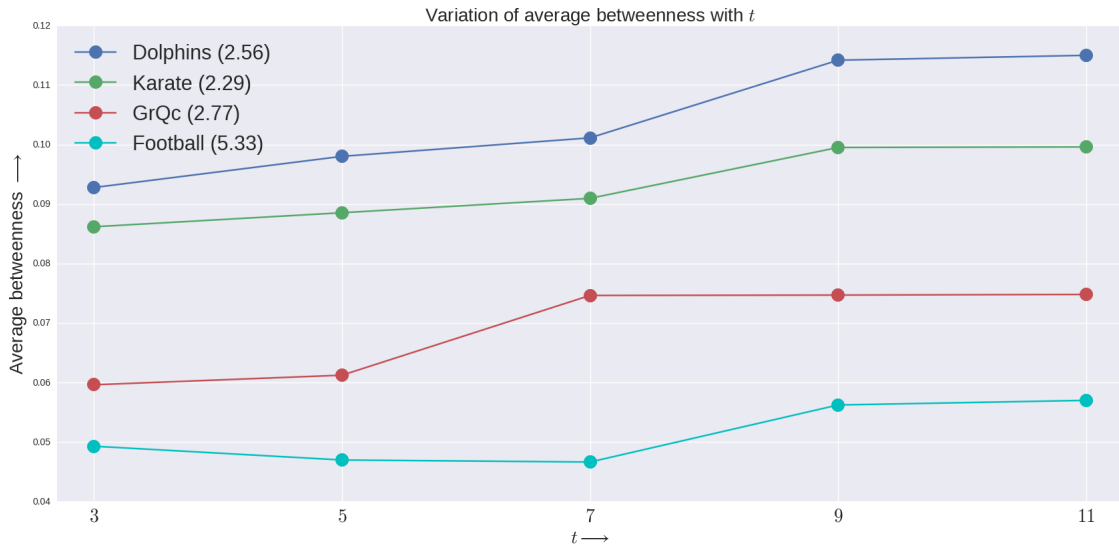


Figure: 9-spanner. $n = 62$, $m = 138$



Name	n	Spanner	#intra-community	#inter-community
Karate	34	Original	59	19
		3	57	19
		5	53	19
		7	51	18
		9	48	19
Dolphin	59	Original	120	39
		3	117	38
		5	102	38
		7	100	38
		9	90	38
Football	115	Original	447	163
		3	385	166
		5	376	166
		7	293	166
		9	286	165

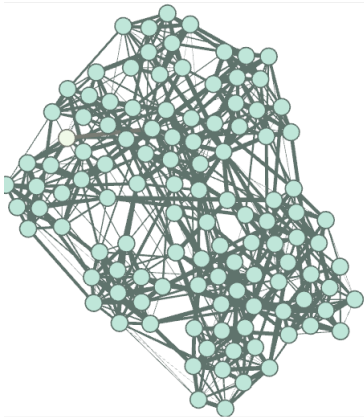


Figure: Original US Football network

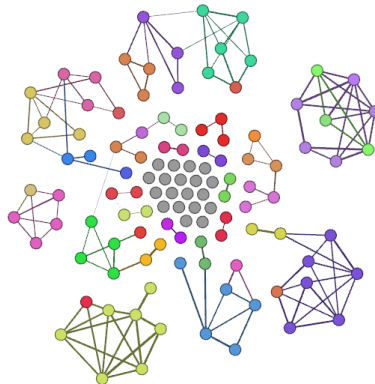


Figure: Sparsified network \mathcal{G}_{comm}

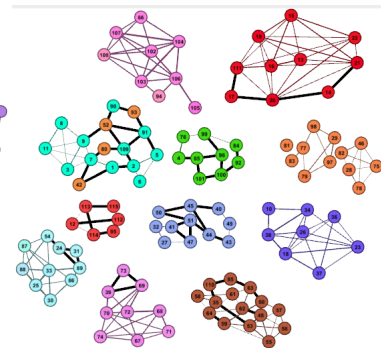
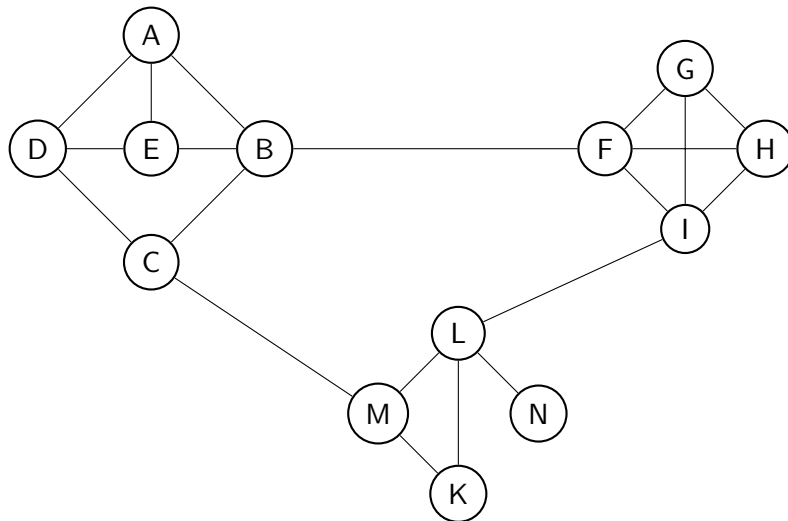


