

Overlapping Community Detection Using Seed Set Expansion

Joyce Jiyoung Whang¹ David F. Gleich² Inderjit S. Dhillon¹

¹The University of Texas at Austin

²Purdue University

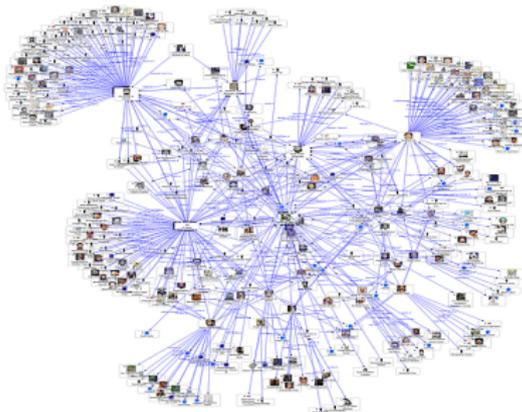
International Conference on Information and Knowledge Management
Oct. 27th - Nov. 1st, 2013.

Contents

- Introduction
 - Overlapping Communities in Real-world Networks
 - Measures of Cluster Quality
 - Graph Clustering and Weighted Kernel k -Means
- The Proposed Algorithm
 - Filtering Phase
 - Seeding Phase
 - Seed Set Expansion Phase
 - Propagation Phase
- Experimental Results
 - Conductance
 - Ground-truth Accuracy
 - Runtime
- Conclusions

Overlapping Communities

- Community (cluster) in a graph $G = (\mathcal{V}, \mathcal{E})$
 - Set of cohesive vertices
 - Communities naturally overlap (e.g. social circles)
- Graph Clustering (Partitioning)
 - k disjoint clusters $\mathcal{C}_1, \dots, \mathcal{C}_k$ such that $\mathcal{V} = \mathcal{C}_1 \cup \dots \cup \mathcal{C}_k$
- Overlapping Community Detection
 - k overlapping clusters such that $\mathcal{C}_1 \cup \dots \cup \mathcal{C}_k \subseteq \mathcal{V}$



Real-world Networks

- Collaboration networks: co-authorship
- Social networks: friendship
- Product network: co-purchasing information

Graph	No. of vertices	No. of edges
<i>Collaboration networks</i>		
HepPh	11,204	117,619
AstroPh	17,903	196,972
CondMat	21,363	91,286
DBLP	317,080	1,049,866
<i>Social networks</i>		
Flickr	1,994,422	21,445,057
Myspace	2,086,141	45,459,079
LiveJournal	1,757,326	42,183,338
<i>Product network</i>		
Amazon	334,863	925,872

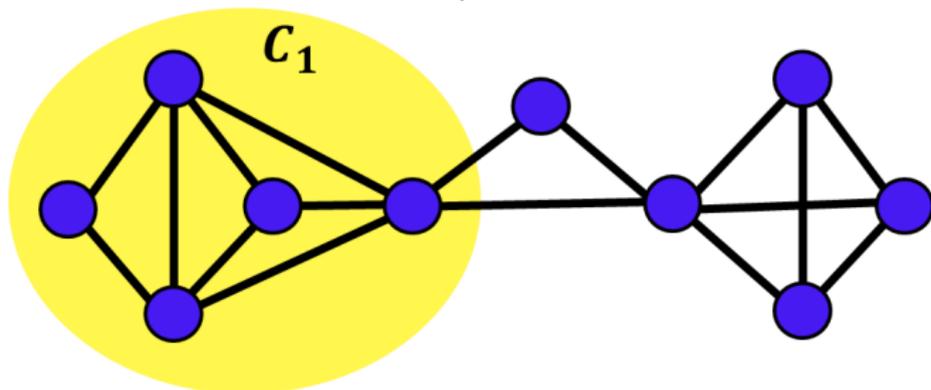
Measures of cluster quality

- Normalized Cut of a cluster

$$\text{ncut}(\mathcal{C}_i) = \frac{\text{links}(\mathcal{C}_i, \mathcal{V} \setminus \mathcal{C}_i)}{\text{links}(\mathcal{C}_i, \mathcal{V})}.$$

- Conductance

$$\text{conductance}(\mathcal{C}_i) = \frac{\text{links}(\mathcal{C}_i, \mathcal{V} \setminus \mathcal{C}_i)}{\min \left(\text{links}(\mathcal{C}_i, \mathcal{V}), \text{links}(\mathcal{V} \setminus \mathcal{C}_i, \mathcal{V}) \right)}.$$



$$\text{links}(\mathcal{C}_1, \mathcal{V} \setminus \mathcal{C}_1) = 2, \text{links}(\mathcal{C}_1, \mathcal{V}) = 10, \text{links}(\mathcal{V} \setminus \mathcal{C}_1, \mathcal{V}) = 9$$

Graph Clustering and Weighted Kernel k -Means

- **A general weighted kernel k -means objective** is equivalent to a **weighted graph clustering objective** (Dhillon et al. 2007).
- Weighted kernel k -means
 - Objective

$$J = \sum_{c=1}^k \sum_{\mathbf{x}_i \in \pi_c} w_i \|\varphi(\mathbf{x}_i) - \mathbf{m}_c\|^2, \text{ where } \mathbf{m}_c = \frac{\sum_{\mathbf{x}_i \in \pi_c} w_i \varphi(\mathbf{x}_i)}{\sum_{\mathbf{x}_i \in \pi_c} w_i}.$$

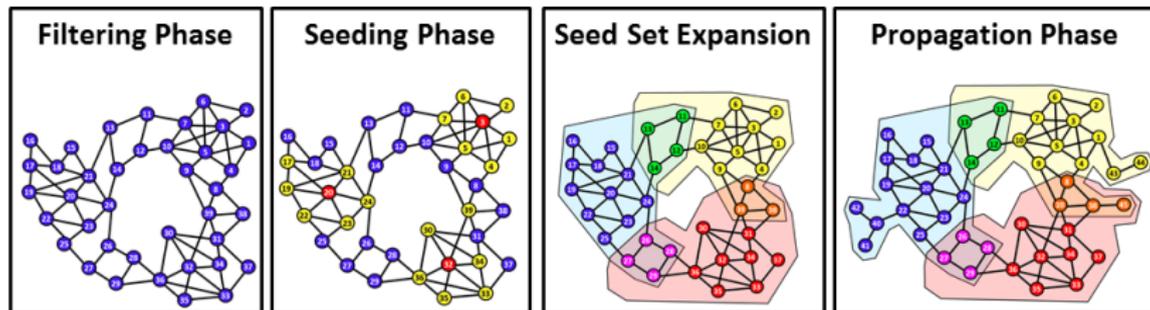
- Distance between a vertex $v \in \mathcal{C}_i$ and cluster \mathcal{C}_i

$$\text{dist}(v, \mathcal{C}_i) = -\frac{2 \text{links}(v, \mathcal{C}_i)}{\text{deg}(v) \text{deg}(\mathcal{C}_i)} + \frac{\text{links}(\mathcal{C}_i, \mathcal{C}_i)}{\text{deg}(\mathcal{C}_i)^2} + \frac{\sigma}{\text{deg}(v)} - \frac{\sigma}{\text{deg}(\mathcal{C}_i)}$$

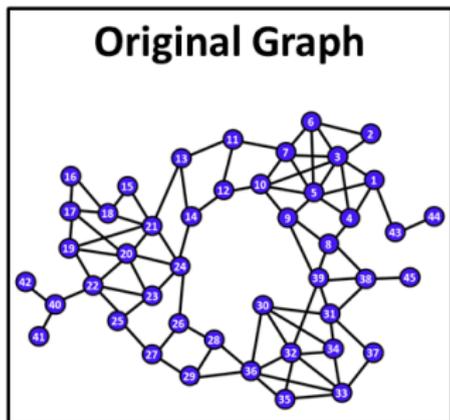
The Proposed Algorithm

Proposed Algorithm

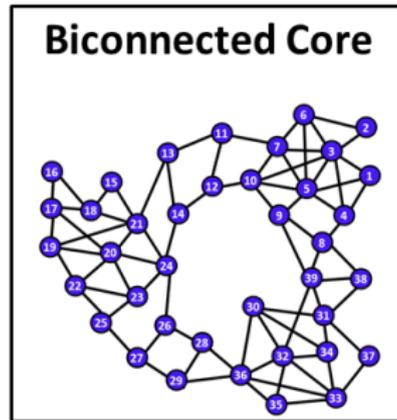
- Seed Set Expansion
 - Carefully select seeds
 - Greedily expand communities around the seed sets
- The algorithm
 - Filtering Phase
 - Seeding Phase
 - Seed Set Expansion Phase
 - Propagation Phase



