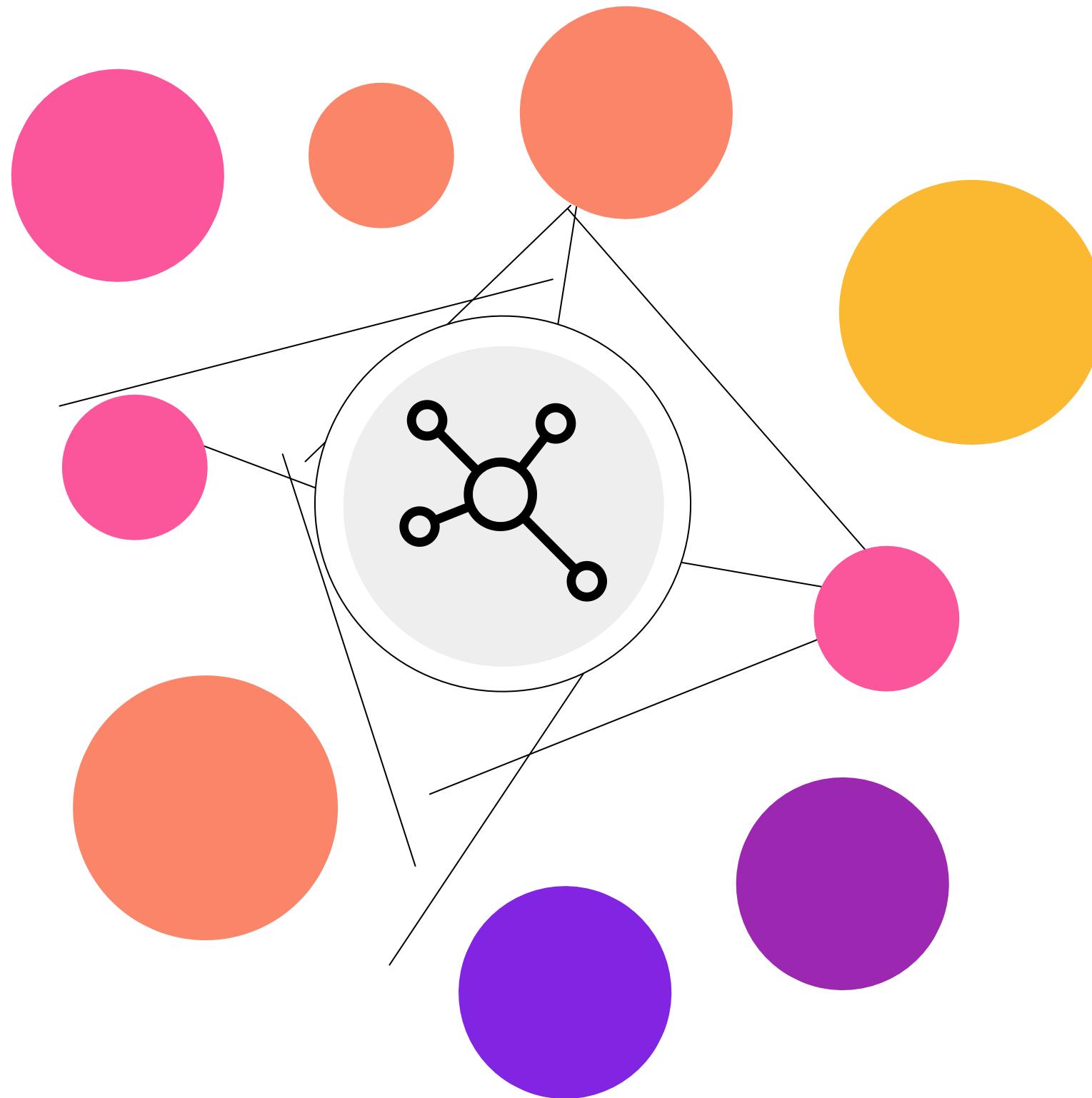


CS328



Sparsifying Graphs While Preserving Communities

Guntas Singh Saran (22110089)
Hrriday V. Ruparel (22110099)
Yajurvedh Bodala (19110077)



Context and Problem Statement



Data



Communities



Relevance



Complexity

Modelling of large data as graphs and networks is a common technique. Doing so is beneficial to capture the nature of connections between entities.

Among multiple regimes of graph analysis algorithms, community structure detection are much sought after. They extract the close relationships of social, biological and physical networks.

These algorithms are used for targeted advertising, identifying influential nodes, protein-protein interactions, brain connectivity networks and recommender systems.

However, as data size increases, the time taken to run community detection algorithms blows up. Not to mention, the storage costs also increases significantly.

Broad Problem Statement:

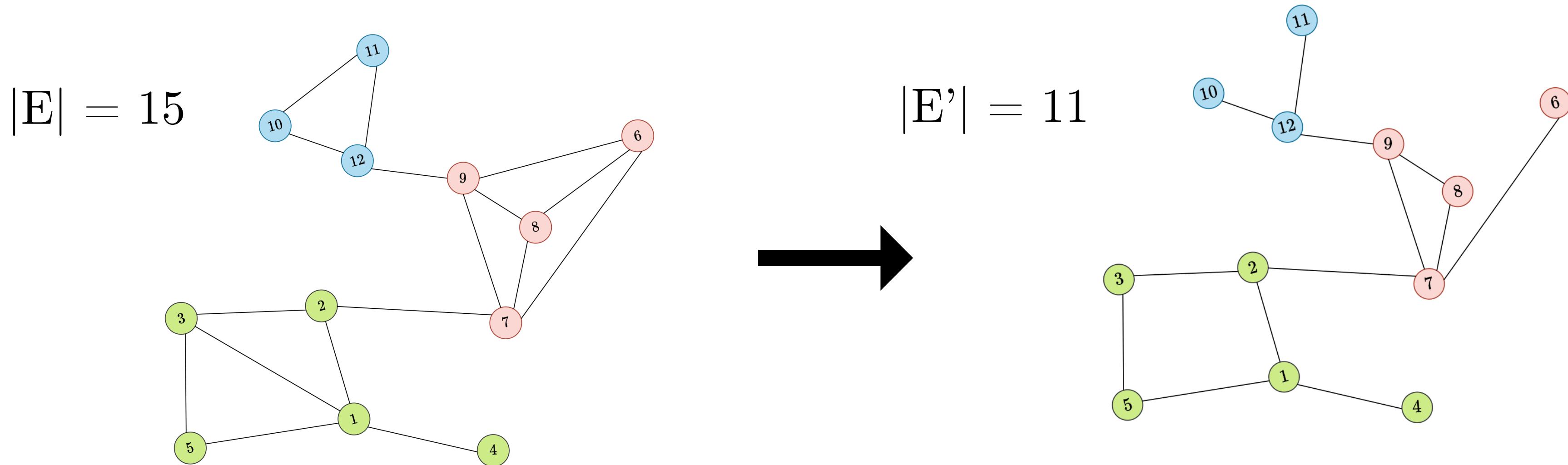
Given a graph $G(V,E)$, where V is the vertex set and E is the edge set, we wish to “sparsify” the graph in a “meaningful” manner and obtain a graph $G'(V,E')$, where V is the vertex set similar to original graph G and $E' \subset E$ is the edge set of sparsified graph G' .

Context and Problem Statement

Broad Problem Statement:

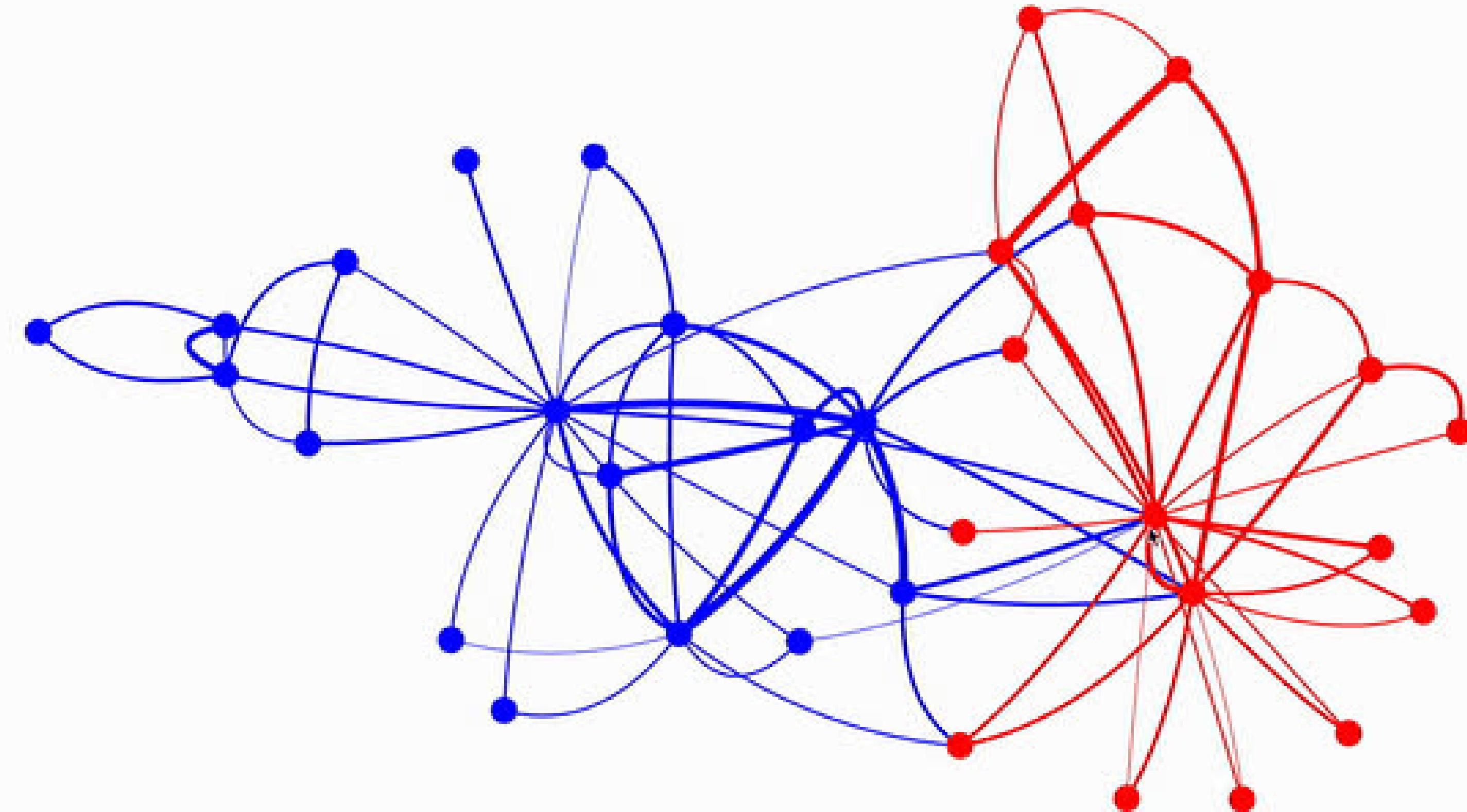
Given a graph $G(V,E)$, where V is the vertex set and E is the edge set, we wish to “sparsify” the graph in a “meaningful” manner and obtain a graph $G'(V,E')$, where V is the vertex set similar to original graph G and $E' \subset E$ is the edge set of sparsified graph G' .

By “sparsifying” a graph in a “meaningful manner”, we intend to preserve the community structure of the original graph G by sampling a subset of edges from G so that community detection algorithms can be implemented on the “sparsified” graph G' with minimum loss of generality and accuracy.