Figure: Final network with communities marked as separate components

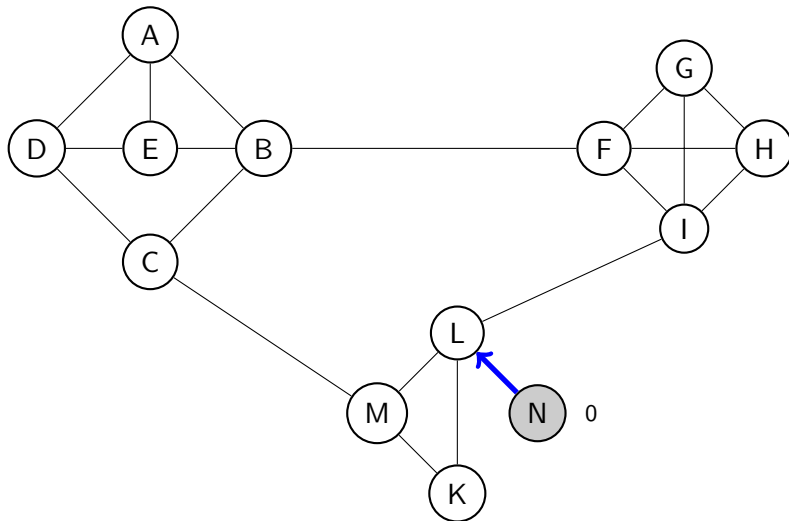
- At every timestamp, a node spreads influence to all of its neighbors.
- Nodes are assigned a *Influenced Neighbors Score*(INS) value at every timestamp

$$INS_t(v) = \frac{\# \text{ neighbors influenced at time } t}{Degree(v)}$$

- A threshold(THR) is set, and the nodes are classified accordingly. If $INS(v) \leq THR$, v is a **broker node**, otherwise v is a **community node**
- Broker nodes are placed in a **stack**, and community nodes in a **queue**. This results in a mix of breadth first and depth first traversals.
- Running time $O(m + n)$

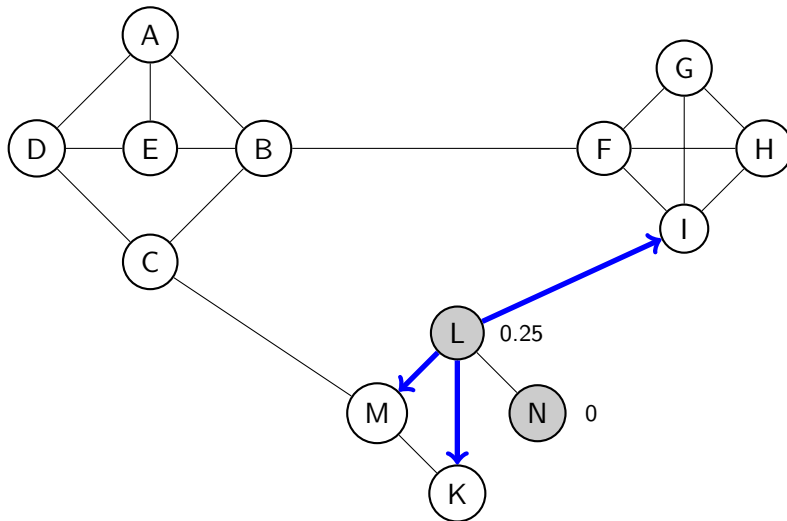


An Example - Threshold is taken to be 0.66.