Filtering Phase

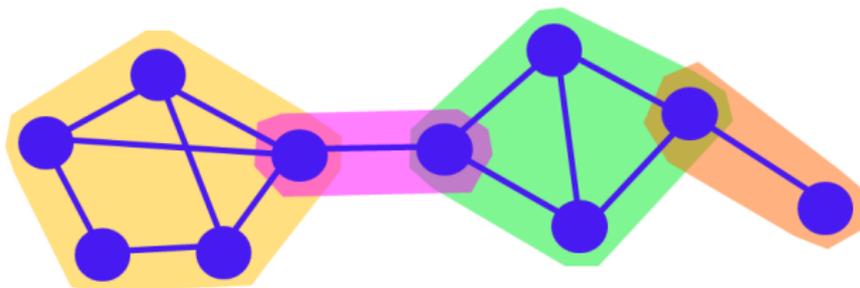


Filtering

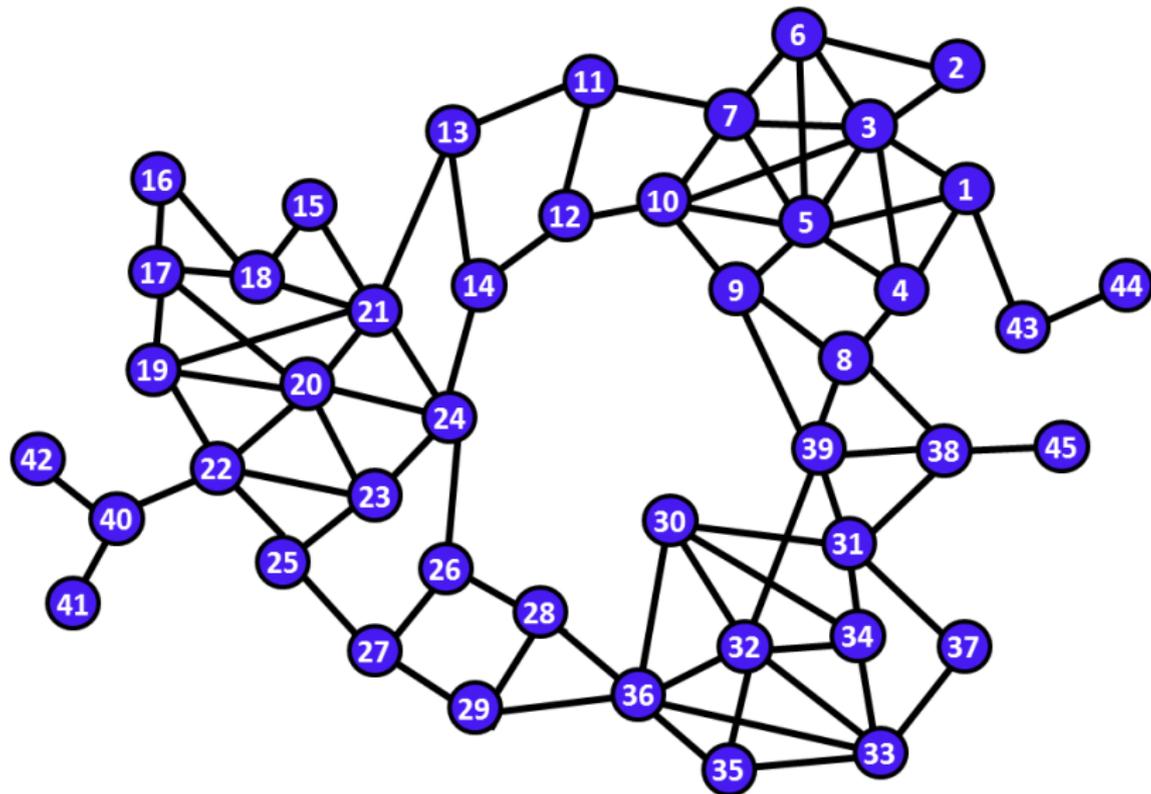


Filtering Phase

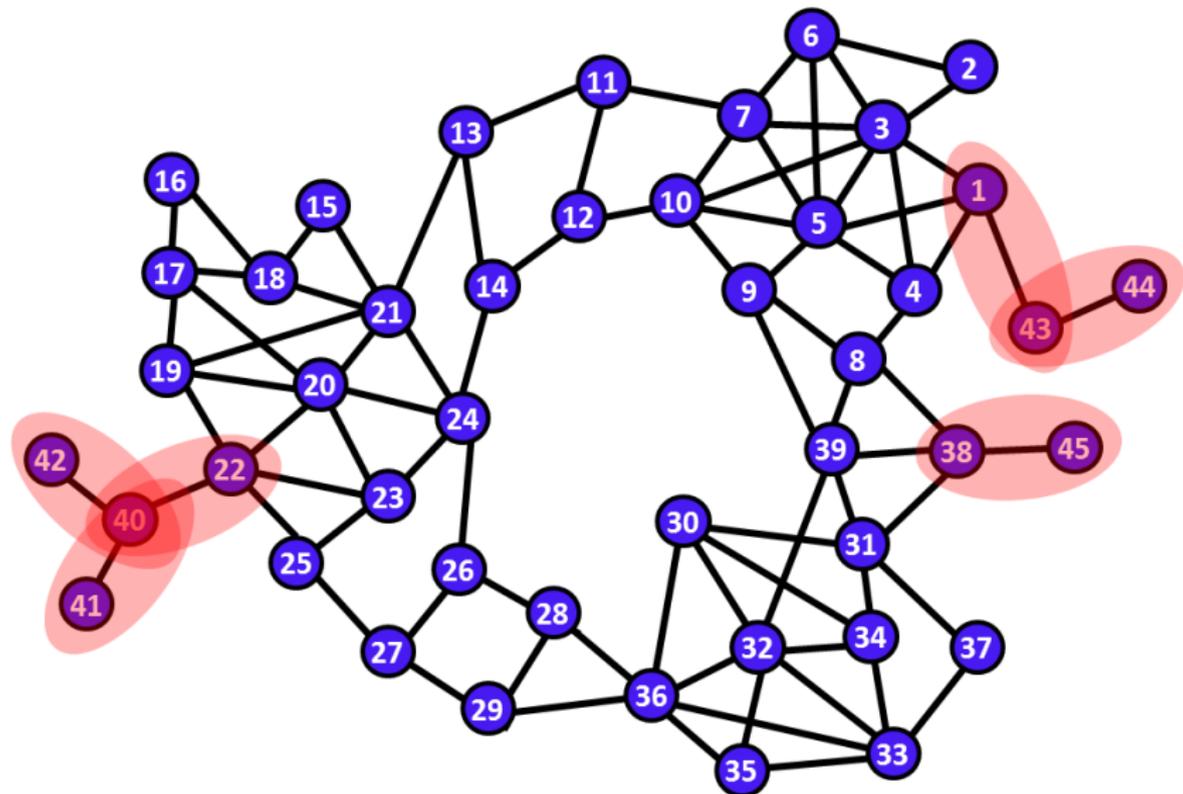
- Remove unimportant regions of the graph
 - Trivially separable from the rest of the graph
 - Do not participate in overlapping clustering
- Our filtering procedure
 - Remove all single-edge biconnected components (remain connected after removing any vertex and its adjacent edges)
 - Compute the largest connected component (LCC)