Graphs

Zachary's Karate Club Network



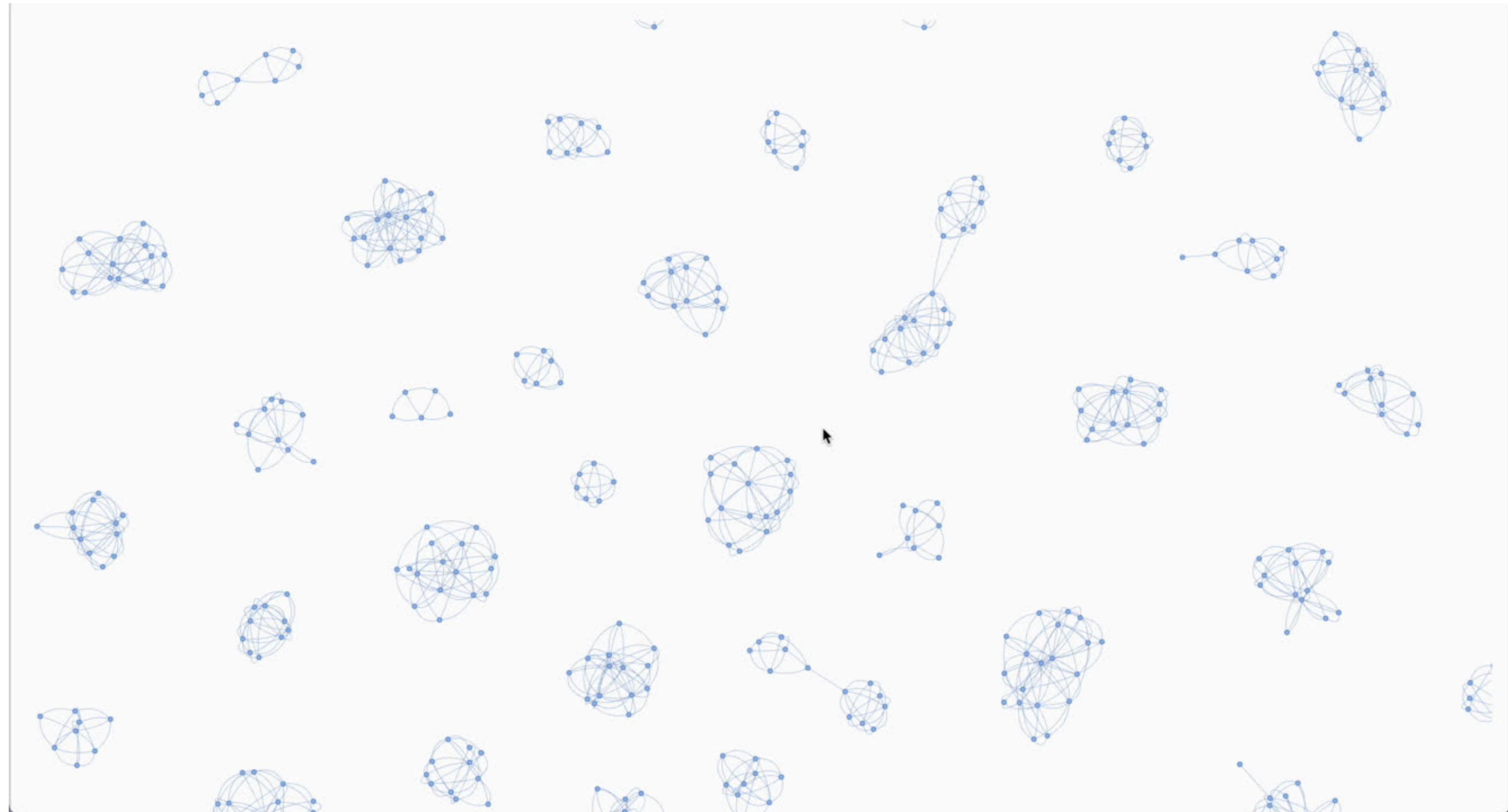
Communities in Networks

DBLP Network



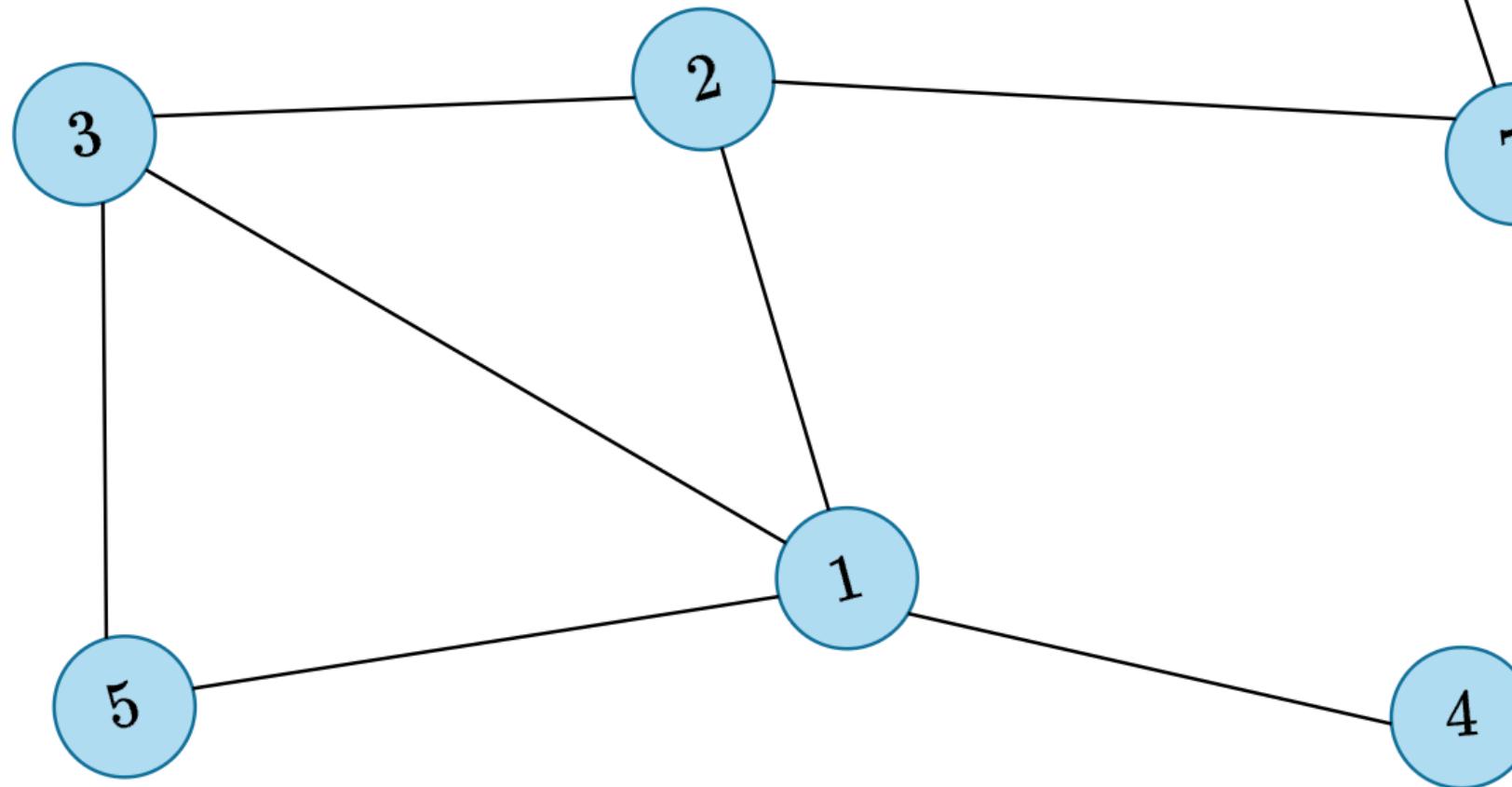
Communities in Networks

Amazon Co-Purchase



$$|V| = 12$$

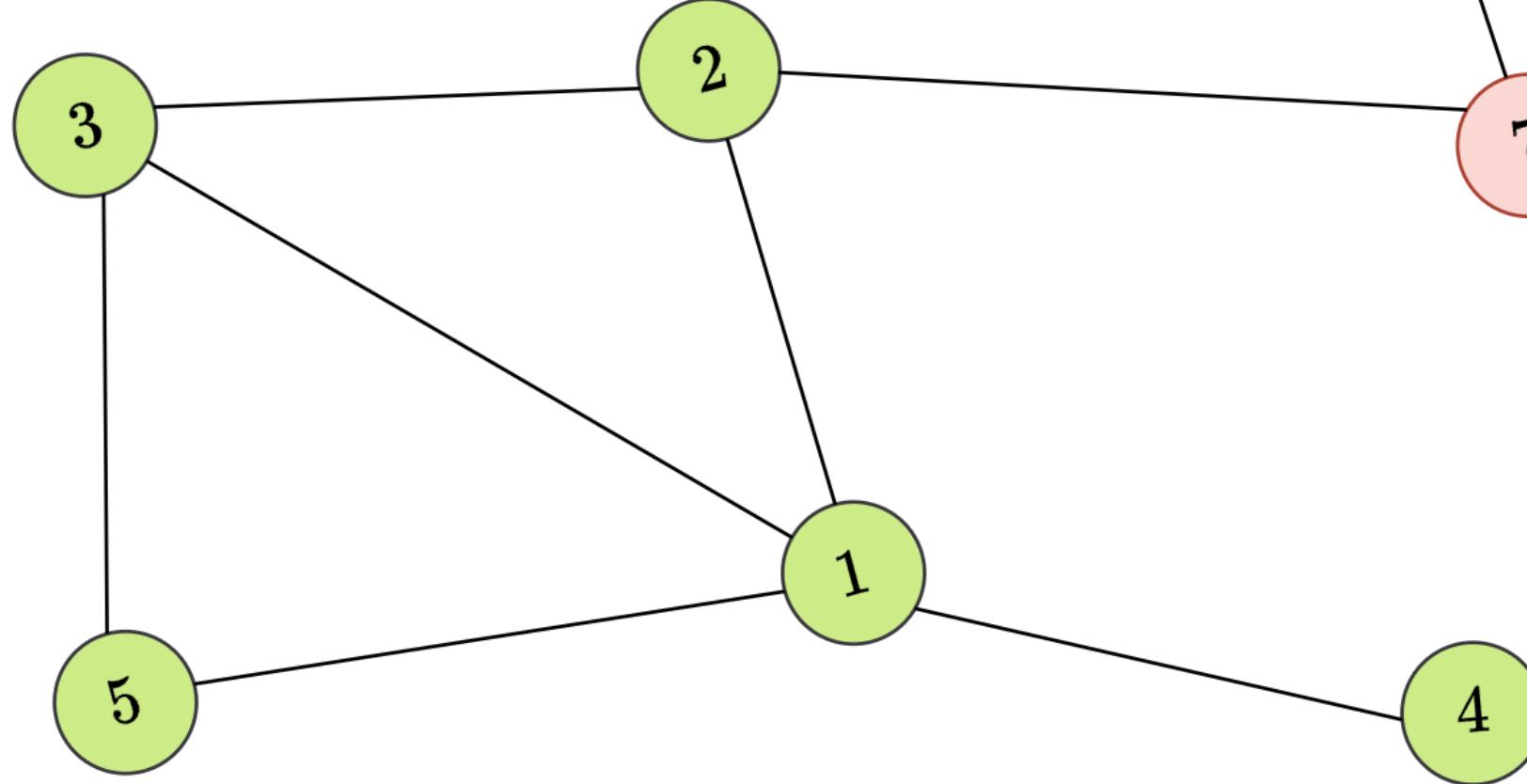
$$|E| = 15$$



Communities

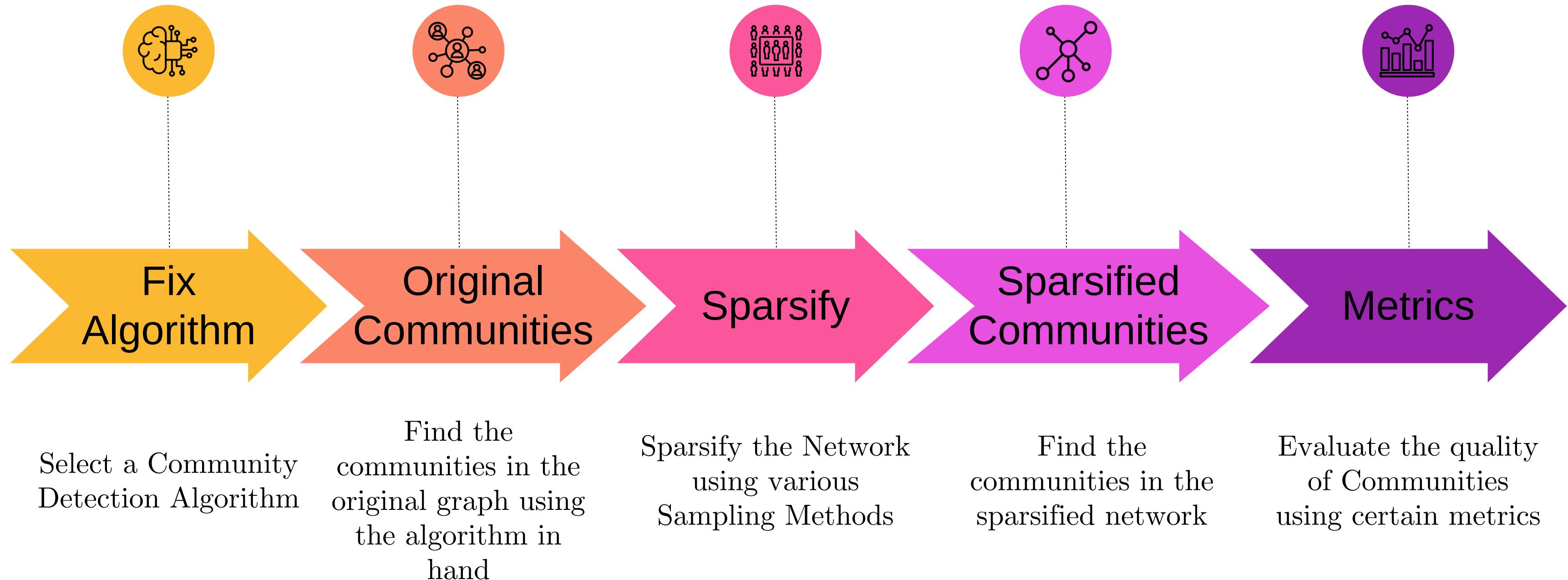
$$|V| = 12$$

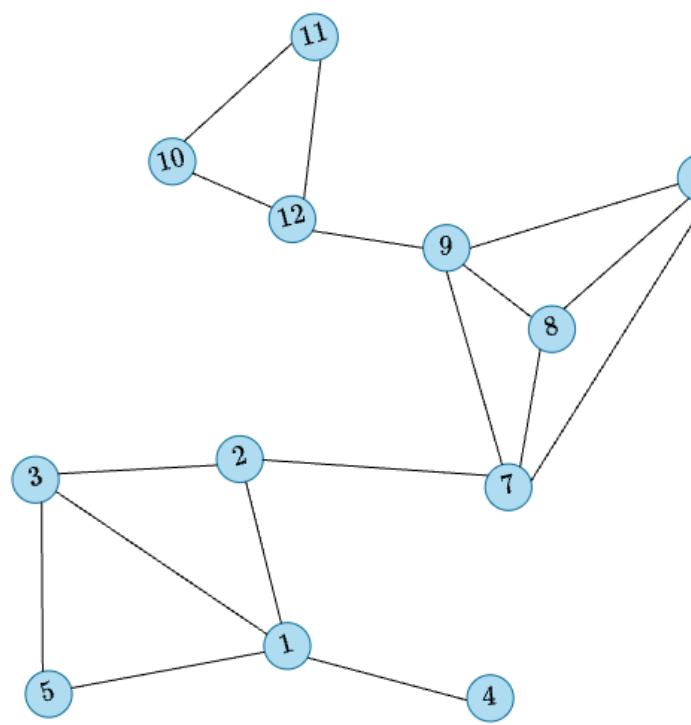
$$|E| = 15$$



Communities

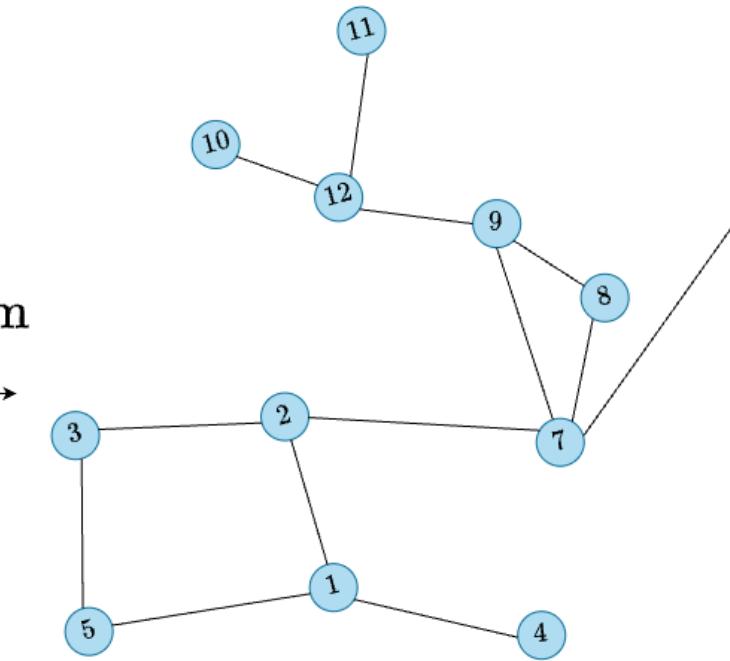
Evaluation Pipeline