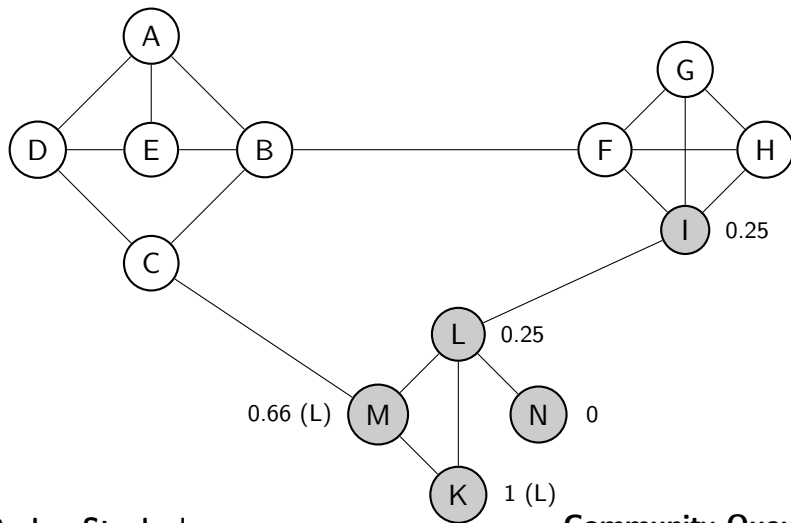
Broker Stack: **N**

Community Queue: \emptyset



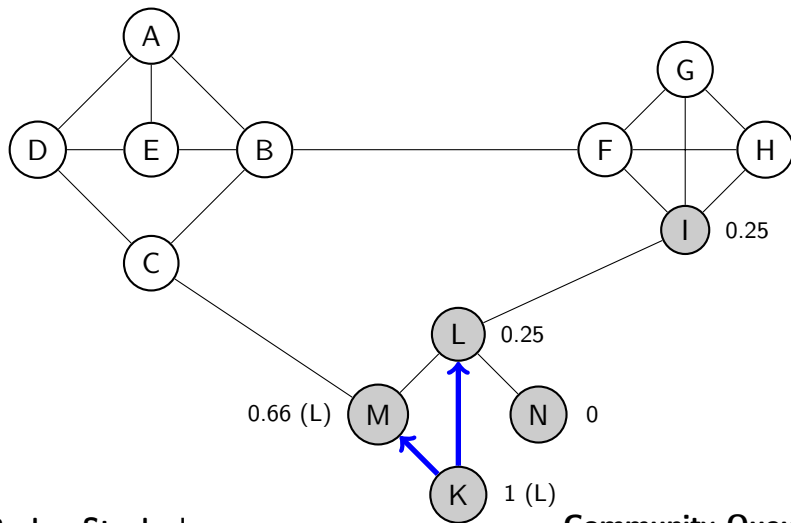
Broker Stack: **L**

Community Queue: \emptyset



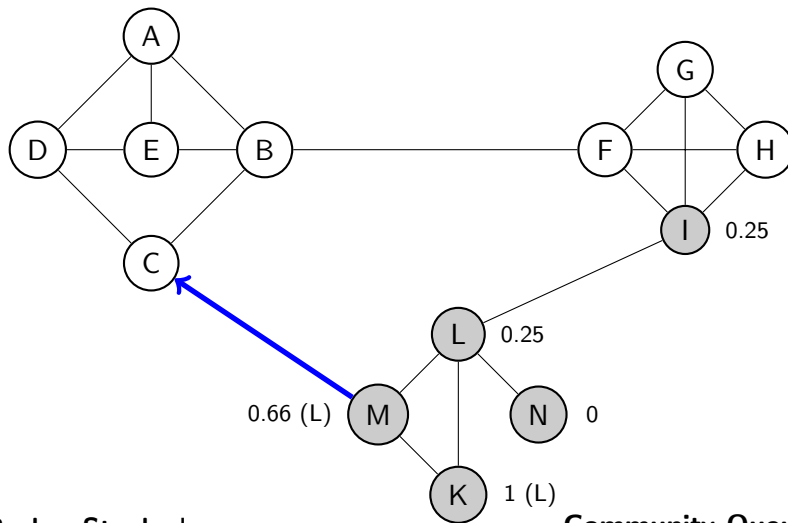
Broker Stack: I

Community Queue: K M



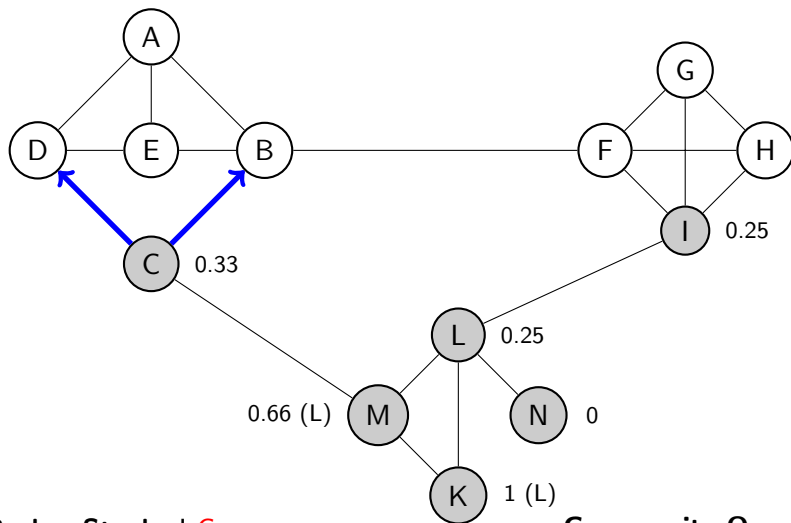