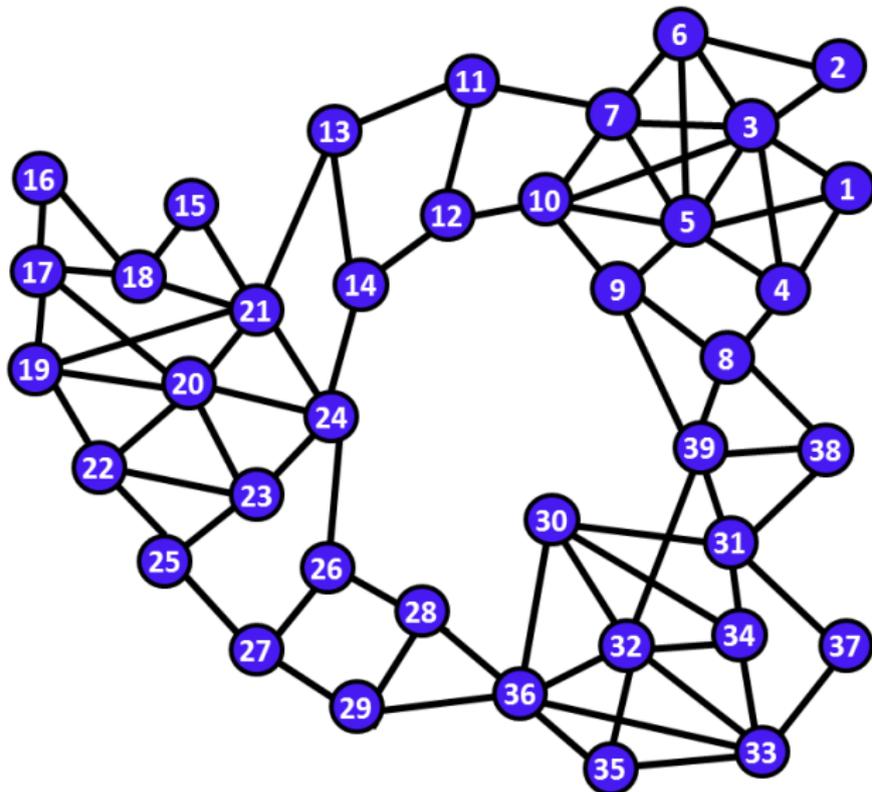
Filtering Phase



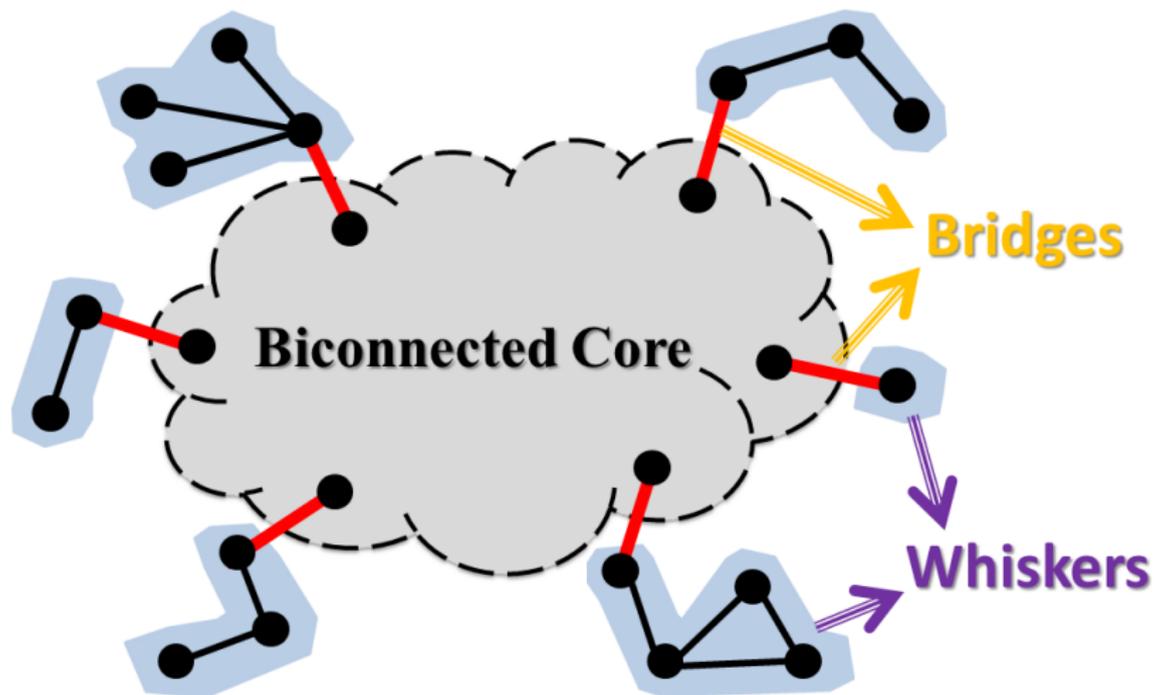
Filtering Phase



Filtering Phase



Filtering Phase

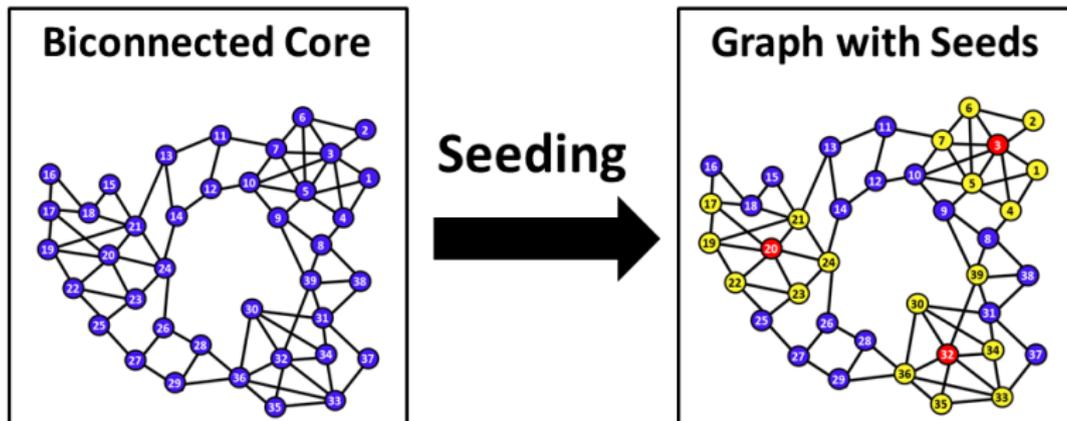


Filtering Phase

	Biconnected core		Detached graph	
	No. of vertices (%)	No. of edges (%)	No. of components	Size of LCC (%)
HepPh	9,945 (88.8%)	116,099 (98.7%)	1,123	21 (0.0019%)
AstroPh	16,829 (94.0%)	195,835 (99.4%)	957	23 (0.0013%)
CondMat	19,378 (90.7%)	89,128 (97.6%)	1,669	12 (0.00056%)
DBLP	264,341 (83.4%)	991,125 (94.4%)	43,093	32 (0.00010%)
Flickr	954,672 (47.9%)	20,390,649 (95.1%)	864,628	107 (0.000054%)
Myspace	1,724,184 (82.7%)	45,096,696 (99.2%)	332,596	32 (0.000015%)
LiveJournal	1,650,851 (93.9%)	42,071,541 (99.7%)	101,038	105 (0.000060%)
Amazon	291,449 (87.0%)	862,836 (93.2%)	25,835	250 (0.00075%)

- The biconnected core – substantial portion of the edges
- Detached graph – likely to be disconnected
- Whiskers – separable from each other, no significant size

Seeding Phase



Seeding Phase

- Graclus centers
 - Graclus: a high quality and efficient graph partitioning scheme

Algorithm 1 Seeding by Graclus Centers

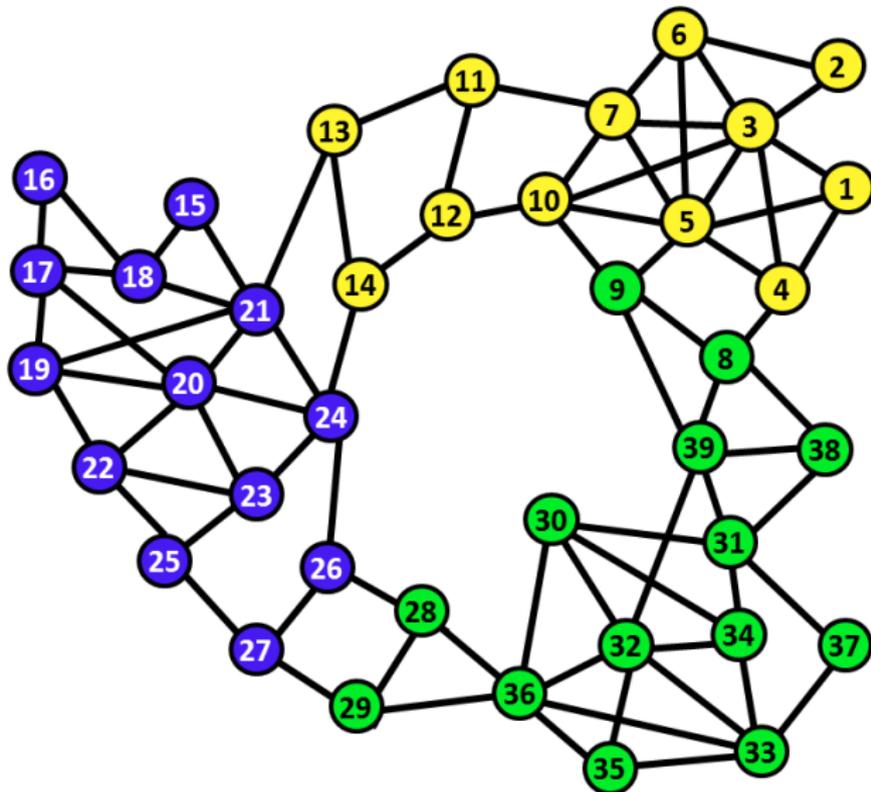
Input: graph G , the number of seeds k .