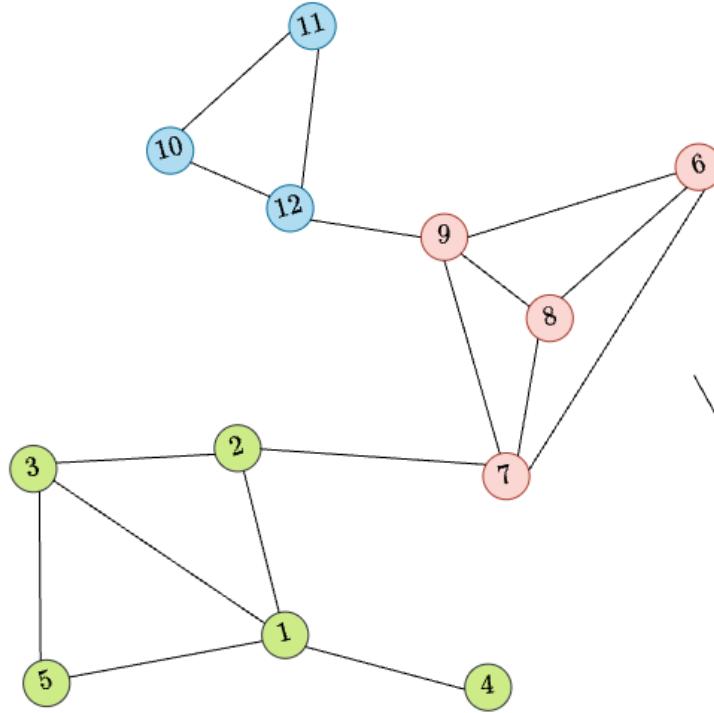
Original Graph

Sparsification Algorithm



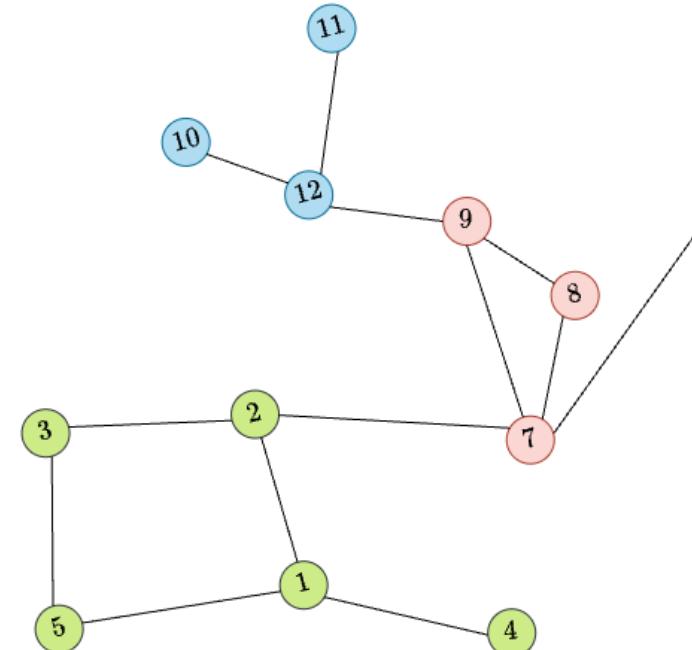
Sparsified Graph

Community Detection Algorithm

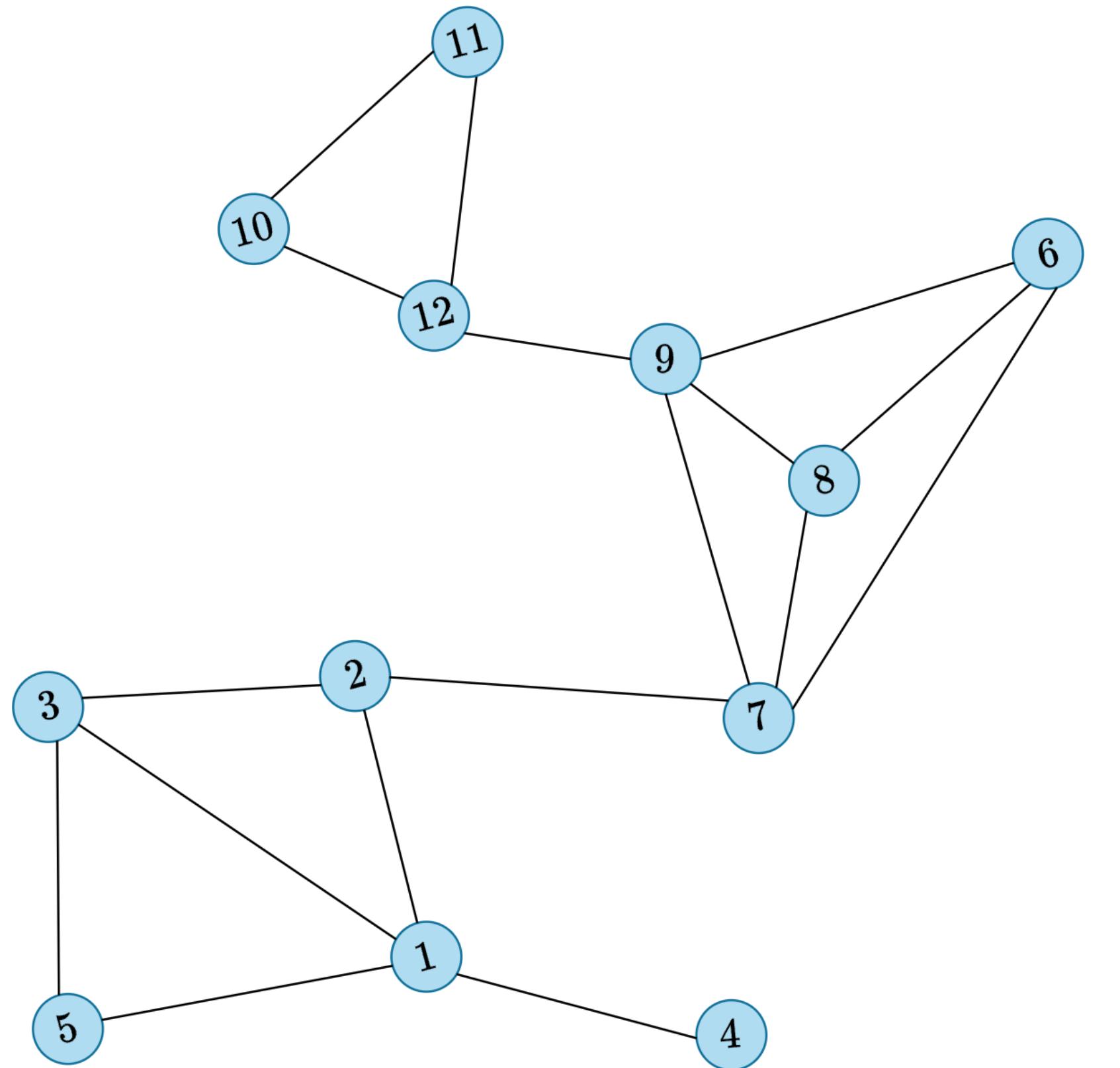


Original Graph Communities

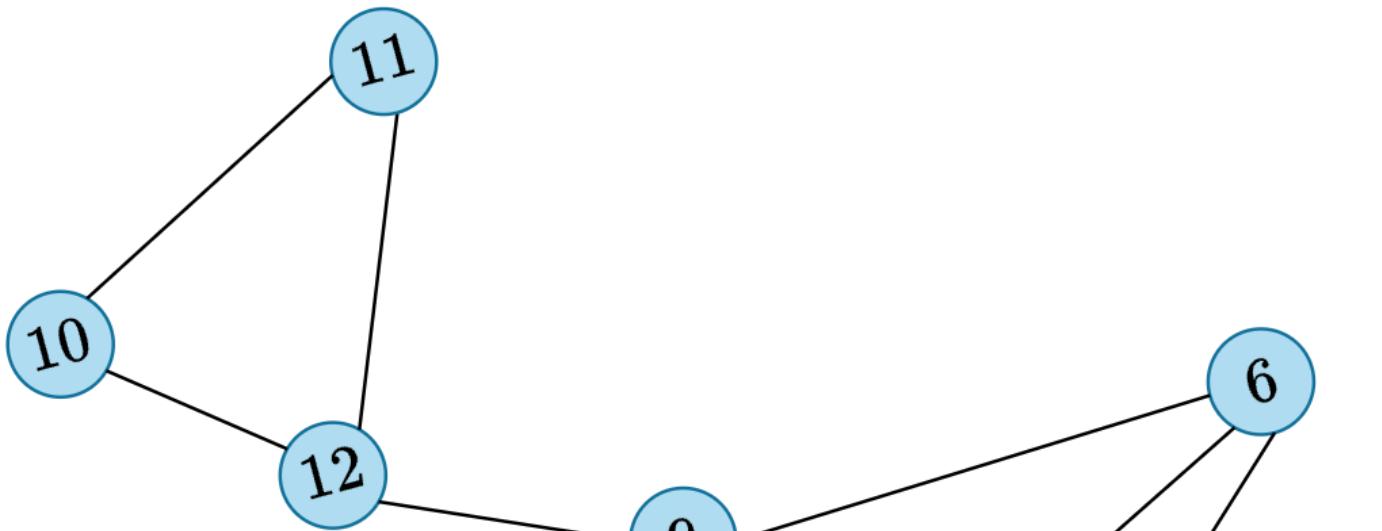
Metric



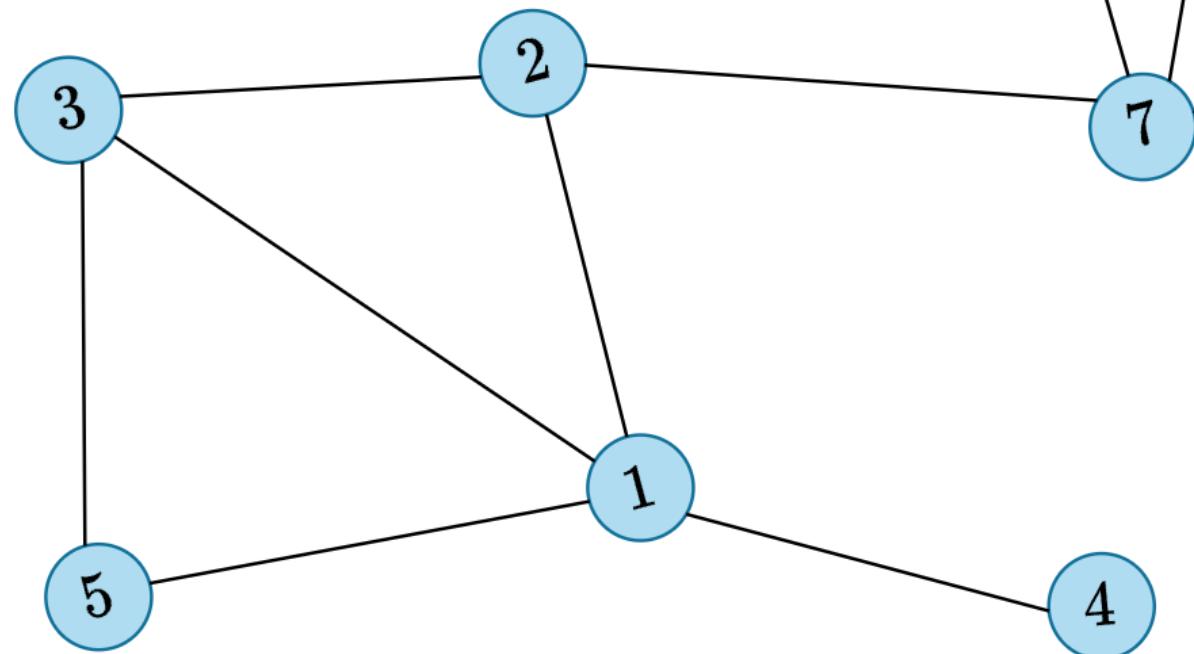
Sparsified Graph Communities



Original Graph

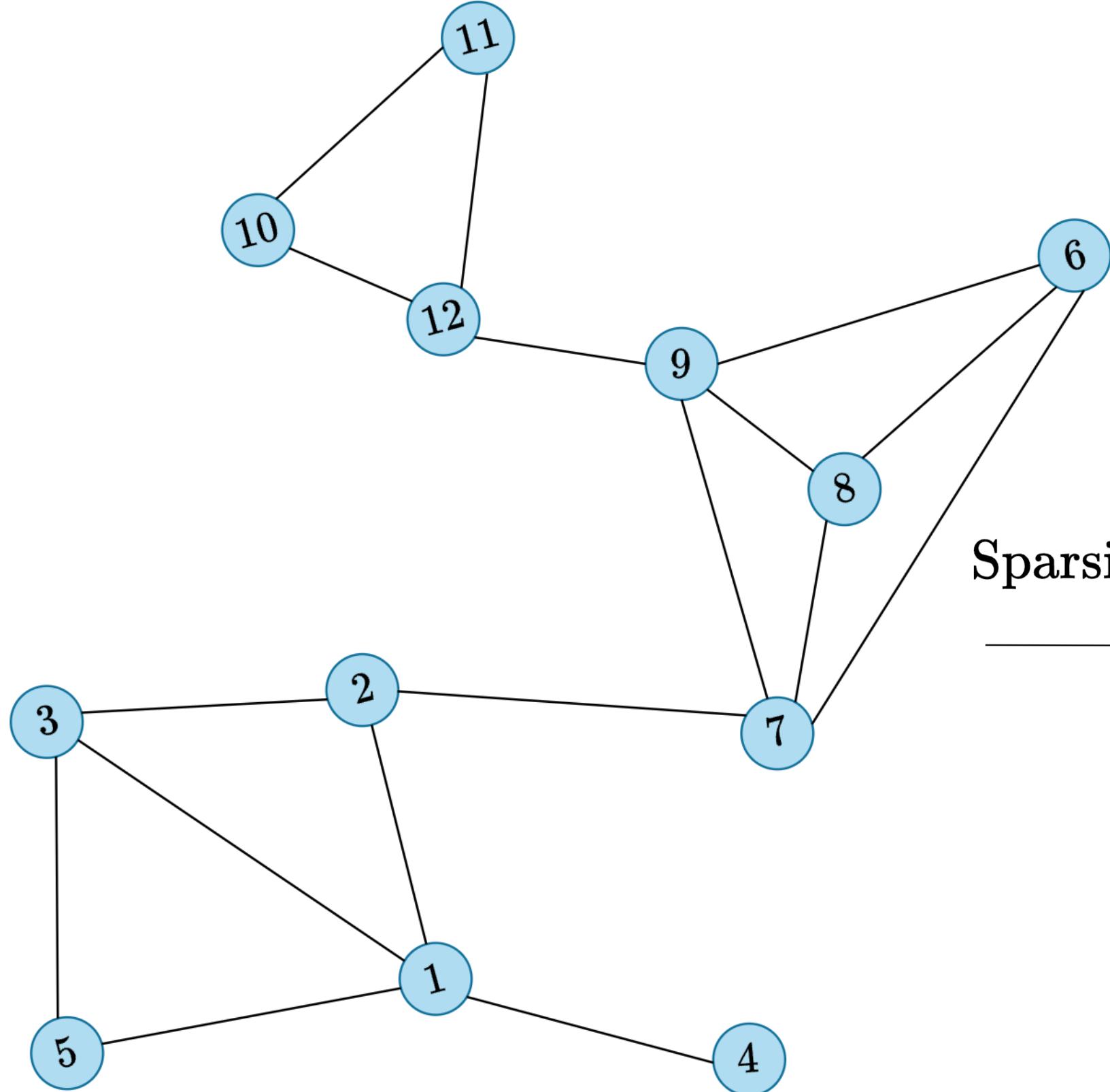


Original Graph



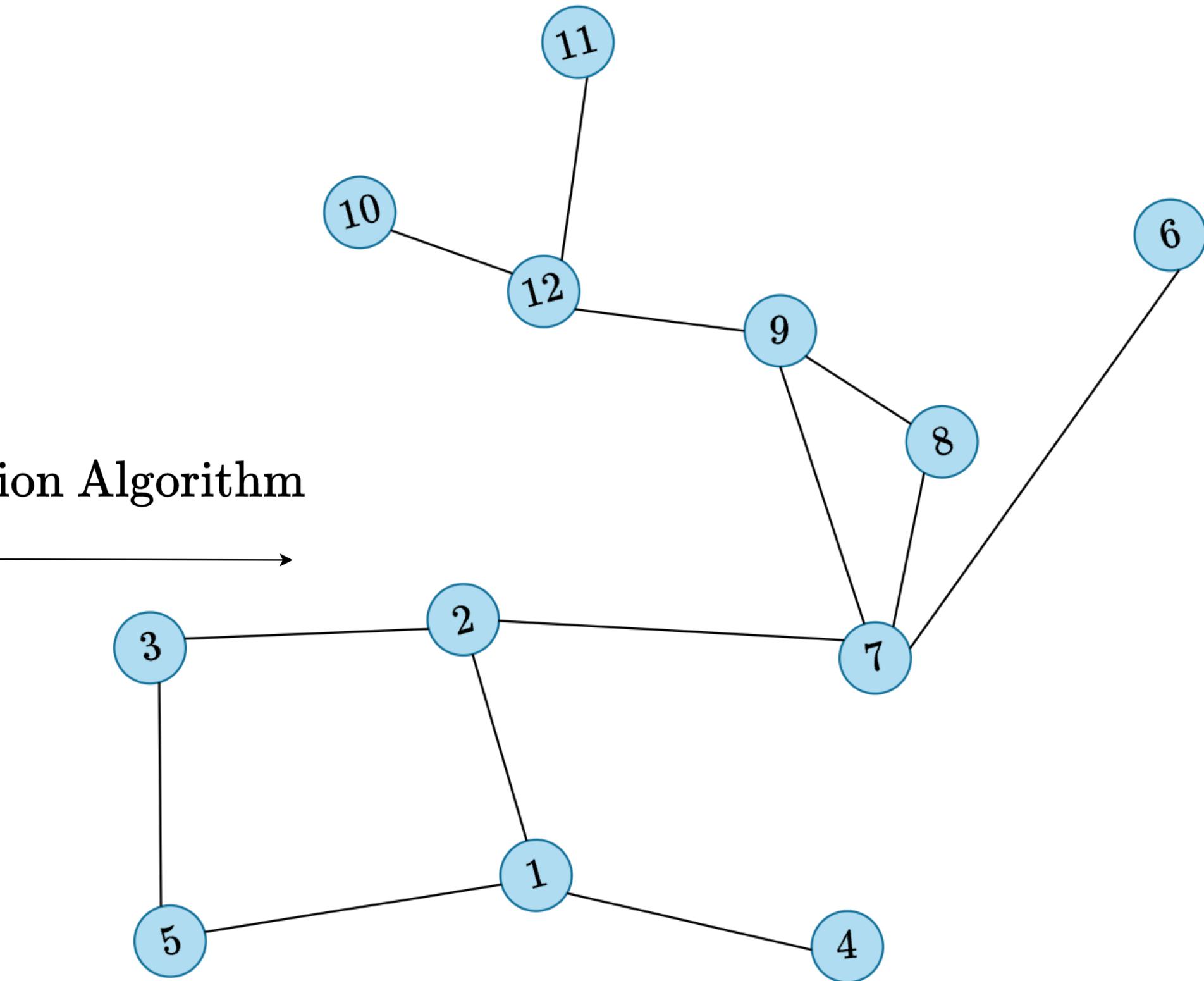
Sparsification Algorithm