Broker Stack: I

Community Queue: K M



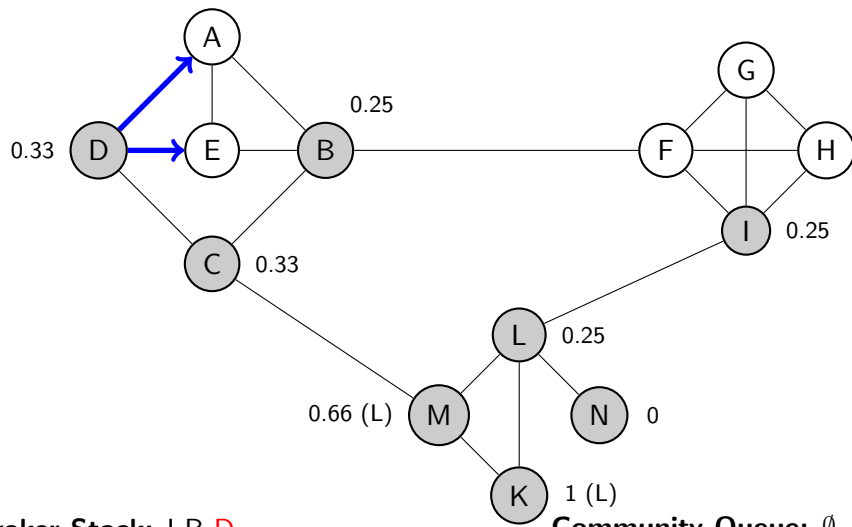
Broker Stack: I

Community Queue: M



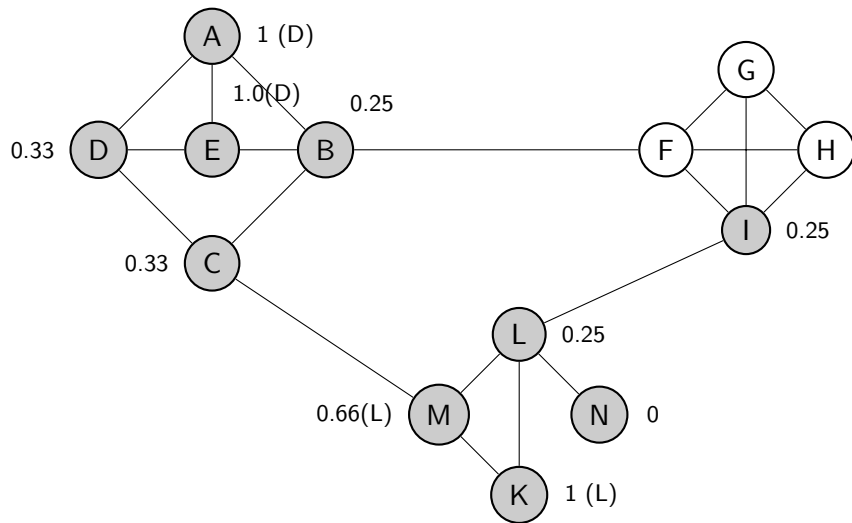
Broker Stack: | C

Community Queue: \emptyset



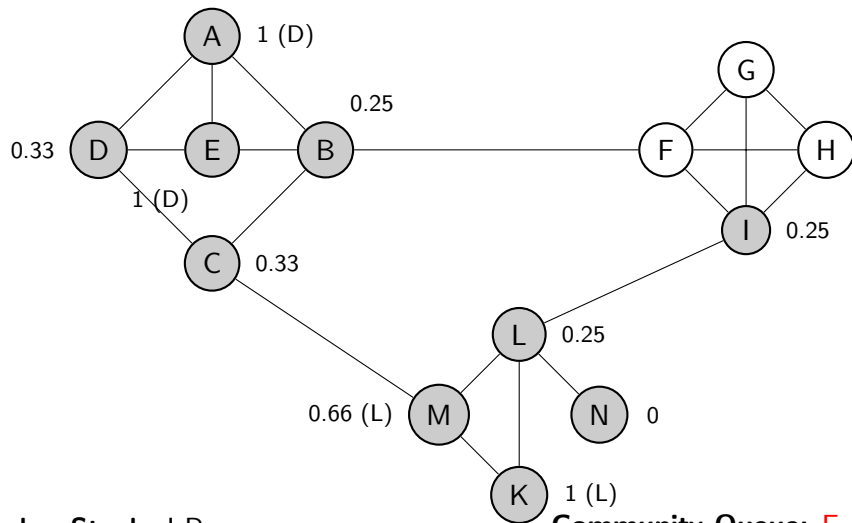
Broker Stack: I B **D**

Community Queue: \emptyset