Output: the seed set \mathcal{S} .

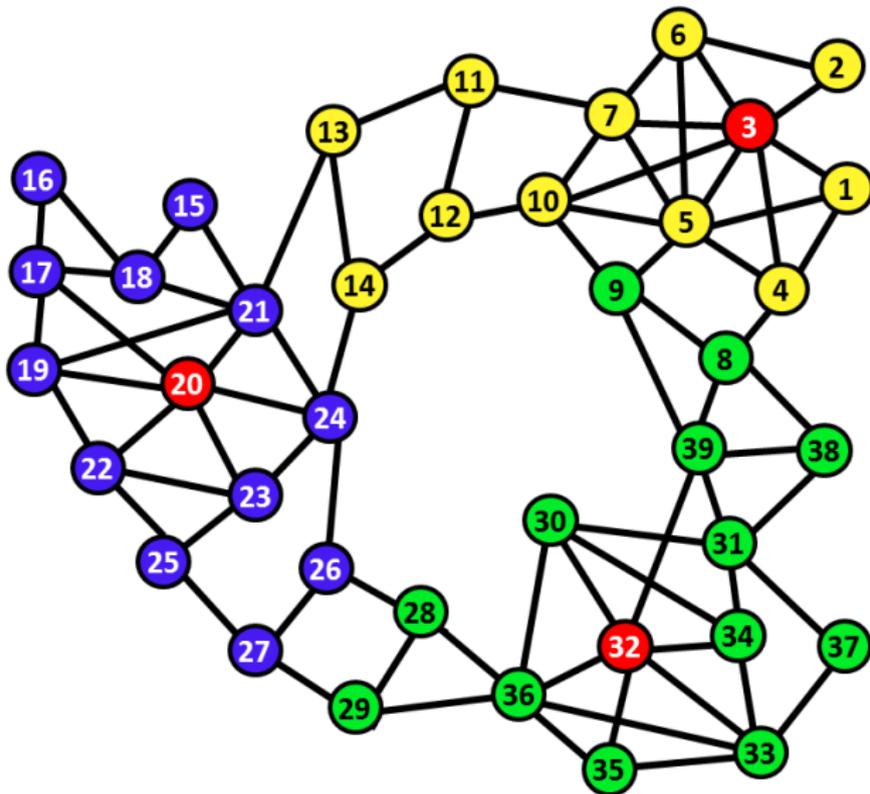
- 1: Compute exhaustive and non-overlapping clusters \mathcal{C}_i ($i=1, \dots, k$) on G .
- 2: Initialize $\mathcal{S} = \emptyset$.
- 3: **for** each cluster \mathcal{C}_i **do**
- 4: **for** each vertex $v \in \mathcal{C}_i$ **do**
- 5: Compute $\text{dist}(v, \mathcal{C}_i)$.
- 6: **end for**
- 7: $\mathcal{S} = \{\underset{v}{\text{argmin}} \text{dist}(v, \mathcal{C}_i)\} \cup \mathcal{S}$.
- 8: **end for**

Find the most central vertex
in cluster \mathcal{C}_i

Seeding Phase



Seeding Phase



Seeding Phase

- Spread Hubs
 - Independent set of high-degree vertices

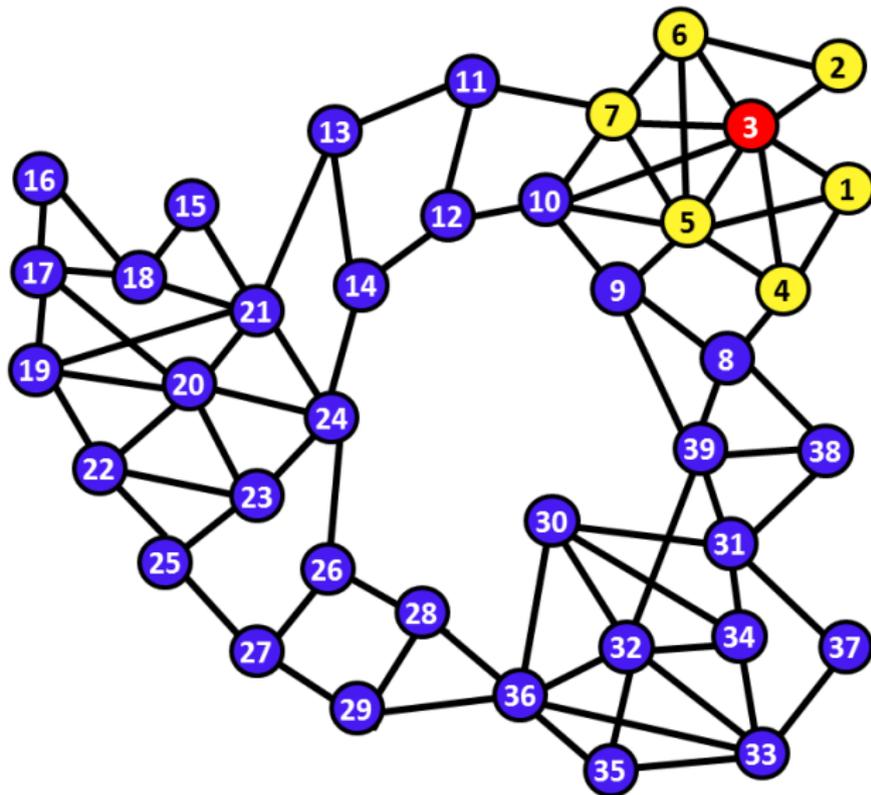
Algorithm 1 Seeding by Spread Hubs

Input: graph $G = (\mathcal{V}, \mathcal{E})$, the number of seeds k .

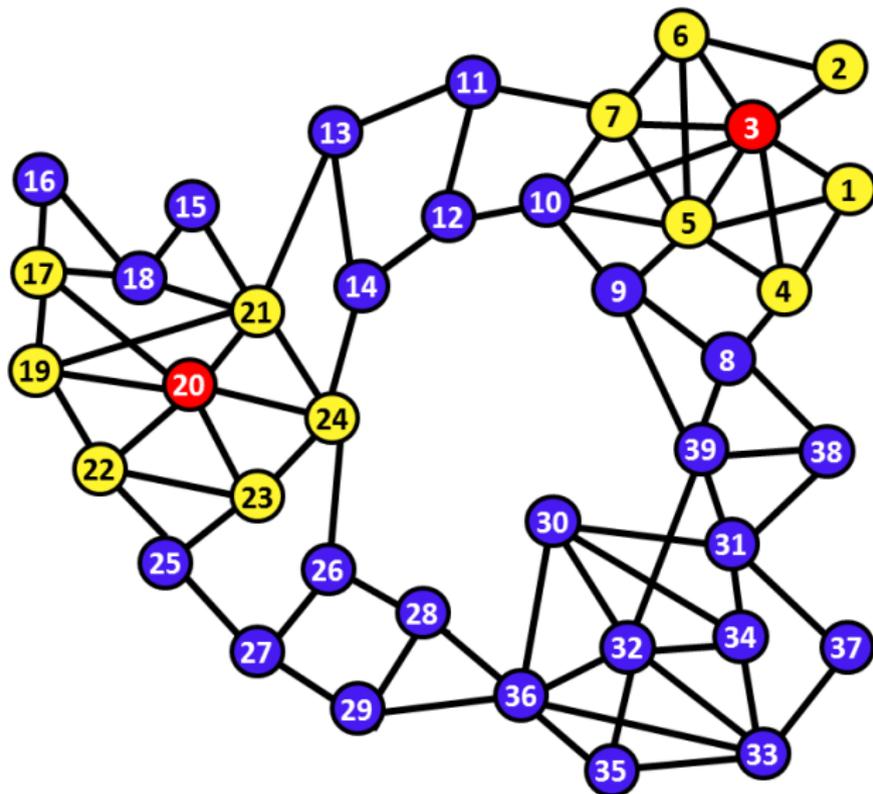
Output: the seed set \mathcal{S} .