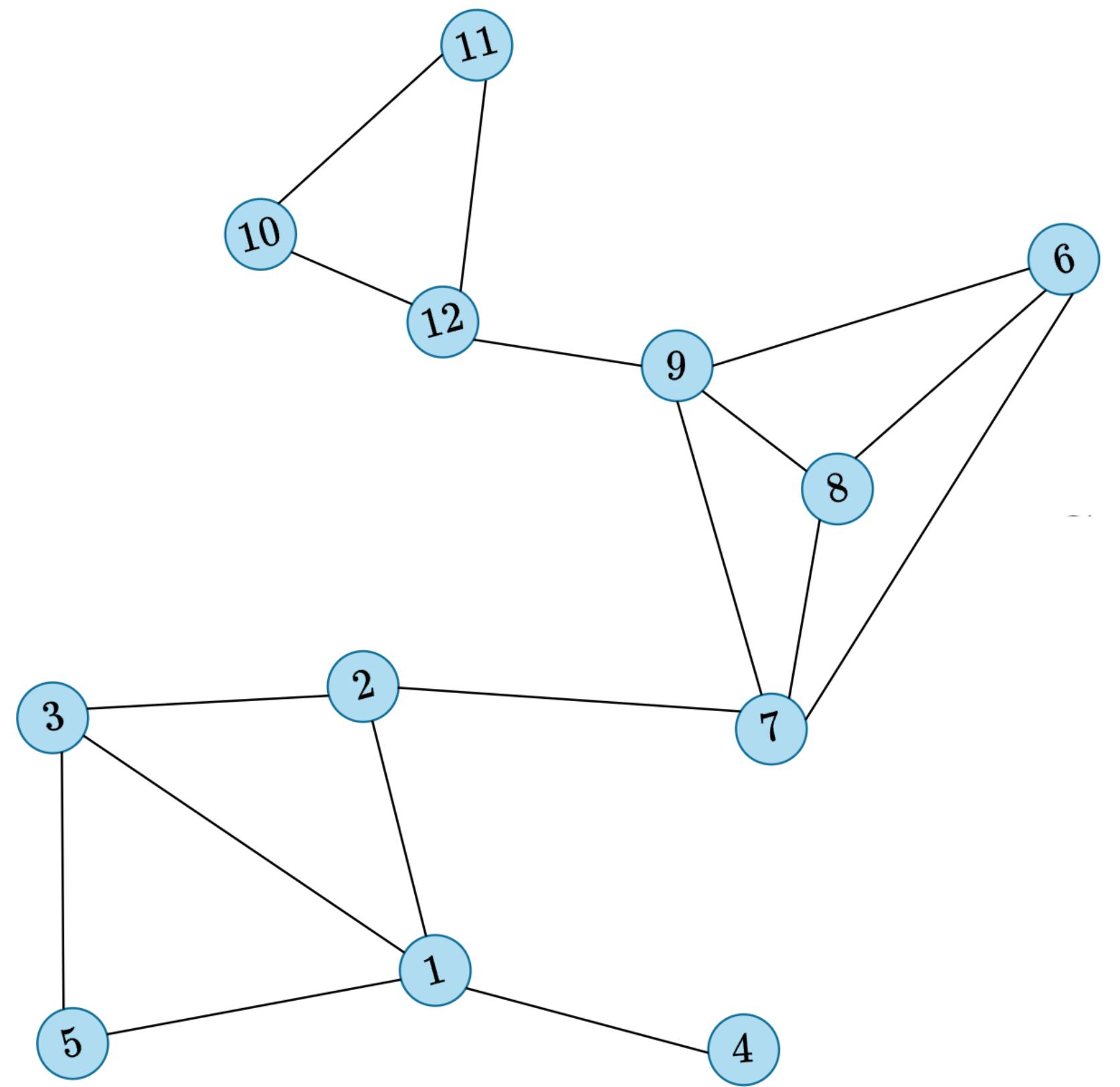


Original Graph

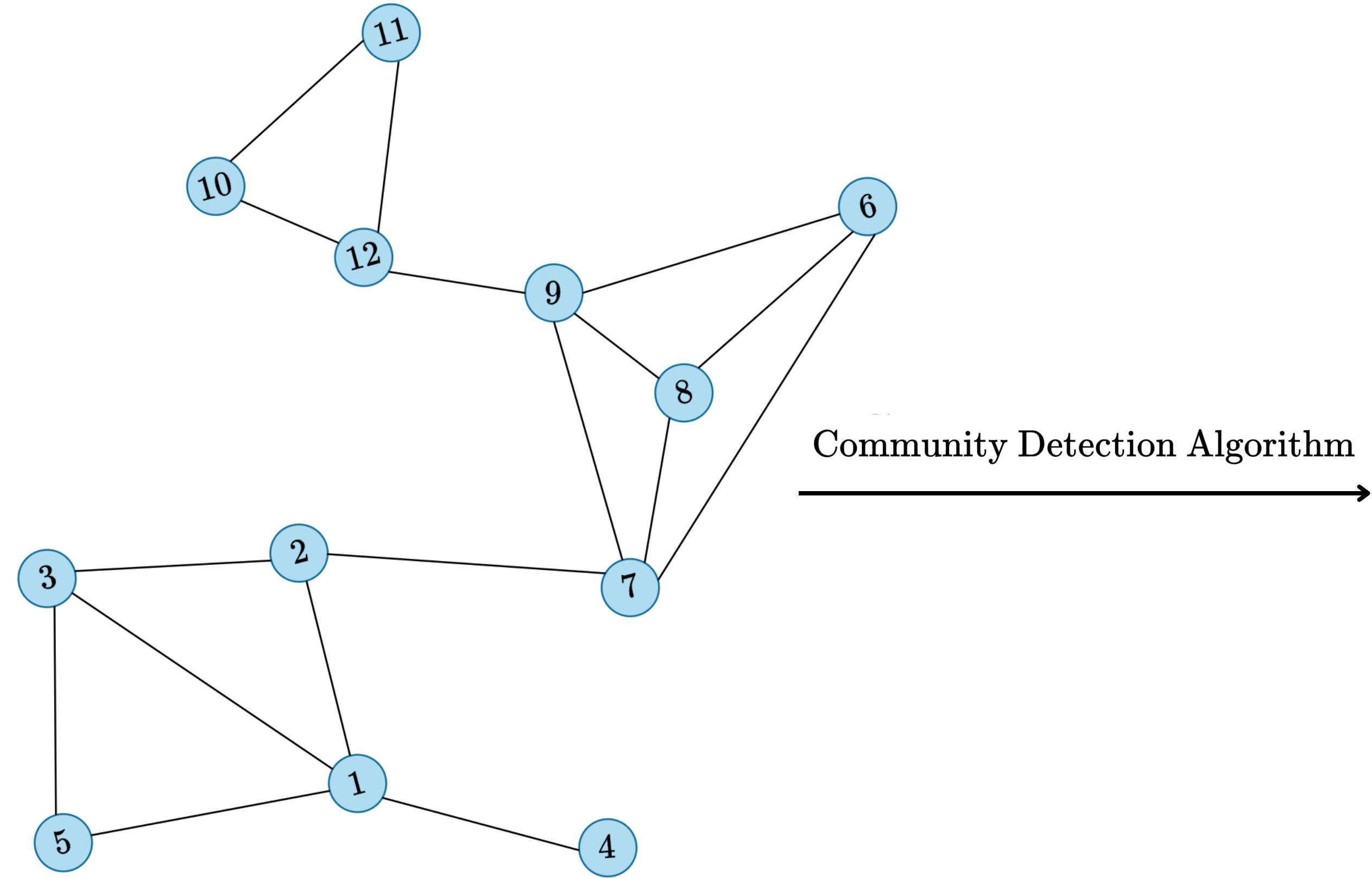
Sparsification Algorithm

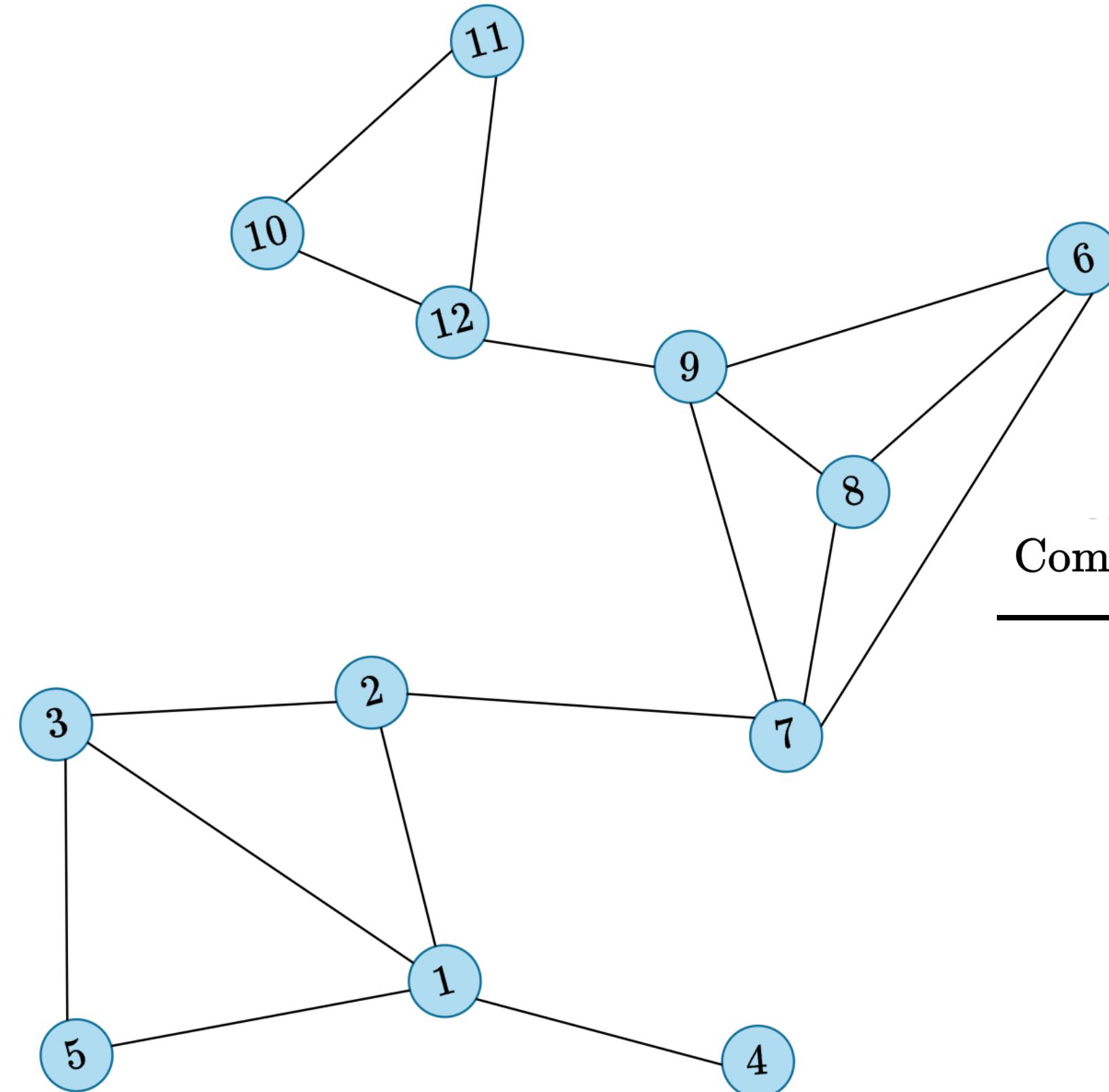


Sparsified Graph



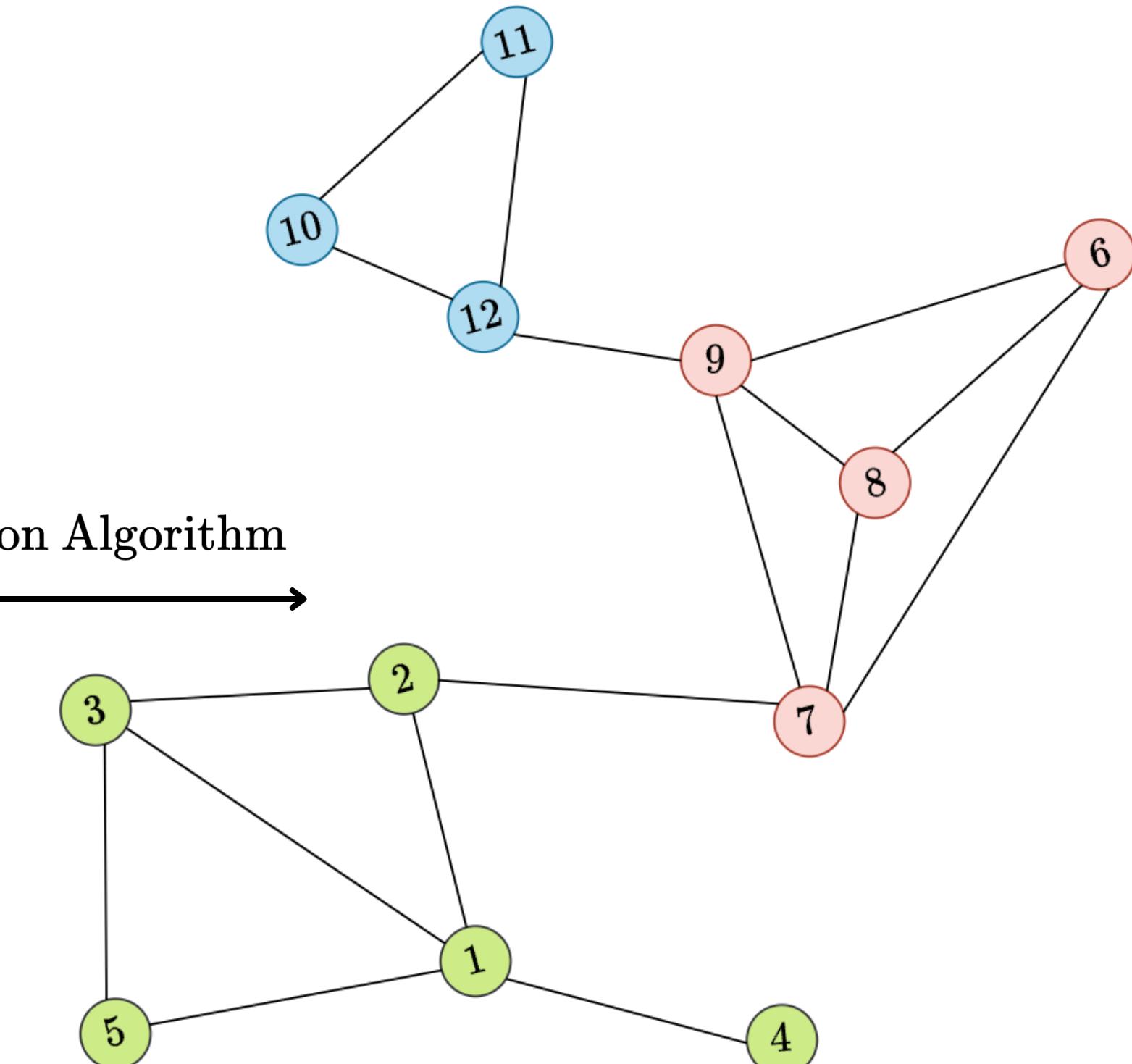
Original Graph



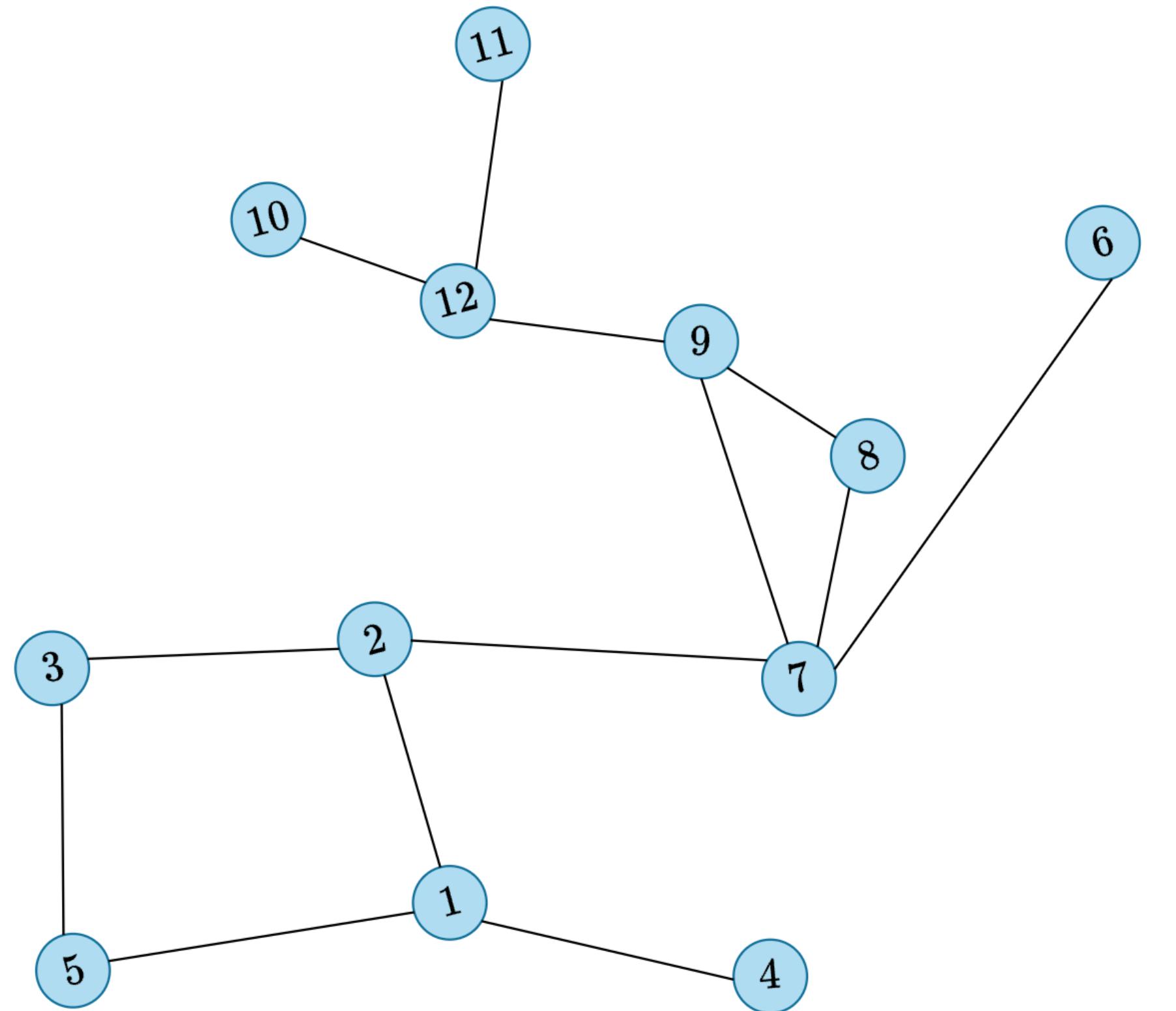


Original Graph

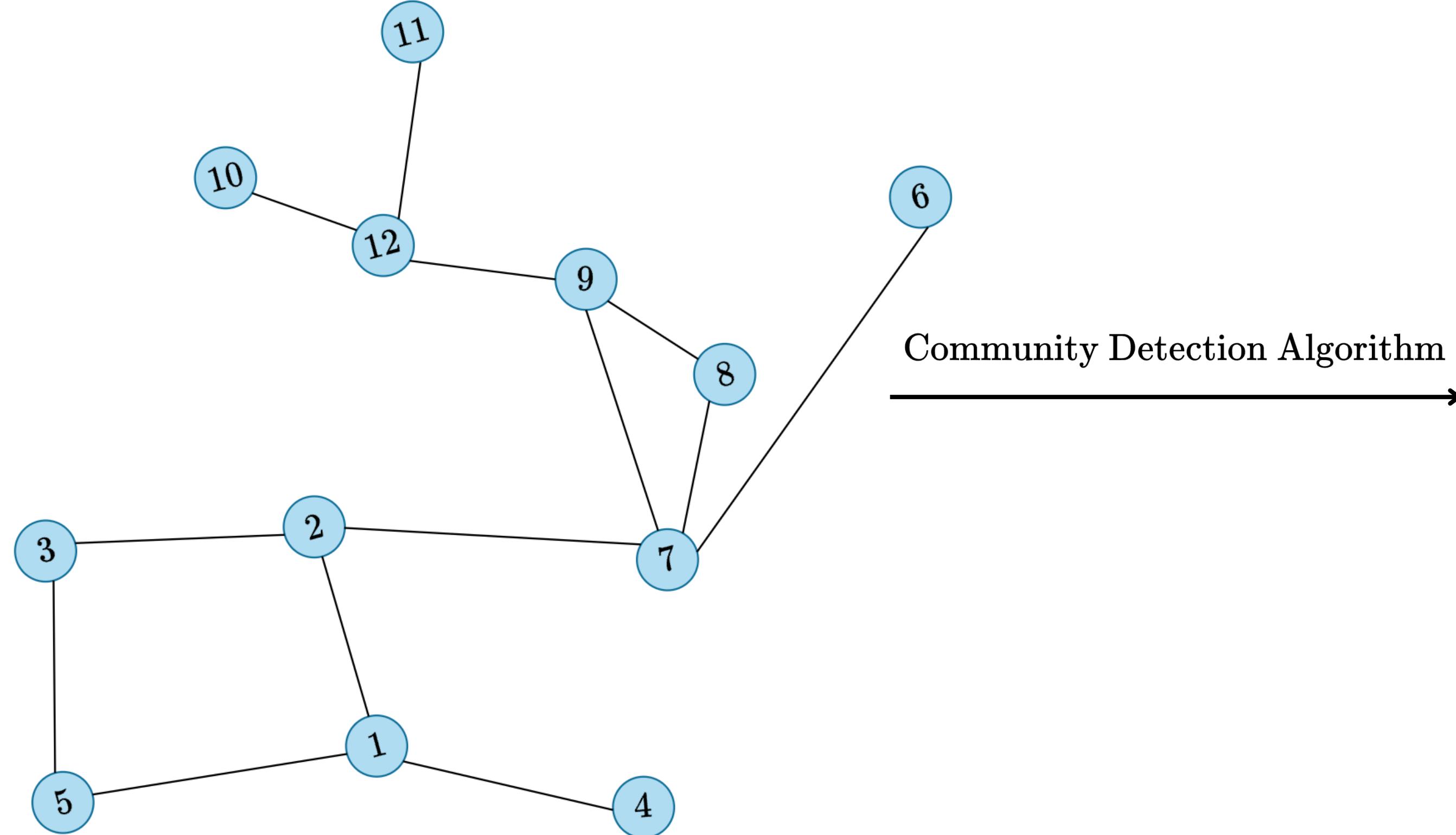
Community Detection Algorithm



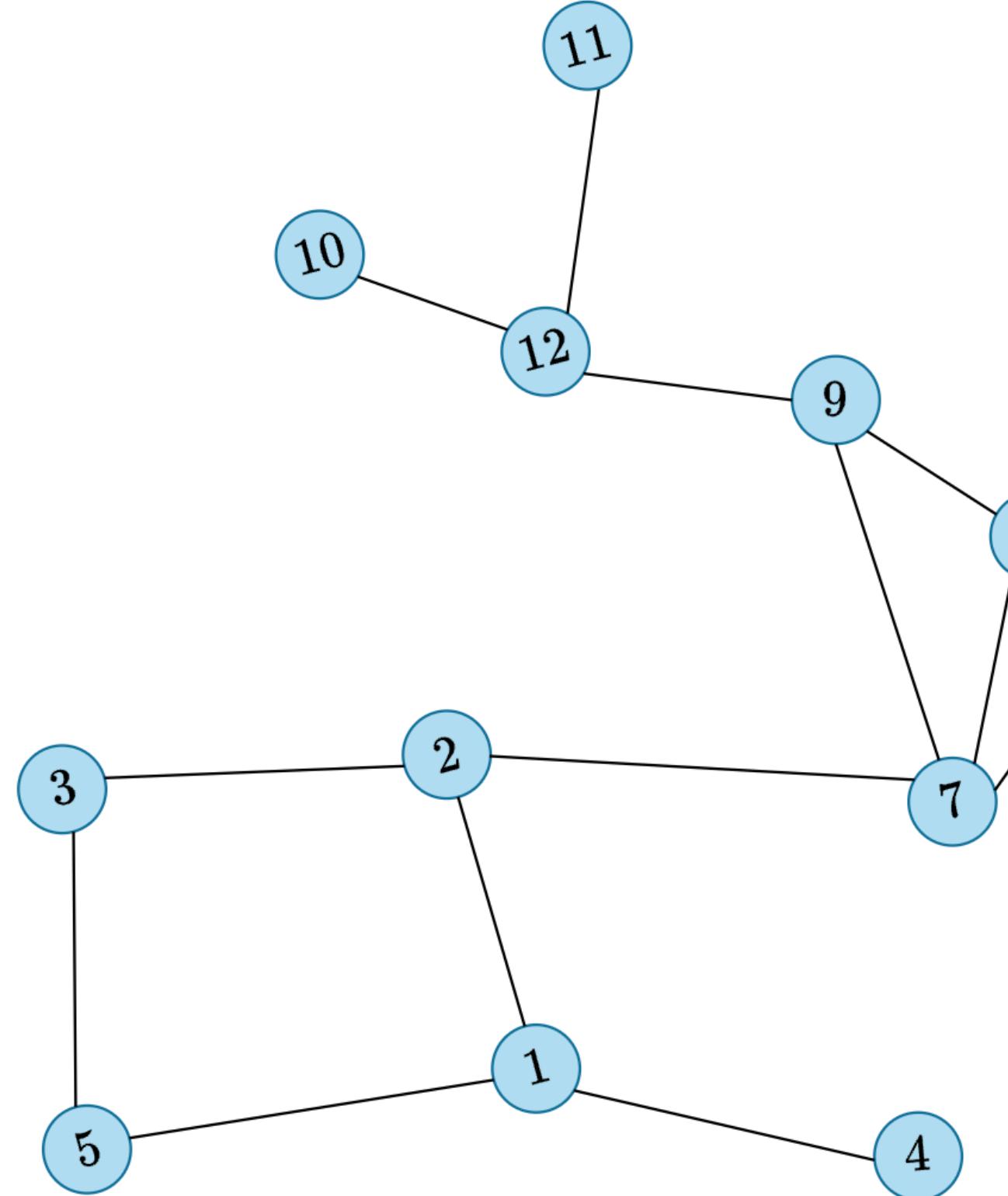
Original Graph Communities



Sparsified Graph

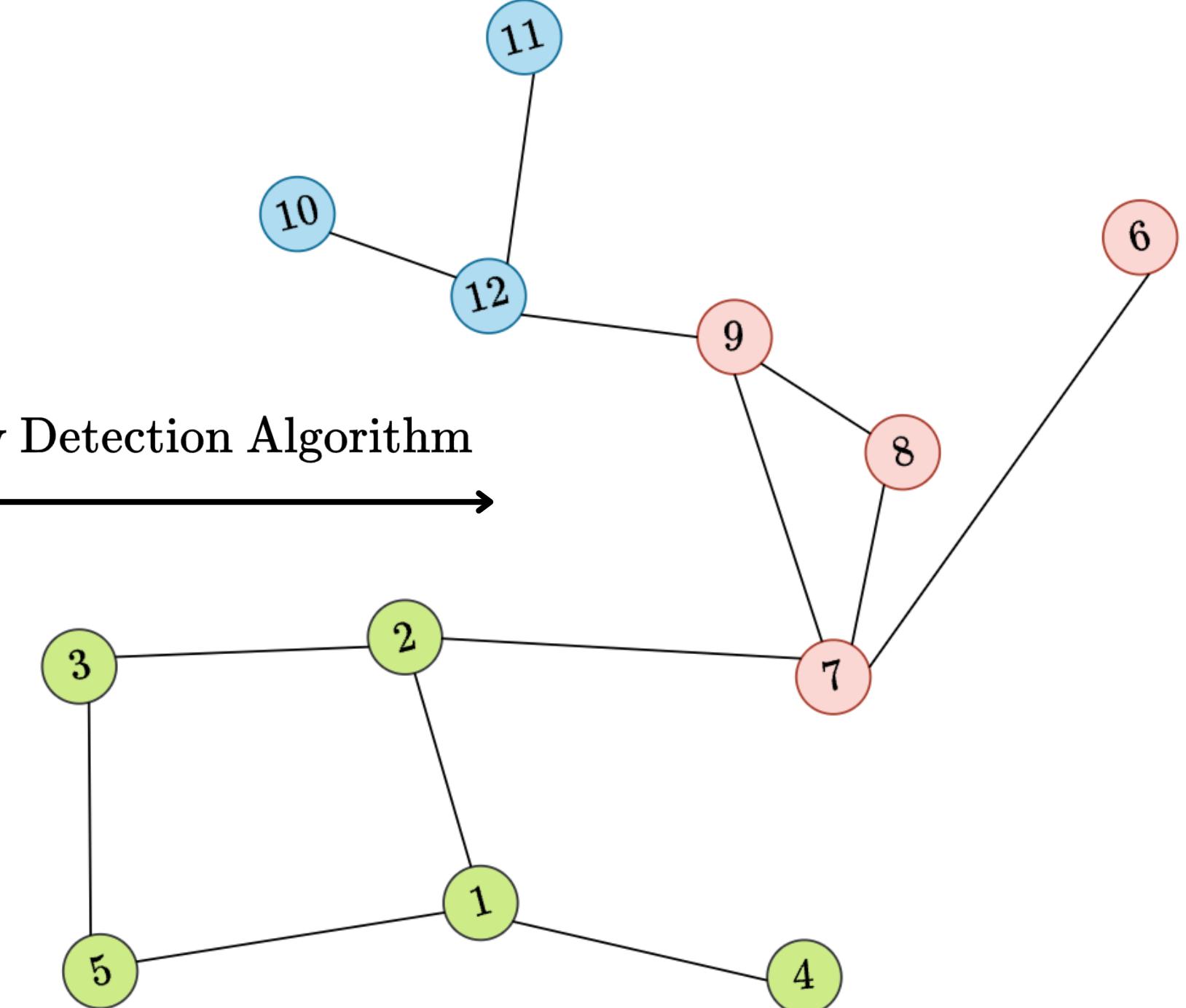