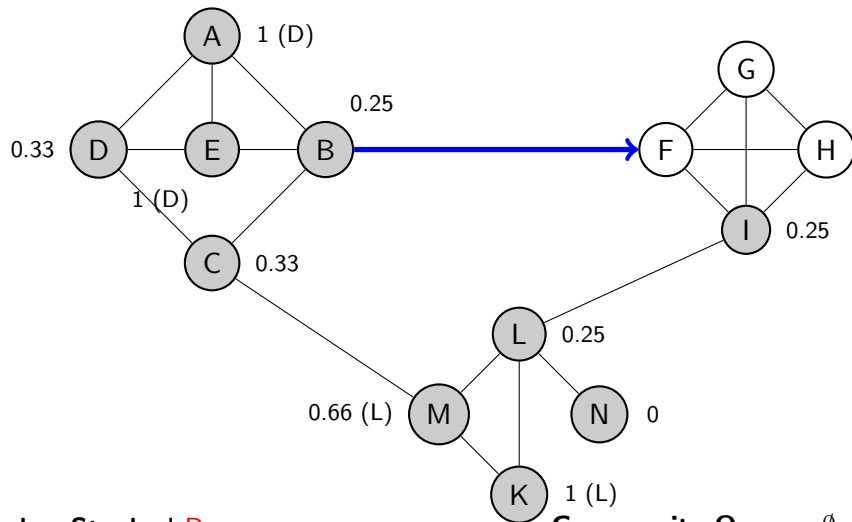
Broker Stack: I B

Community Queue: A E



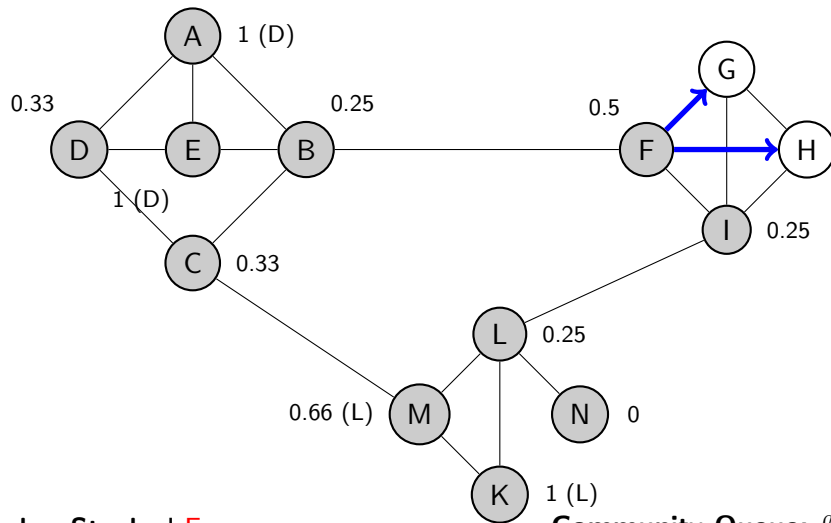
Broker Stack: I B

Community Queue: E



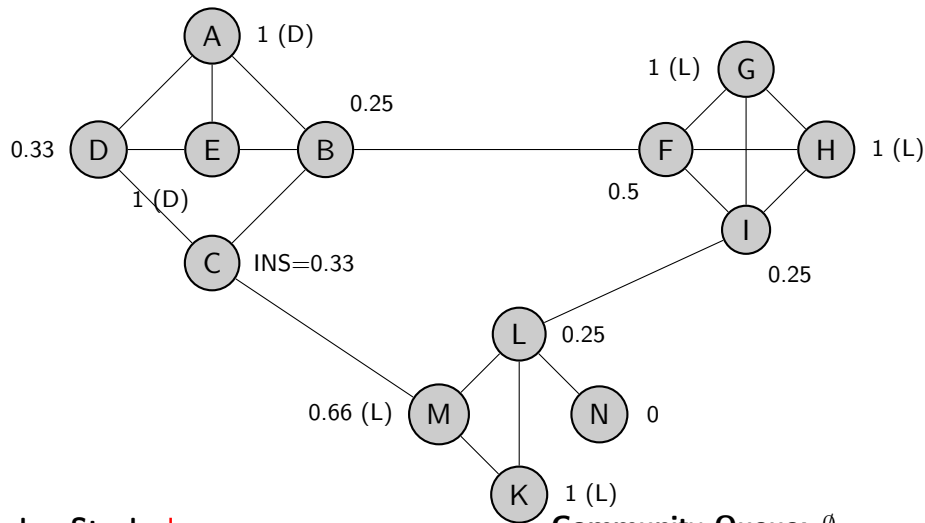
Broker Stack: | B

Community Queue: \emptyset



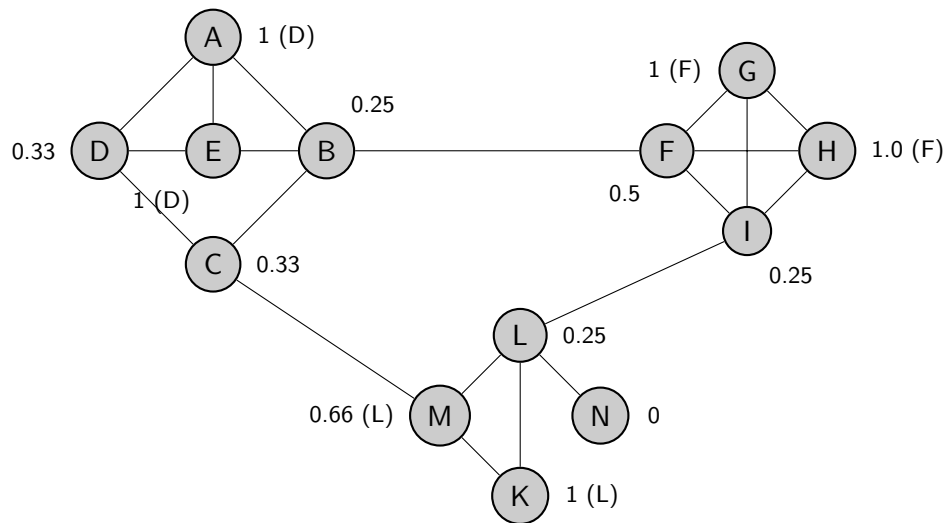