- 1: Initialize $\mathcal{S} = \emptyset$.
 - 2: All vertices in \mathcal{V} are unmarked.
 - 3: **while** $|\mathcal{S}| < k$ **do**
 - 4: Let \mathcal{T} be the set of unmarked vertices with max degree.
 - 5: **for each** $t \in \mathcal{T}$ **do**
 - 6: **if** t is unmarked **then**
 - 7: $\mathcal{S} = \{t\} \cup \mathcal{S}$.
 - 8: Mark t and its neighbors.
 - 9: **end if**
 - 10: **end for**
 - 11: **end while**
-

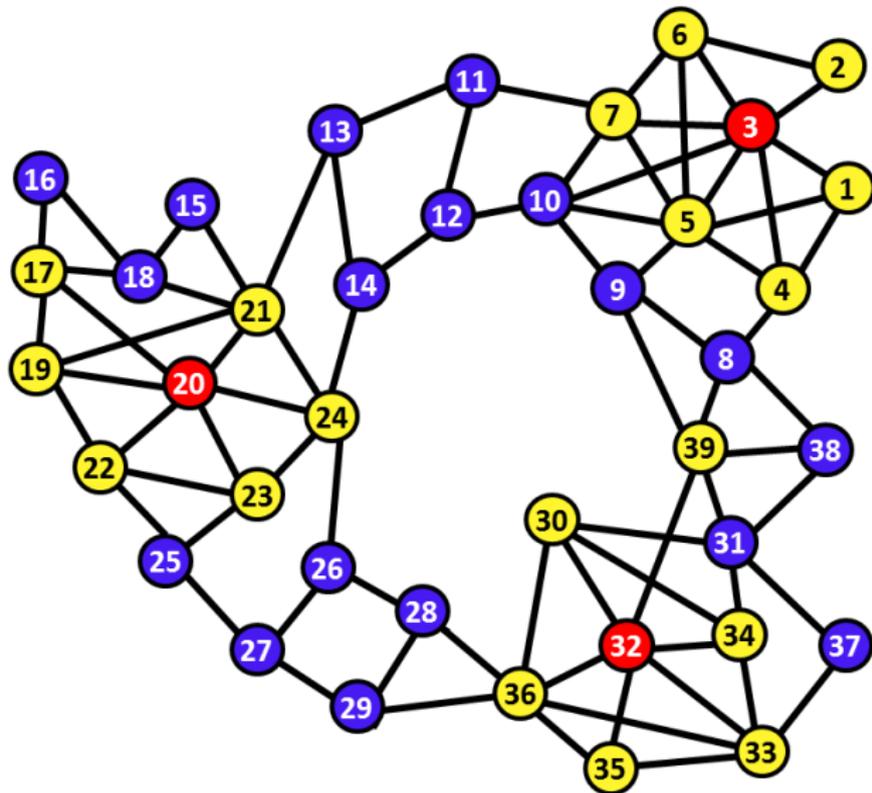
Seeding Phase



Seeding Phase



Seeding Phase



Seeding Phase

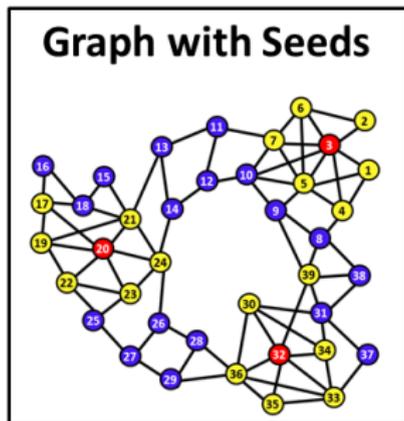
- Other seeding strategies
 - **Local Optimal Egonets.** (Gleich and Seshadhri 2012)
 - $\text{ego}(s)$: the egonet of vertex s .
 - Select a seed s such that

$$\text{conductance}(\text{ego}(s)) \leq \text{conductance}(\text{ego}(v))$$

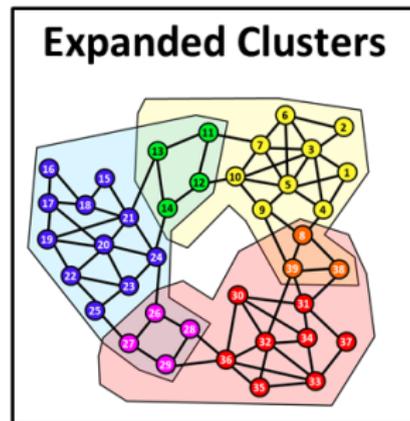
for all v adjacent to s .

- **Random Seeds.** (Andersen and Lang 2006)
 - Randomly select k seeds.