Sparsified Graph

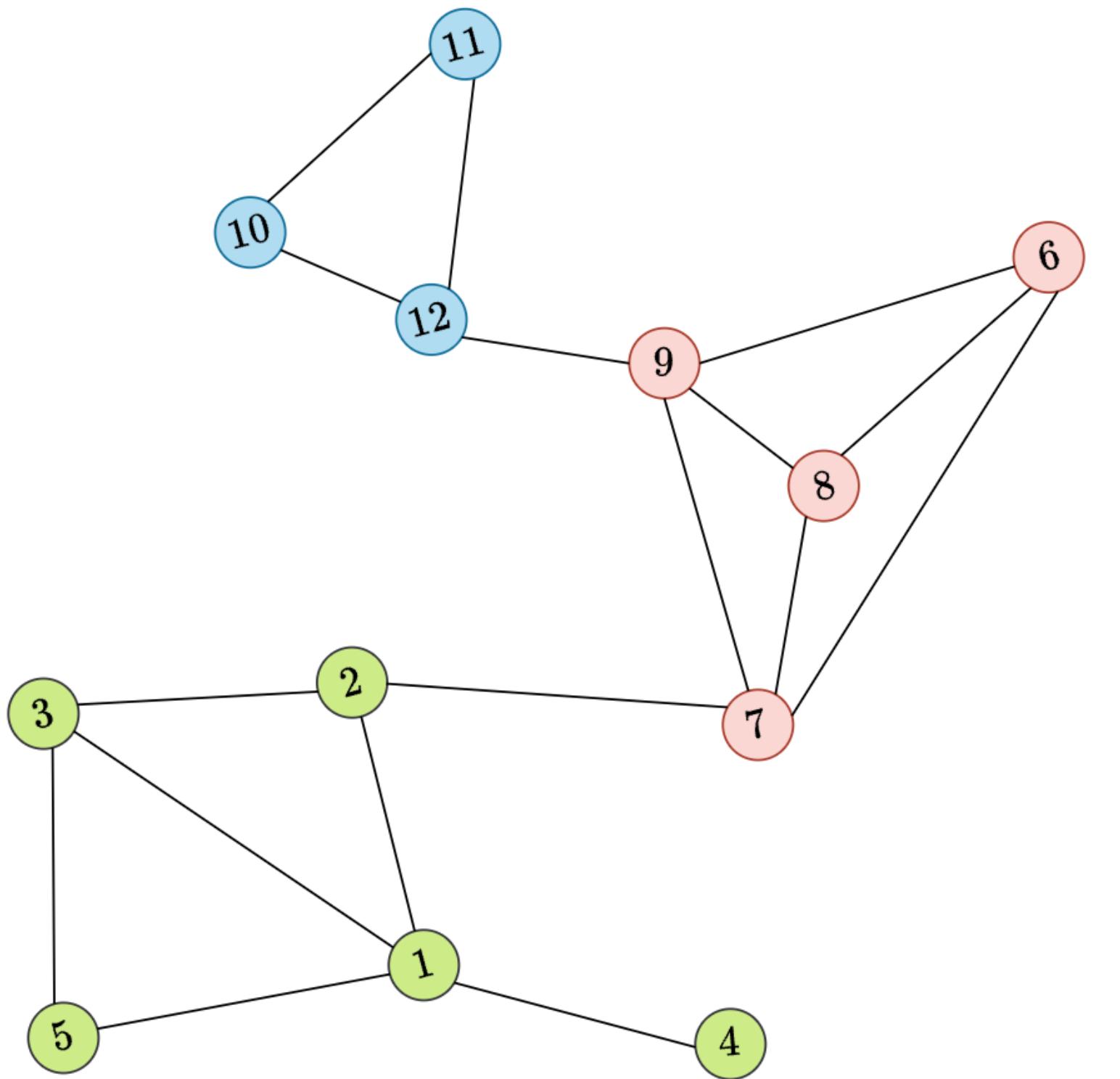


Sparsified Graph

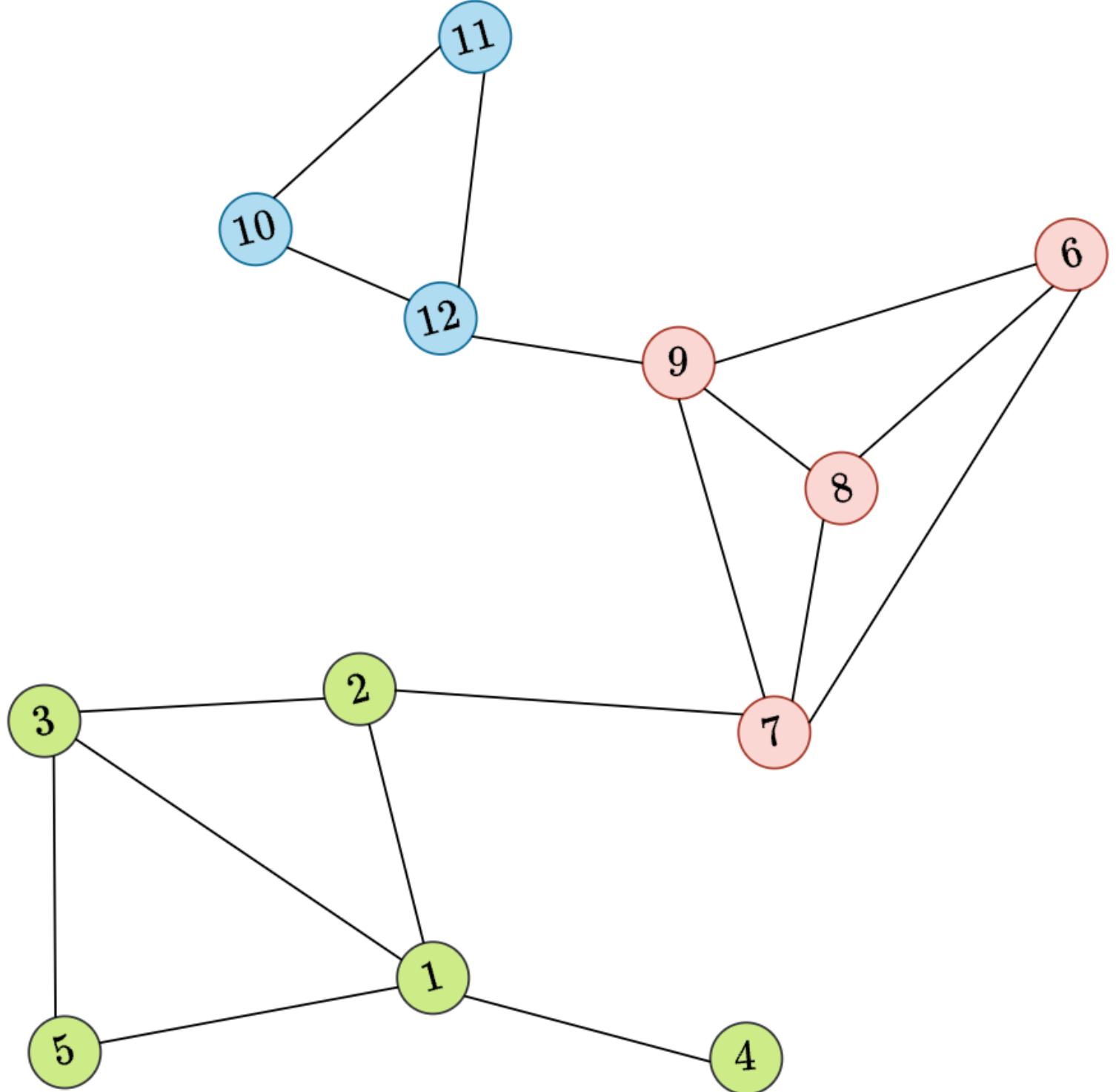
Community Detection Algorithm



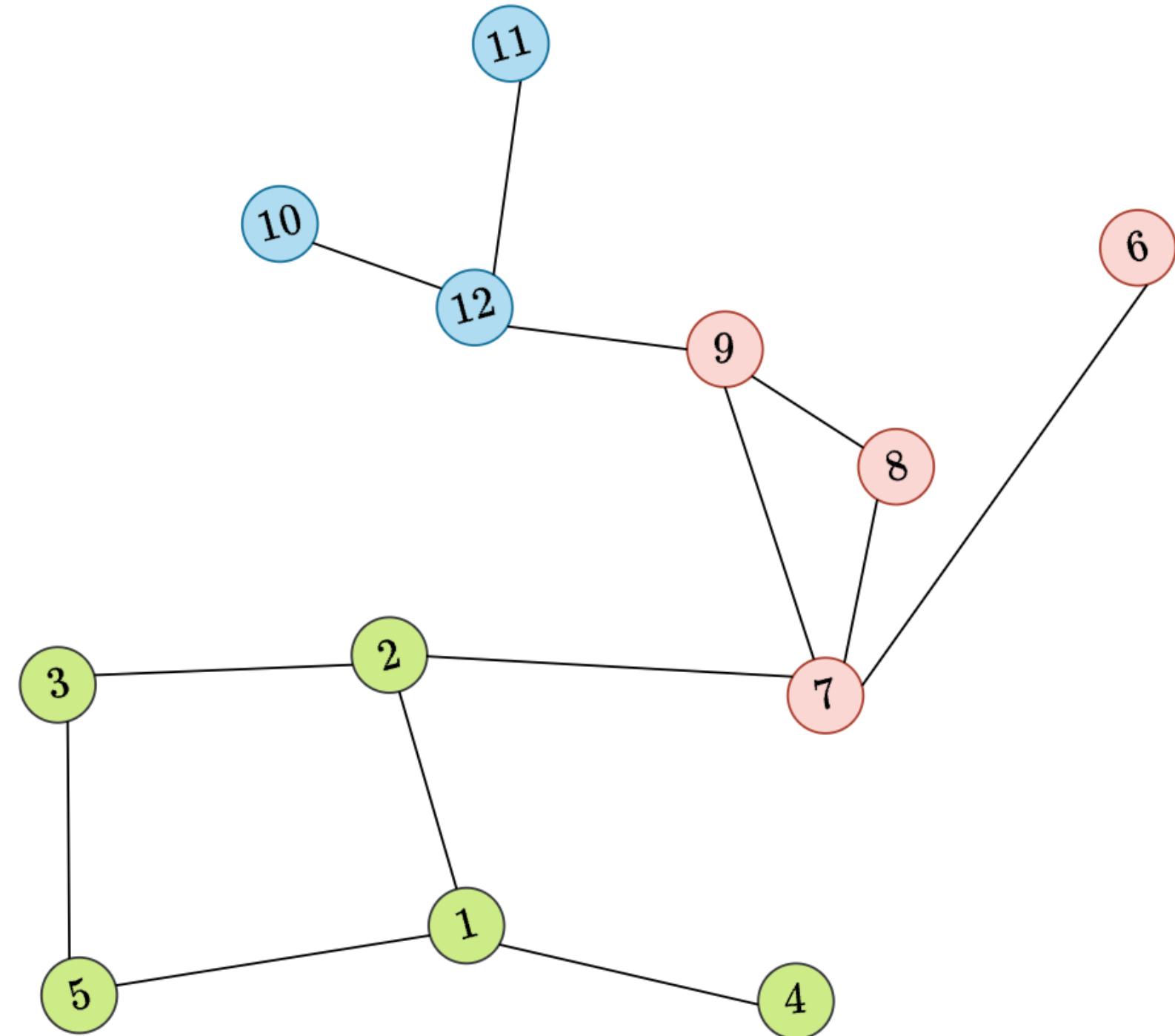
Sparsified Graph Communities



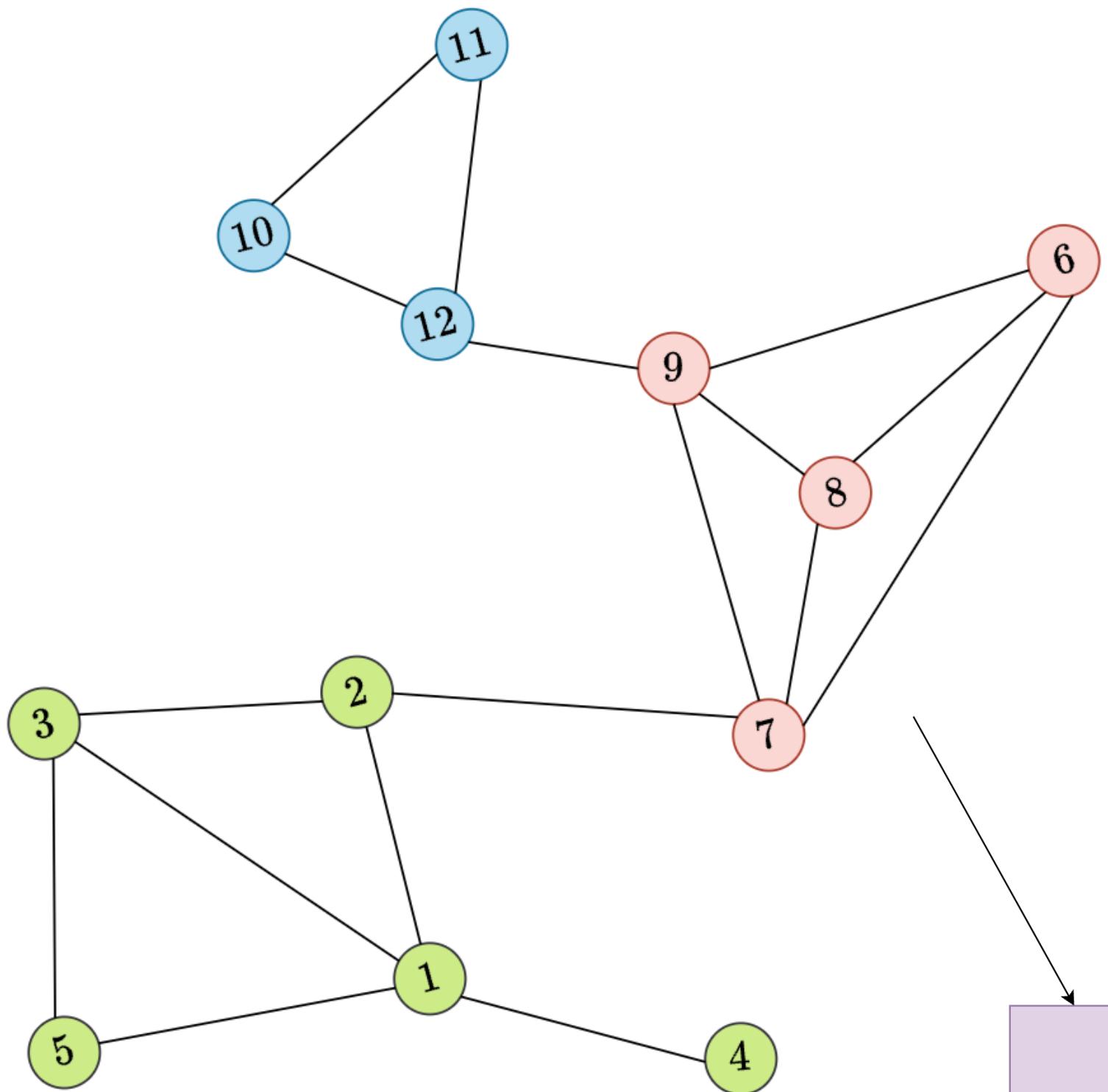
Original Graph Communities



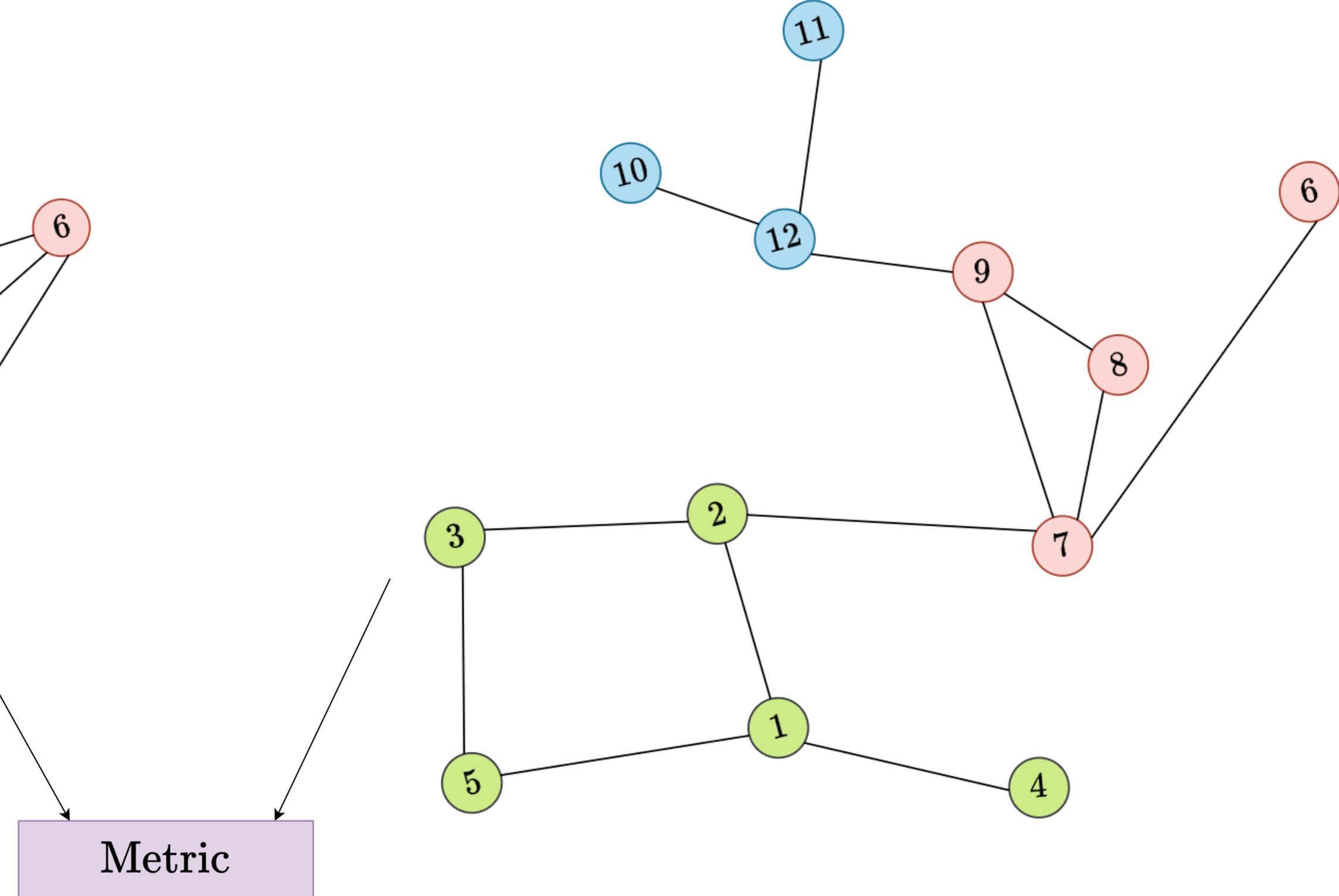
Original Graph Communities



Sparsified Graph Communities



Original Graph Communities



Sparsified Graph Communities

Metric

Salient Community Structure Properties