Broker Stack: | F

Community Queue: \emptyset

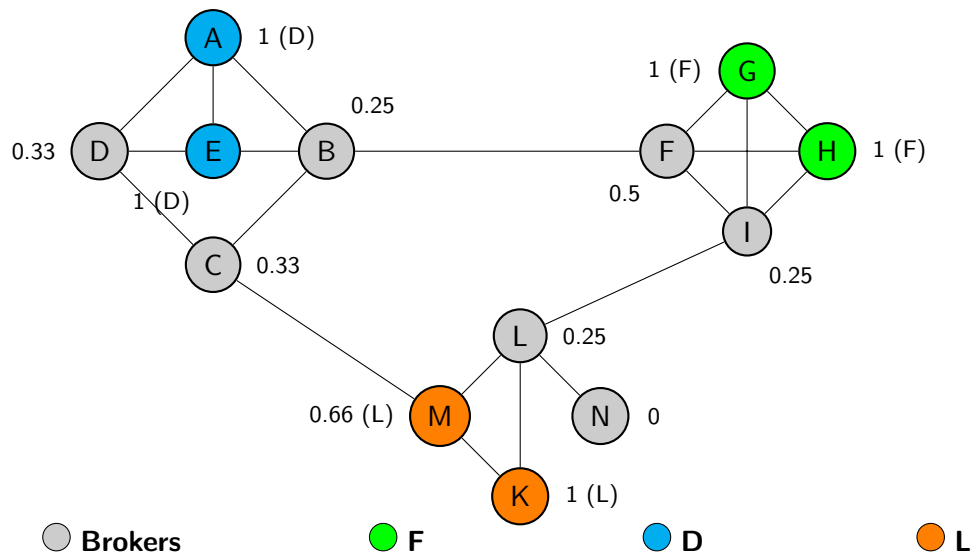


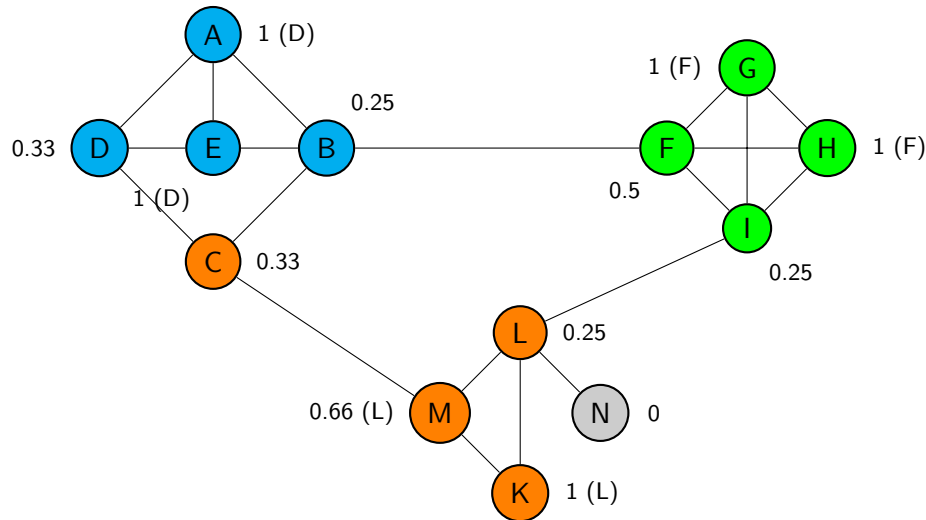
Broker Stack: |

Community Queue: \emptyset

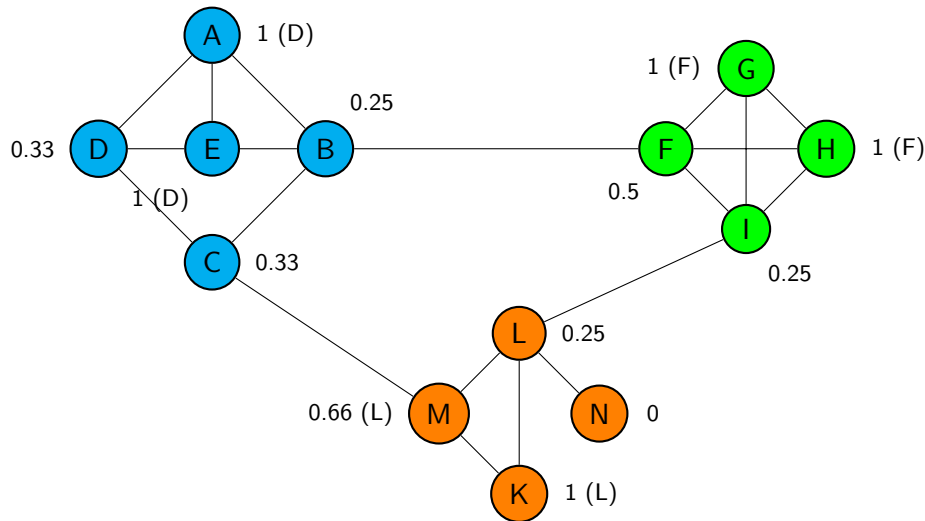


Both the **Broker Stack** and **Community Queue** are empty, so, the process stops.