Seed Set Expansion Phase



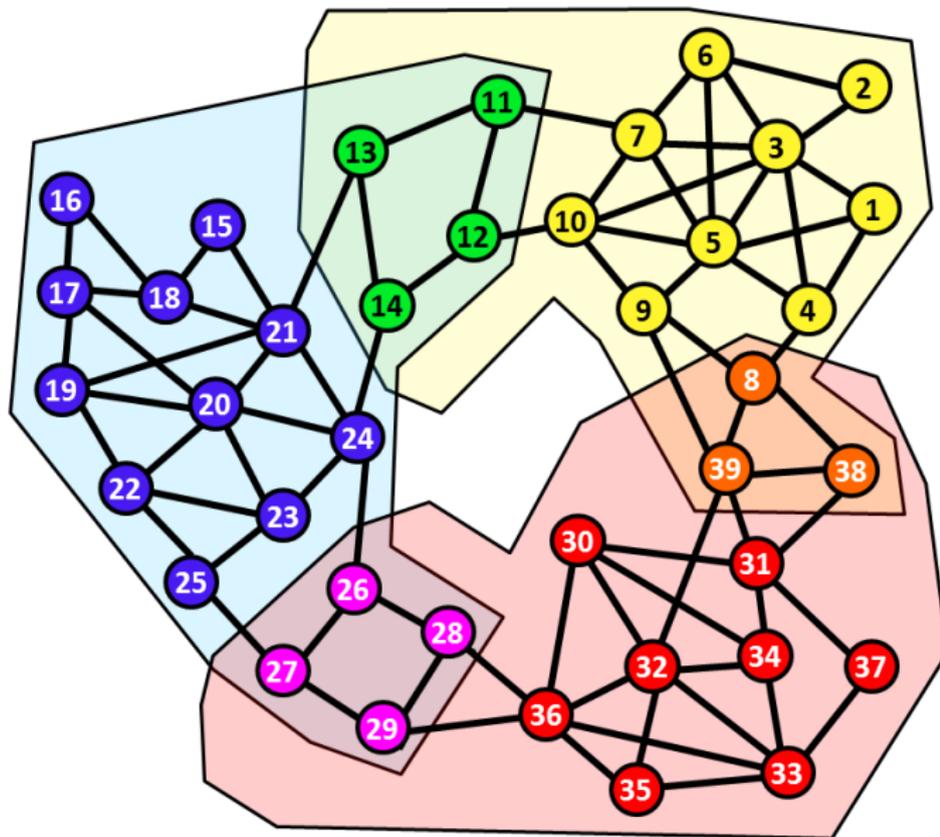
Expansion



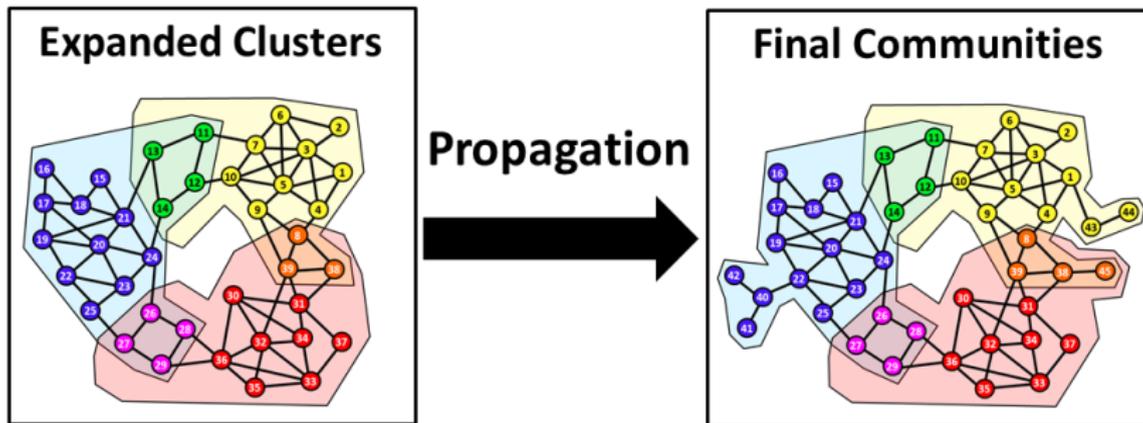
Seed Set Expansion Phase

- Personalized PageRank clustering scheme (Andersen et al. 2006)
 - 1 Given a seed node, compute an approximation of the stationary distribution of a random walk.
 - 2 Divide the stationary distribution scores by the degree of each node (technical detail needed to remove bias towards high-degree nodes).
 - 3 Sort the vector, and examine nodes in order of highest to lowest score and compute the conductance score for each threshold cut.
- Returns a good conductance cluster
- Remarkably efficient when combined with appropriate data structures
- For each seed, we use the entire vertex neighborhood as the restart for the personalized PageRank routine.

Seed Set Expansion Phase



Propagation Phase



Propagation Phase

- Each community is further expanded.
- Add whiskers to communities via bridge.

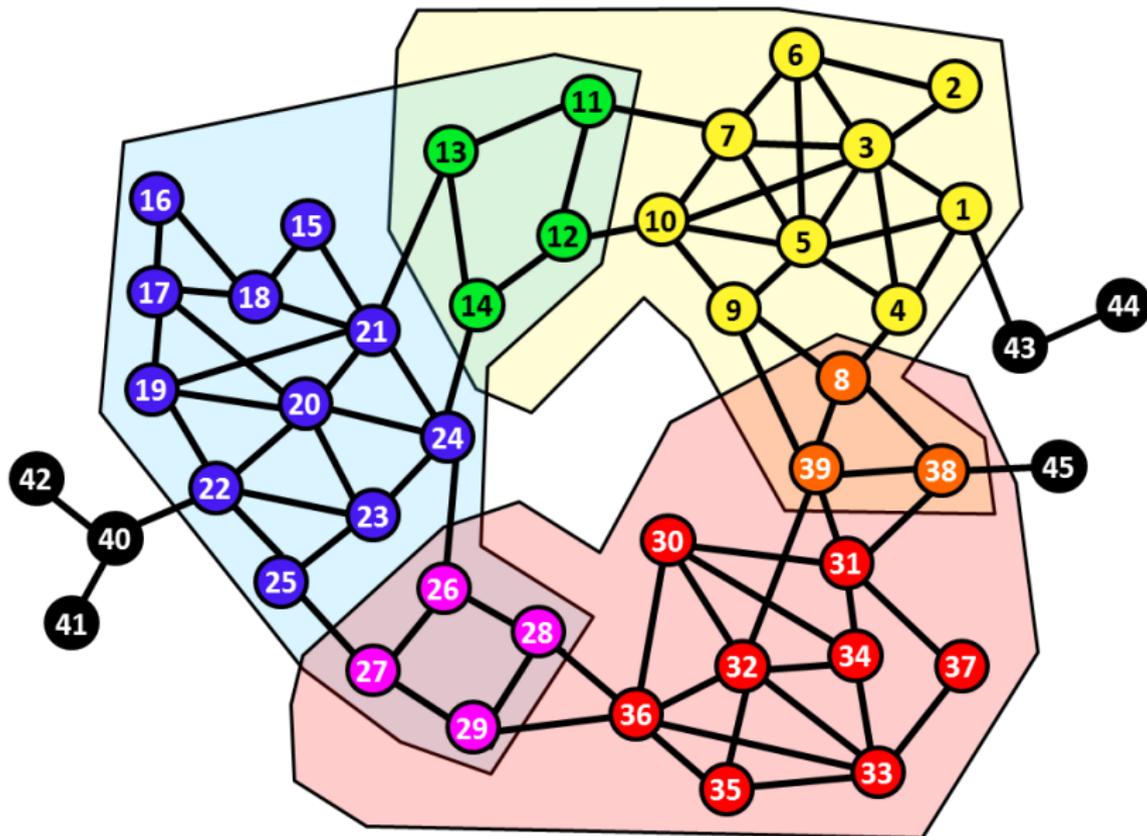
Algorithm 2 Propagation Module

Input: graph $G = (\mathcal{V}, \mathcal{E})$, biconnected core $G_C = (\mathcal{V}_C, \mathcal{E}_C)$, communities of $G_C : \mathcal{C}_i$ ($i = 1, \dots, k$) $\in \mathcal{C}$.

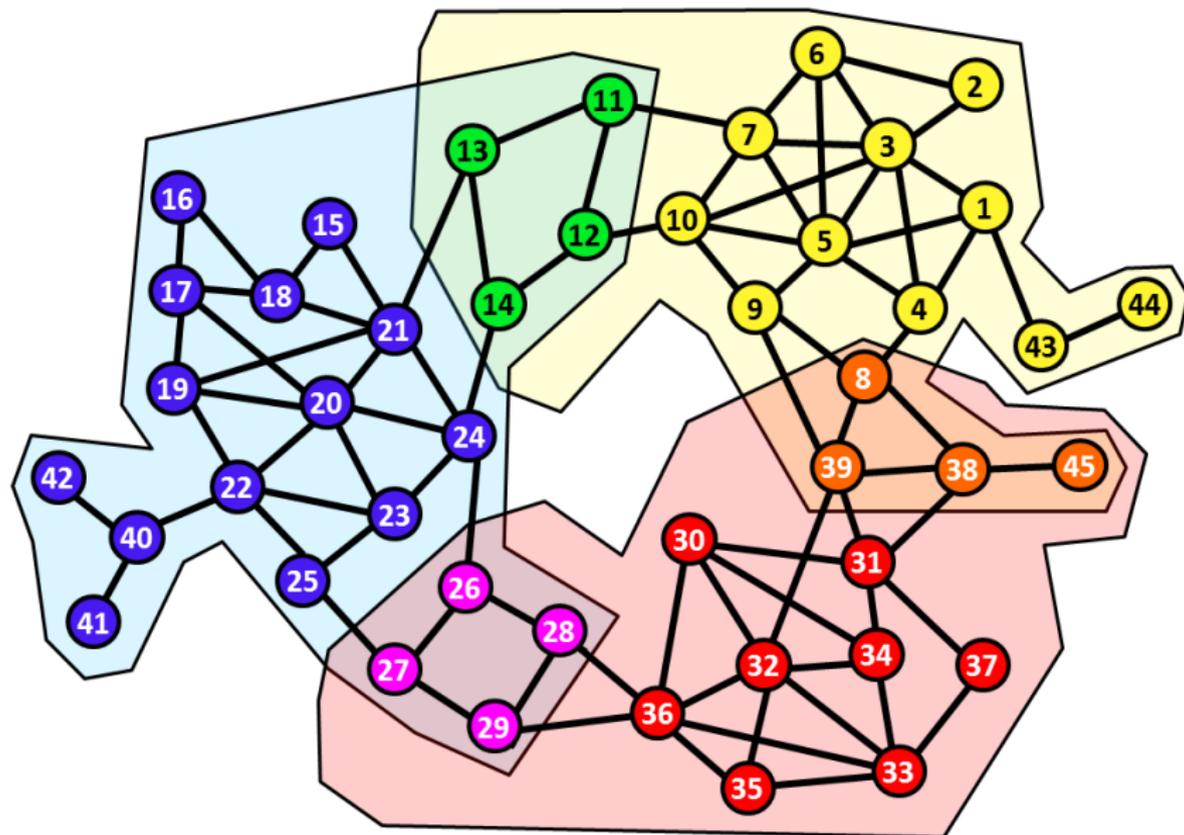
Output: communities of G .