The number of shortest paths between pairs of nodes in a graph that pass through a particular edge. It quantifies the importance of an edge in connecting different parts of the network.

Edge Betweenness

Quantifies the degree to which a network is partitioned into communities or modules. Compares the number of edges within communities to the number of edges expected in a random network with the same degree distribution, indicating the presence of densely connected communities.

Jaccard Similarity

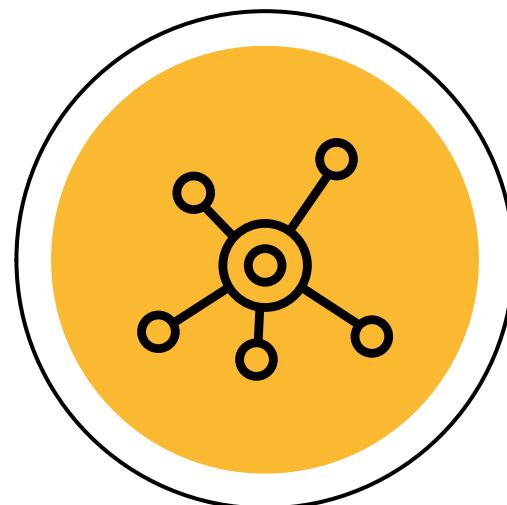
Measures similarity between two sets by comparing their intersection to their union. It is often used to measure the similarity between the sets of neighbors of two nodes.

Modularity

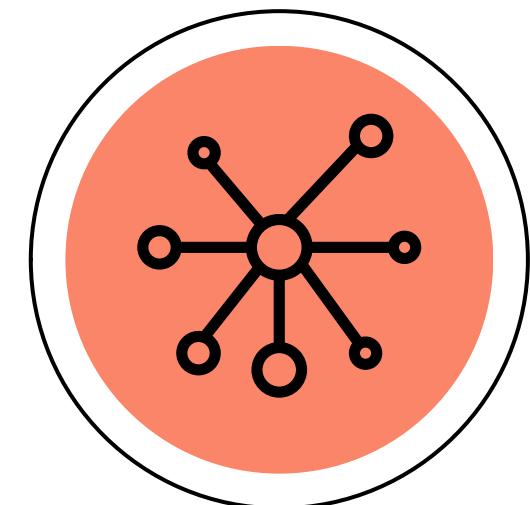
Clustering Coefficients

Measures the degree to which nodes in a graph tend to cluster together. It quantifies the likelihood that two neighbors of a node are connected, providing insights into the local clustering structure of the network.