Post-processing - Brokers are assigned to communities (if possible).



Algorithm 3: Traversal-based Linear Time Community Detection (LINCOT(G, r))

Input : Undirected, unwt. graph $G(V, E)$, threshold (r)
Output: Cover of k communities,
 $G_s = \{G_{s_1}, G_{s_2}, \dots, G_{s_k}\}$

```

1 begin
2   find  $v_{start} \in V$ ,  $\exists v_{start}$  is the node with lowest
   degree
3   forall the  $v \in V$  do
4     nodeType( $v$ )  $\leftarrow$  community( $v$ )  $\leftarrow$  covered( $v$ )  $\leftarrow$ 
       0
5    $S \leftarrow Q \leftarrow \emptyset$   $\triangleright$ brokerStack( $S$ ), communityQueue( $Q$ )
6   coverCount  $\leftarrow$  1
7    $v \leftarrow v_{start}$ 
8   NODE-CAT( $G, v, Q, S$ )
9   while coverCount  $< n$  do
10    if  $Q$  is non-empty then
11       $v \leftarrow$  dequeue( $Q$ )
12    else
13       $v \leftarrow$  pop( $S$ )
14    NODE-CAT( $G, v, Q, S$ )
15   forall the  $v \in V$  do
16     check community( $v$ ) to form  $G_s$ 
17   POST-PROCESS( $G_s$ )
18   CONVERGE( $G, G_s, t$ )
19   return  $G_s$ 

```

Algorithm 4: Categorizing uncovered nodes in $Nghb(v)$ (NODE-CAT(G, v, Q, S))

Input : Undirected, unwt. graph $G(V, E)$, threshold (r), brokerStack (S), communityQueue (Q)
Output: Update Q, S , nodeType list, community list

```

1 begin
2   SPREAD( $v$ )
3   for  $i = 1$  to  $deg(v)$  do
4      $u \leftarrow adjList(v, i)$ 
5     calculate INS( $u$ )
6     if nodeType( $u$ ) is NON-ZERO then continue
7
8     if INS( $u$ )  $< r$  then
9       nodeType( $u$ )  $\leftarrow$  1  $\triangleright$ marking the broker
       nodes
10      push( $S, u$ )
11      community( $u$ )  $\leftarrow$   $u$ 
12    else
13      nodeType( $u$ )  $\leftarrow$  2  $\triangleright$ marking the community
       nodes
14      enqueue( $Q, u$ )
15      community( $u$ )  $\leftarrow$  community( $v$ )
16   return

```

Algorithm 5: Spreading Influence to the Neighbors (SPREAD(v))

Input : Undirected, unwtcd. graph $G(V, E)$, root node (v)

Output: Update coverCount, covered list

```

1 begin
2   for  $i = 1$  to  $\text{deg}(v)$  do
3     if  $\text{covered}(\text{adjList}(v, i)) \text{ EQUALS TO } 0$  then
4       covered( $\text{adjList}(v, i)$ )  $\leftarrow 1$ 
5       coverCount  $\leftarrow$  coverCount + 1
6   return

```

Algorithm 6: Influenced Neighbors Score (INS(v))

Input : Undirected, unwtcd. graph $G(V, E)$, root node (v)

Output: INS(v)

```

1 begin
2   insTotal( $v$ )  $\leftarrow 0$ 
3   for  $i = 1$  to  $\text{deg}(v)$  do
4     if  $\text{covered}(\text{adjList}(v, i)) \text{ EQUALS TO } 1$  then
5       insTotal( $v$ )  $\leftarrow$  insTotal( $v$ ) + 1
6   INS( $v$ )  $\leftarrow$  insTotal( $v$ )/ $\text{deg}(v)$ 
7   return INS( $v$ )

```

Algorithm 7: Post-processing of the Broker Nodes (POST-PROCESS(G_s))

Input : Undirected, unwtcd. graph $G(V, E)$, root node (v)