- 1: **for** each $\mathcal{C}_i \in \mathcal{C}$ **do**
 - 2: Detect bridges \mathcal{E}_{B_i} attached to \mathcal{C}_i .
 - 3: **for** each $b_j \in \mathcal{E}_{B_i}$ **do**
 - 4: Detect the whisker $w_j = (\mathcal{V}_j, \mathcal{E}_j)$ which is attached to b_j .
 - 5: $\mathcal{C}_i = \mathcal{C}_i \cup \mathcal{V}_j$.
 - 6: **end for**
 - 7: **end for**
-

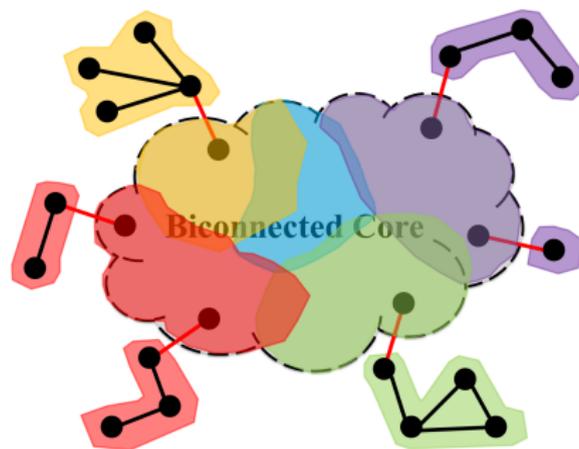
Propagation Phase



Propagation Phase



Propagation Phase



- This process does not increase the cut of each cluster.
- Normalized cut of the expanded cluster is always smaller than equal to that of original cluster.

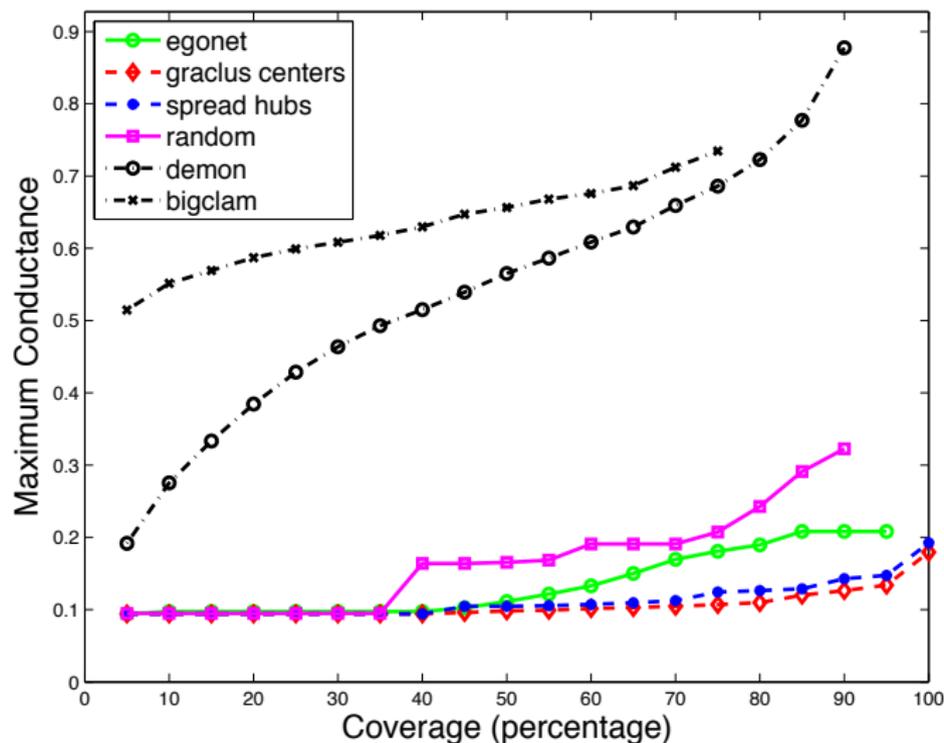
Experimental Results

Experiments

- Comparison with other state-of-the-art methods
 - **Demon** (Coscia et al. 2012)
 - Extracts and computes clustering of ego networks
 - **Bigclam** (Yang and Leskovec 2013)
 - Low-rank non-negative matrix factorization based modeling
- Seed set expansion methods with different seeding strategies
 - **Graclus centers**
 - **Spread hubs**
 - **Local Optimal Egonets** (Gleich and Seshadhri 2012)
 - **Random Seeds** (Andersen and Lang 2006)

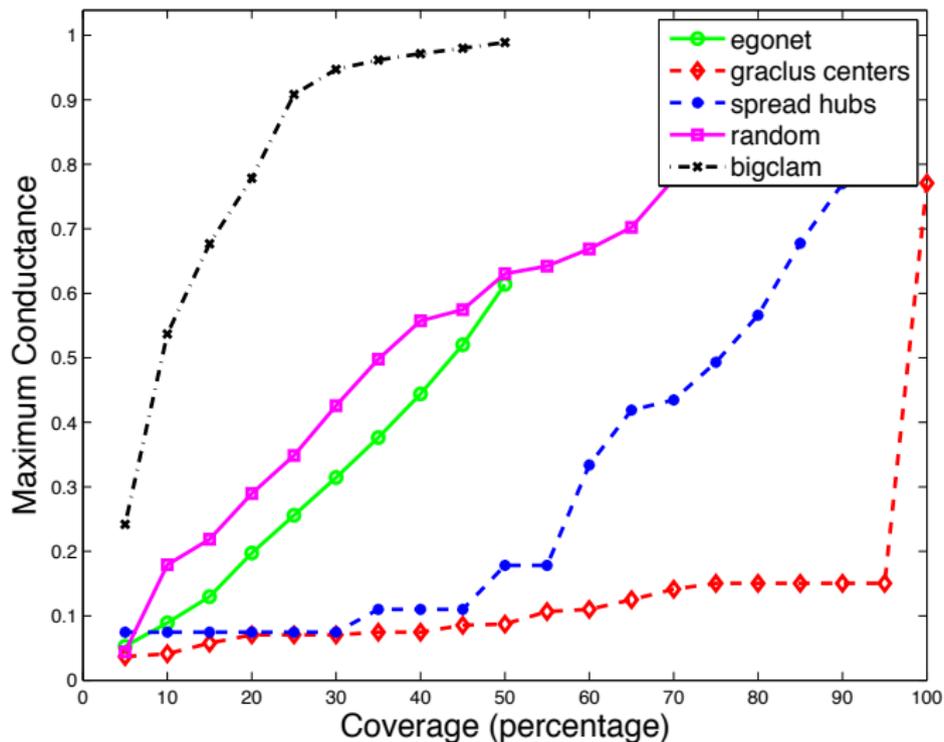
Community Quality using Conductance

- arXiv CondMat collaboration network (21,363 nodes)



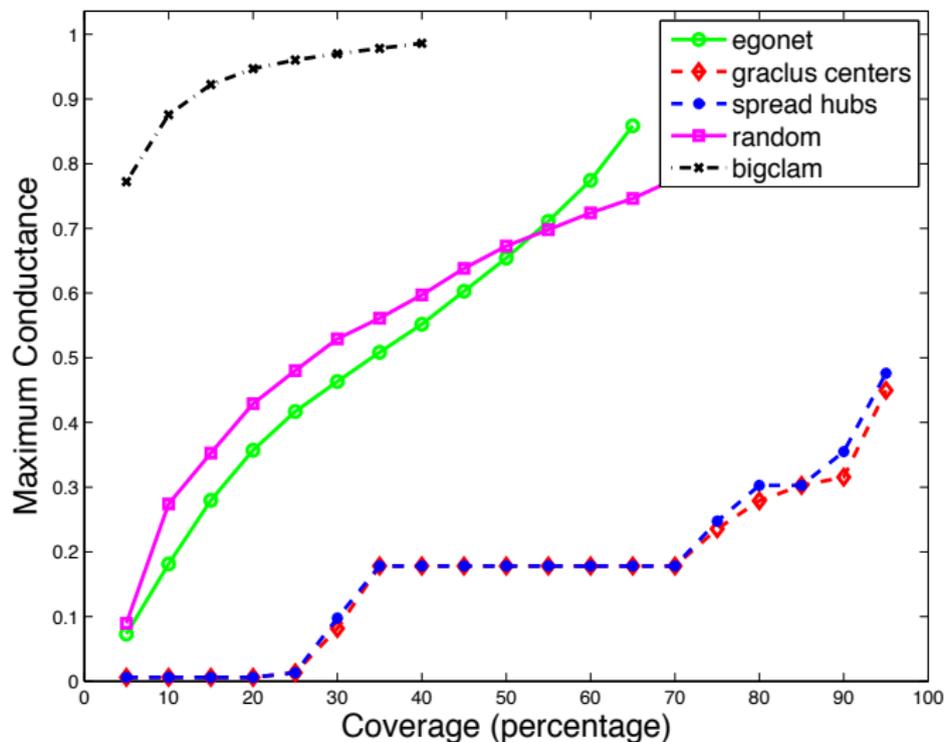
Community Quality using Conductance

- Flickr (1,994,422 nodes)
 - Demon fails on Flickr.