Output: Update community(v) value for broker nodes

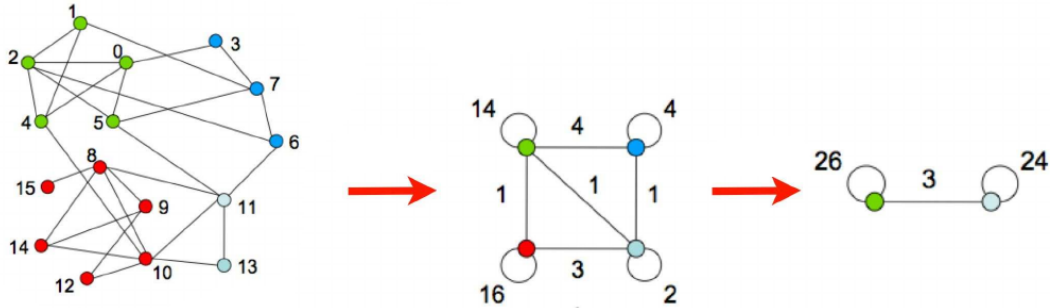
```

1 begin
2   forall the  $v \in V$  do
3      $\text{max}(v) \leftarrow 0$ 
4      $\text{cover}(v) \leftarrow 0$ 
5   forall the  $v \in V$  do
6     if  $\text{nodeType}(v) \text{ EQUALS TO } 1$  then
7       forall the  $G_{s_i} \in G_s$  do
8         if  $\frac{|\text{Neighbors of } v \text{ in } G_{s_i}|}{|G_{s_i}|} > \text{max}(v)$  then
9            $\text{max}(v) \leftarrow \frac{|\text{Neighbors of } v \text{ in } G_{s_i}|}{|G_{s_i}|}$ 
10        if  $\text{max}(v)$  is NON-ZERO then
11          community( $v$ )  $\leftarrow$  cover that leads to  $\text{max}(v)$ 
12   return

```



Modularity Maximization[3]





Modularity

We define modularity [2] as

$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j) \quad Q \in [-1, 1]$$

- The higher the modularity, the better is the community structure*
- Actual social networks have modularity values between 0.30 - 0.60
- Suffers from resolution limit problems, but usually works very well in practice



Performance Comparison

Dataset Information [11]			Louvain Method		Our Algorithm (3-spanner)	
Name	n	m	Modularity	Time(sec)	Modularity	Time(sec)
Karate	34	78	0.415	0	0.589	0.51
Dolphins	62	159	0.518	0	0.676	0.53
Football	115	613	0.604	0	0.8615	0.65
Enron	33,696	180,811	0.596	0.38	0.855	13.13
Epinions	75,877	405,739	0.45	0.97	0.695	27.03
Amazon	334,863	925,872	0.926	6	0.995	78.47
DBLP	317,080	1,049,866	0.819	11	0.961	76.8607



Conclusions and Future Scope

- Jaccard coefficient based edge weight preserves the community structure
 - Spanners can be used to identify inter-community edges (to a large extent)
 - LINCOM is effective in breaking the graph into pure fragments
 - Works very fast in practice - produces good quality clusters
-
- The stretch-factor t and THR is predetermined, may turn out to be network dependent
 - The algorithm can be modified to do overlapping community detection



References I

- [1] Markus M Möbius and Tanya S Rosenblat.
The process of ghetto formation: Evidence from chicago.
Unpublished manuscript. <http://trosenblat.nber.org/papers/Files/Chicago/chicago-dec7.pdf>, 2001.
- [2] M. Girvan and M. E. J. Newman.
Community structure in social and biological networks.
PNAS, 99(12):7821–7826, June 2002.
- [3] V.D. Blondel, J.L. Guillaume, R. Lambiotte, and E.L.J.S. Mech.
Fast unfolding of communities in large networks.
J. Stat. Mech, page P10008, 2008.
- [4] Venu Satuluri, Srinivasan Parthasarathy, and Yiye Ruan.
Local graph sparsification for scalable clustering.
In Proceedings of the 2011 ACM SIGMOD International Conference on Management of data, pages 721–732. ACM, 2011.



References II

- [5] David Liben-Nowell and Jon Kleinberg.
The link-prediction problem for social networks.
Journal of the American society for information science and technology, 58(7):1019–1031, 2007.
- [6] Gerd Lindner, Christian L Staudt, Michael Hamann, Henning Meyerhenke, and Dorothea Wagner.
Structure-preserving sparsification of social networks.
In *Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015*, pages 448–454. ACM, 2015.
- [7] Ingo Althöfer, Gautam Das, David Dobkin, Deborah Joseph, and José Soares.
On sparse spanners of weighted graphs.
Discrete & Computational Geometry, 9(1):81–100, 1993.
- [8] Liam Roditty and Uri Zwick.
On dynamic shortest paths problems.
In *Algorithms–ESA 2004*, pages 580–591. Springer, 2004.



References III

- [9] Liam Roditty, Mikkel Thorup, and Uri Zwick.
Deterministic constructions of approximate distance oracles and spanners.
In *Automata, languages and programming*, pages 261–272. Springer, 2005.
- [10] Surender Baswana and Sandeep Sen.
A simple and linear time randomized algorithm for computing sparse spanners in weighted graphs.
Random Structures & Algorithms, 30(4):532–563, 2007.
- [11] Jure Leskovec and Andrej Krevl.
SNAP Datasets: Stanford large network dataset collection.
<http://snap.stanford.edu/data>, June 2014.

"No one ever complains about a speech being too short!" - Ira Hayes

Thank you for your attention. Any questions?