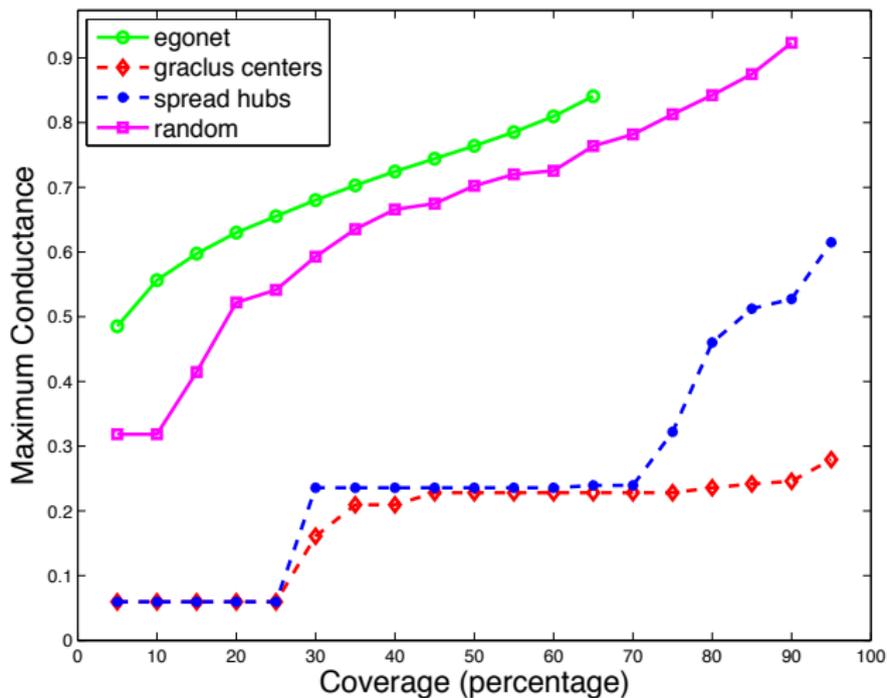
Community Quality using Conductance

- LiveJournal (1,757,326 nodes)
 - Demon fails on LiveJournal.



Community Quality using Conductance

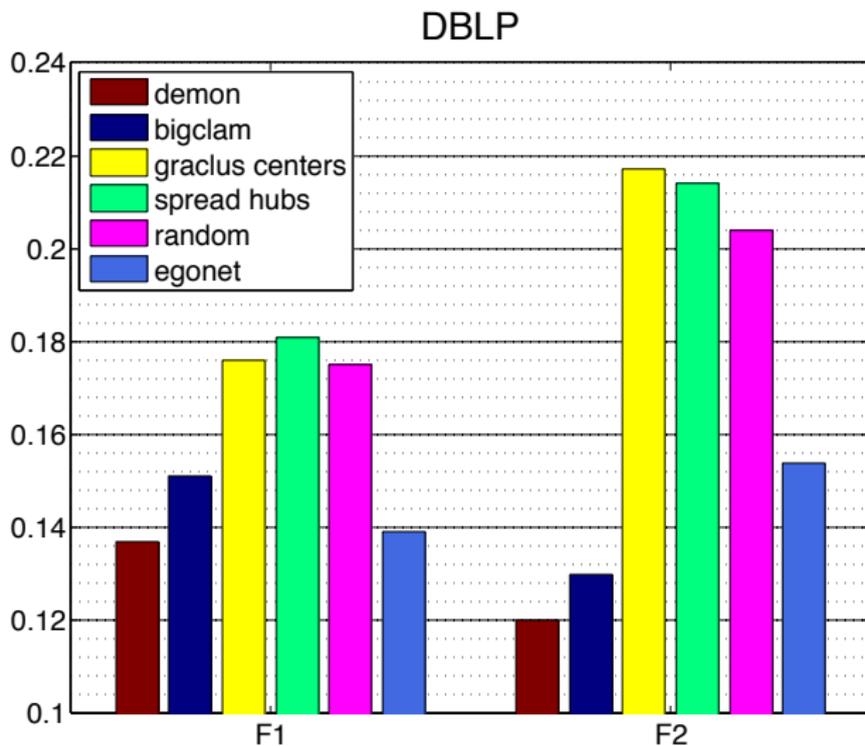
- Myspace (2,086,141 nodes)
 - Demon fails on Myspace.
 - Bigclam does not finish after running for one week.



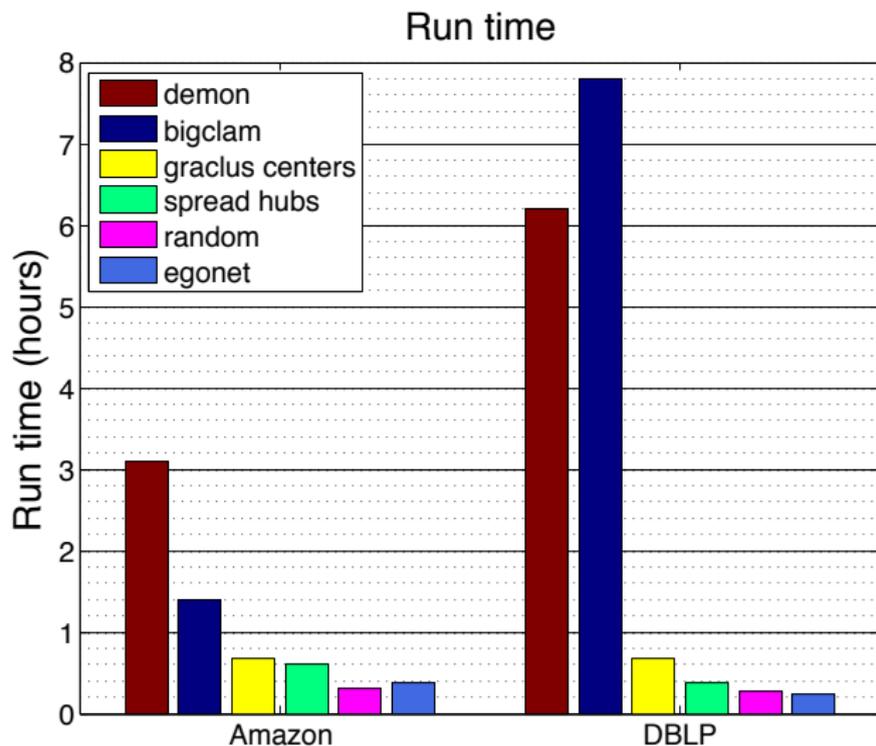
Community Quality via Ground Truth

- Precision
 - how many vertices are actually in the same ground truth community
- Recall
 - how many vertices are predicted to be in the same community in a retrieved community
- Compute F_1 , and F_2 measures
 - The ground truth communities are partially annotated.
 - F_2 measure puts more emphasis on recall than precision

Community Quality via Ground Truth



Comparison of Running Times



Conclusions

Conclusions

- Efficient overlapping community detection algorithm
 - Uses a seed set expansion
- Two seed finding strategies
 - Graclus centers
 - Spread hubs
- Our new seeding strategies are better than other strategies, and are thus effective in finding good overlapping clusters in a graph.
- The seed set expansion approach significantly outperforms other state-of-the-art methods.

References

- I. S. Dhillon, Y. Guan, and B. Kulis. Weighted graph cuts without eigenvectors: A multilevel approach. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 29, no. 11, pp. 1944-1957, 2007.
- R. Andersen, F. Chung and K. Lang. Local graph partitioning using PageRank vectors. In *FOCS*, 2006.
- D. F. Gleich and C. Seshadhri. Vertex neighborhoods, low conductance cuts, and good seeds for local community methods. In *KDD*, pages 597-605, 2012.
- R. Andersen and K. J. Lang. Communities from seed sets. In *WWW*, pages 223-232, 2006.
- J. Yang and J. Leskovec. Overlapping community detection at scale: a nonnegative matrix factorization approach. In *WSDM*, pages 587-596, 2013.
- M. Coscia, G. Rossetti, F. Giannotti, and D. Pedreschi. Demon: a local-first discovery method for overlapping communities. In *KDD*, 2012.