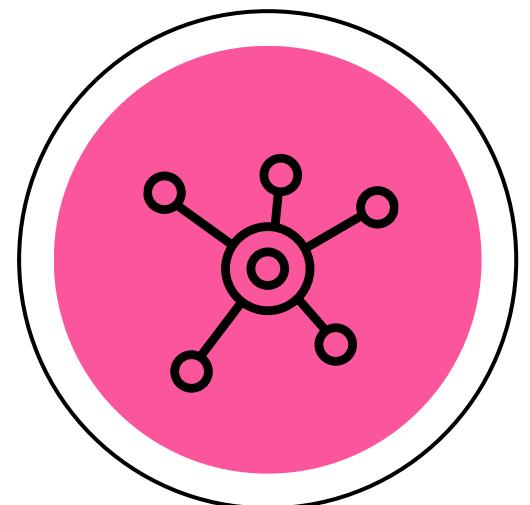
Sparsification Methods



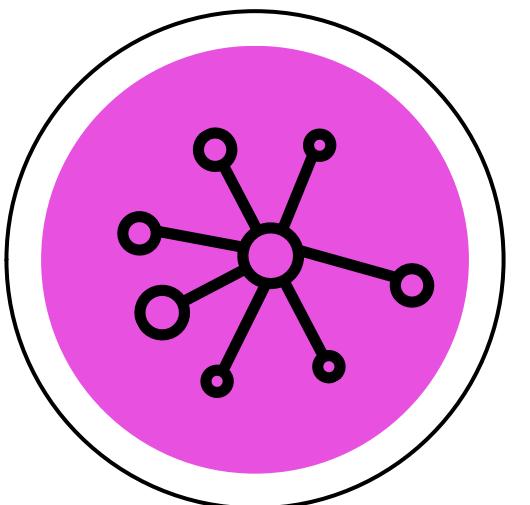
**Random
Sampling**



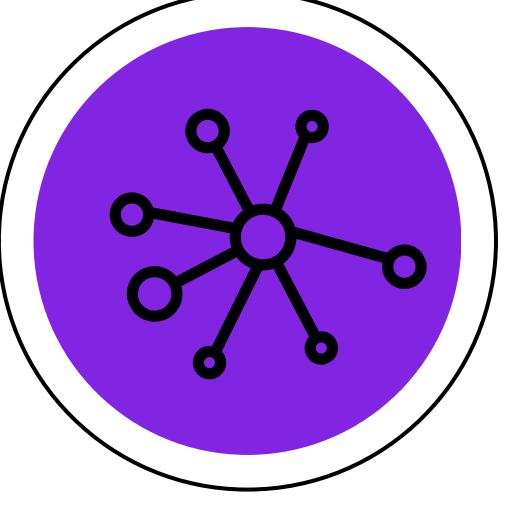
**Edge
Betweenness
Based**



**Global Similarity
Based**



**Local Similarity
Based**



**Clustering
Coefficient
Based**

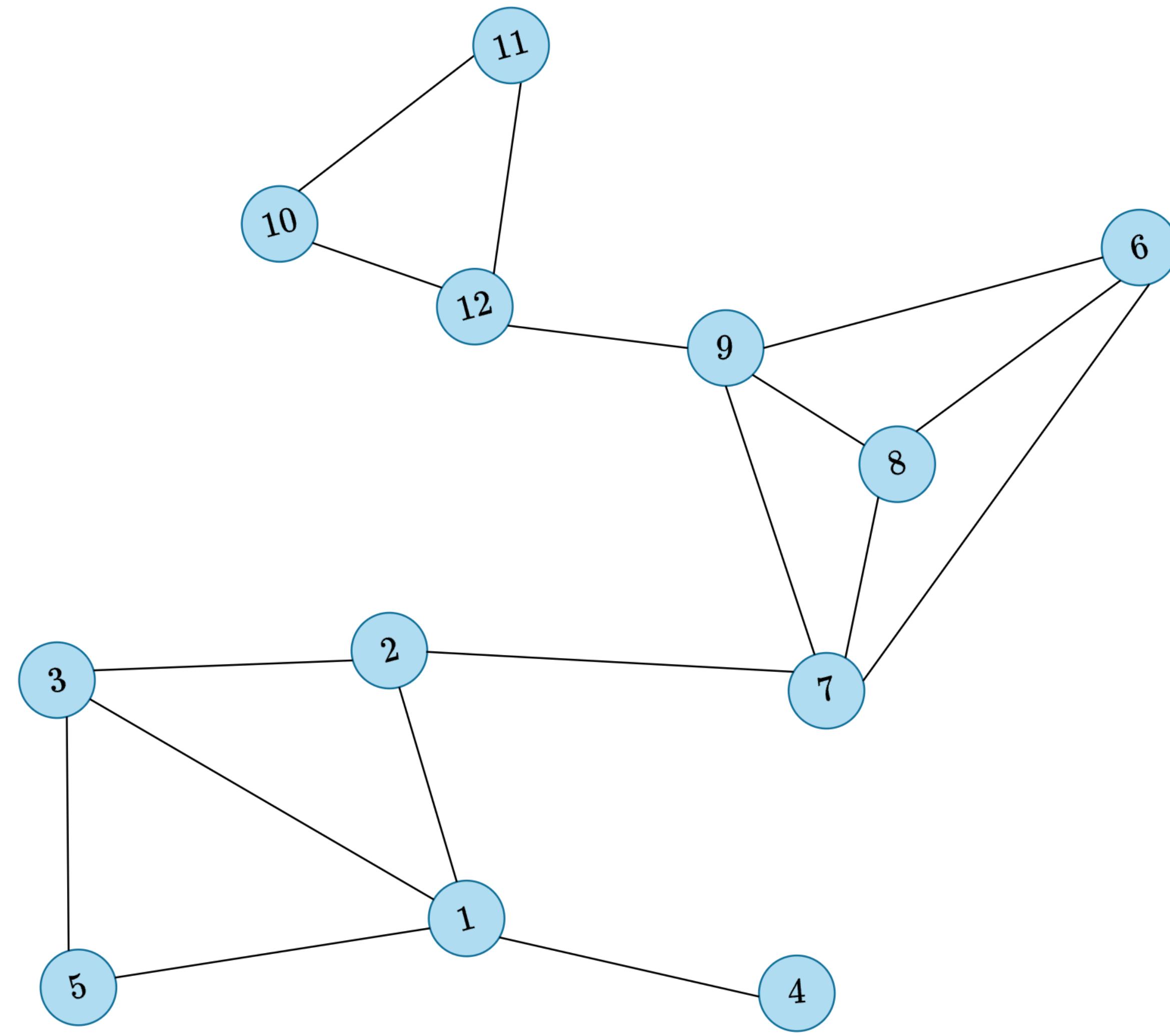
+ some more !

Edge Betweenness

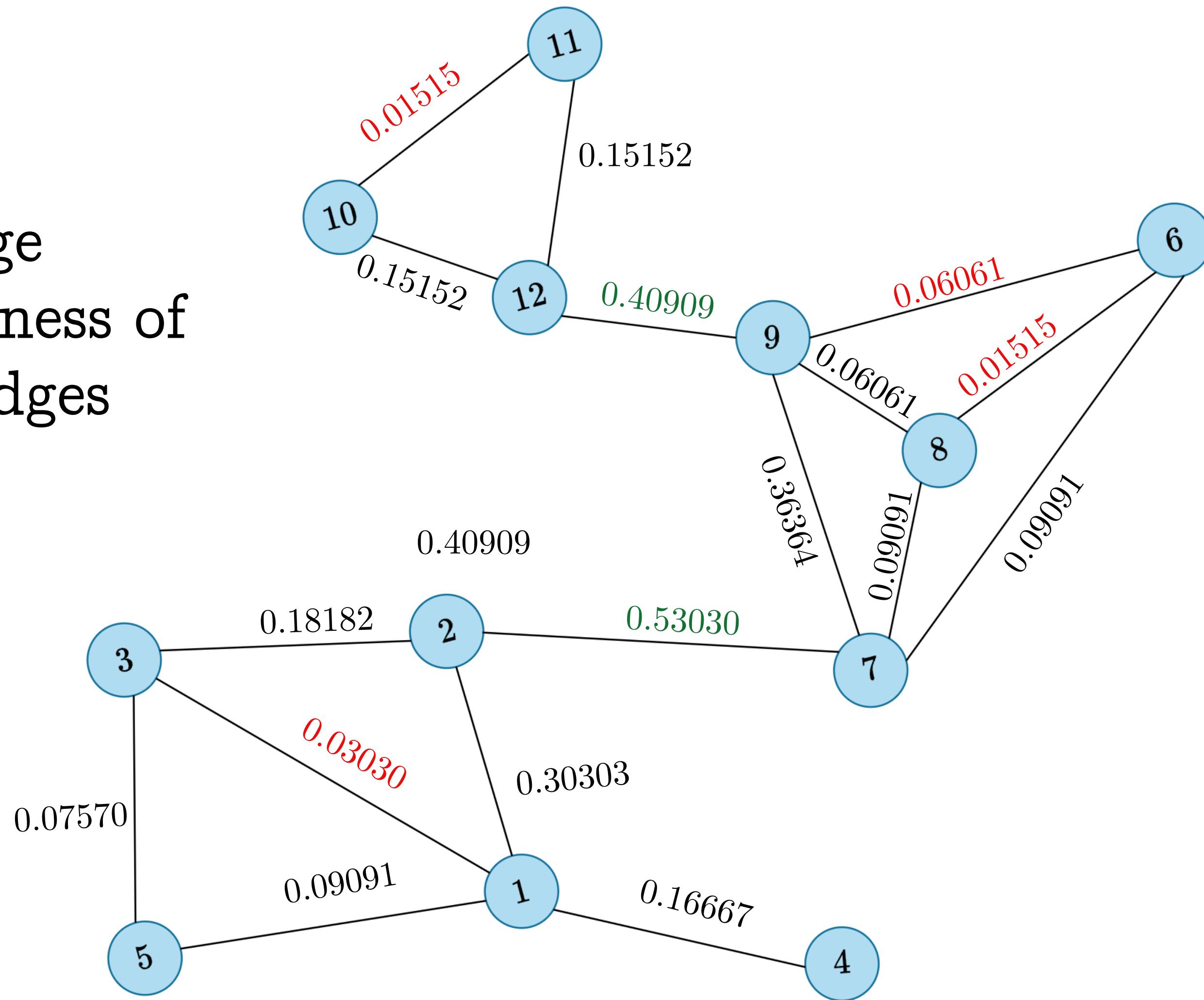
$$c_B(e) = \sum_{s,t \in V} \frac{\sigma(s, t|e)}{\sigma(s, t)} \quad \begin{aligned} \longrightarrow & \text{ No. of shortest } (s, t) \text{-paths crossing e} \\ \longrightarrow & \text{ Total No. of shortest } (s, t) \text{-paths} \end{aligned}$$

Time Complexity : $O(|V| \cdot |E|)$

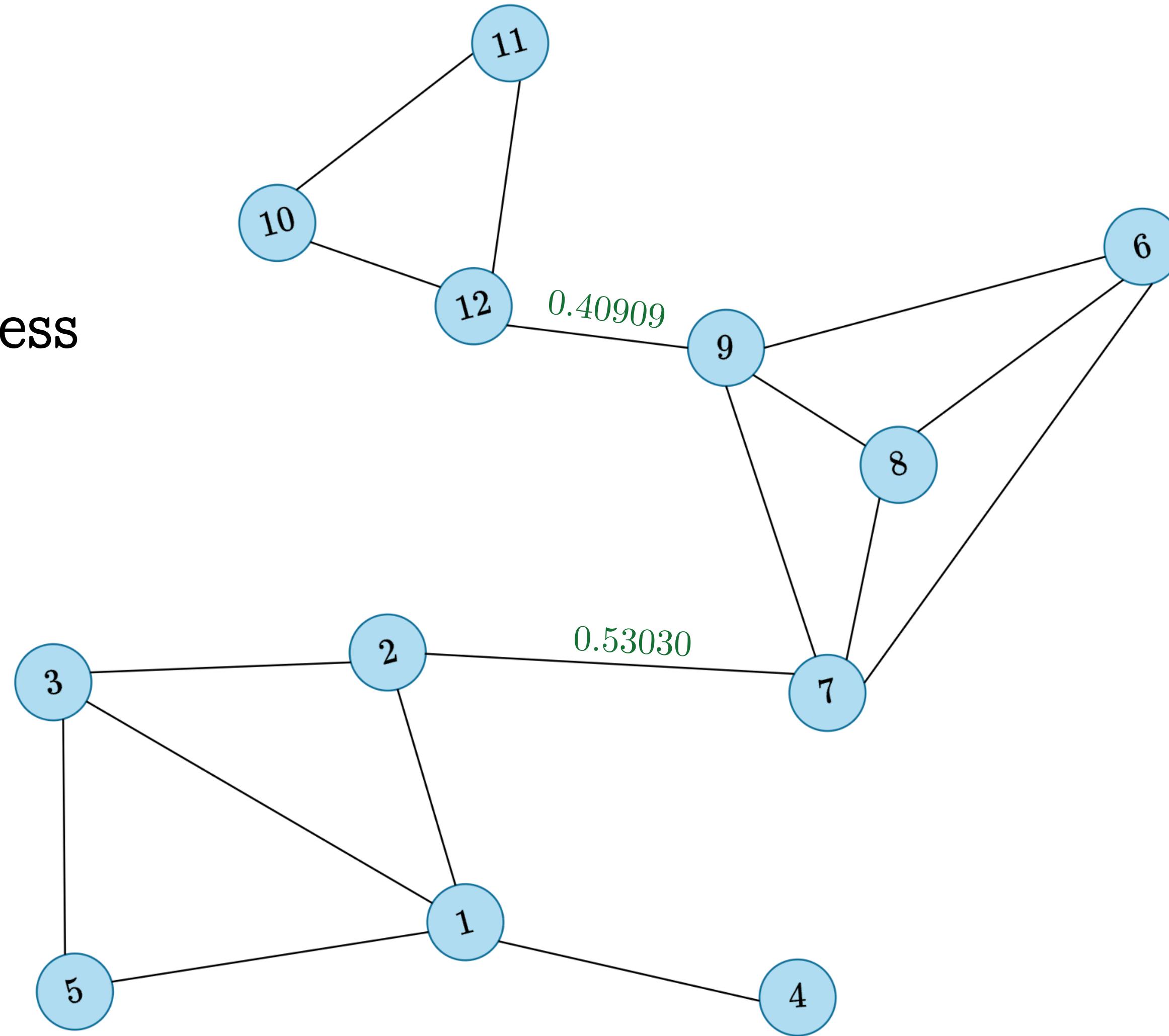
Space Complexity : $O(|V| + |E|)$

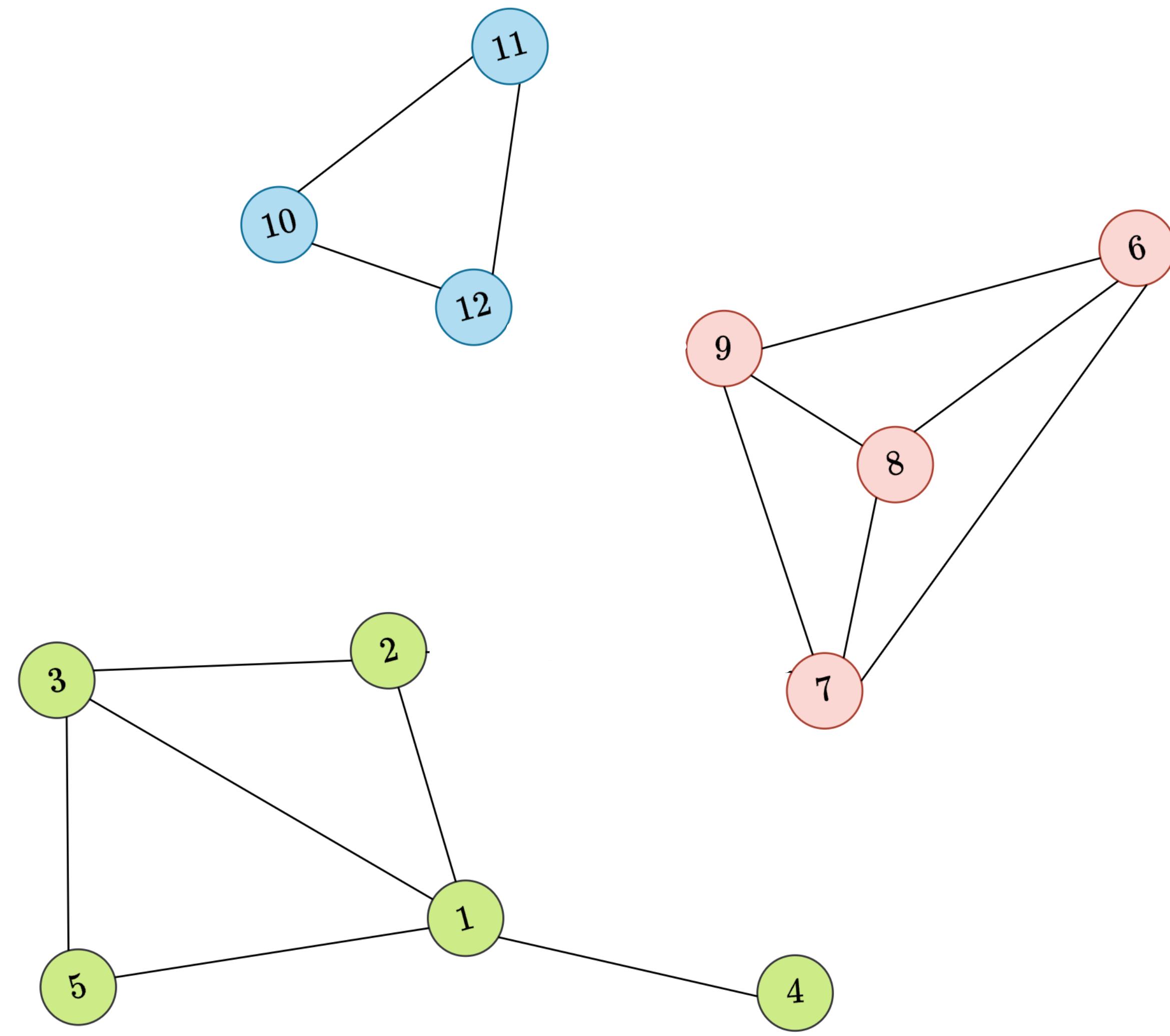


Edge Betweenness of the Edges

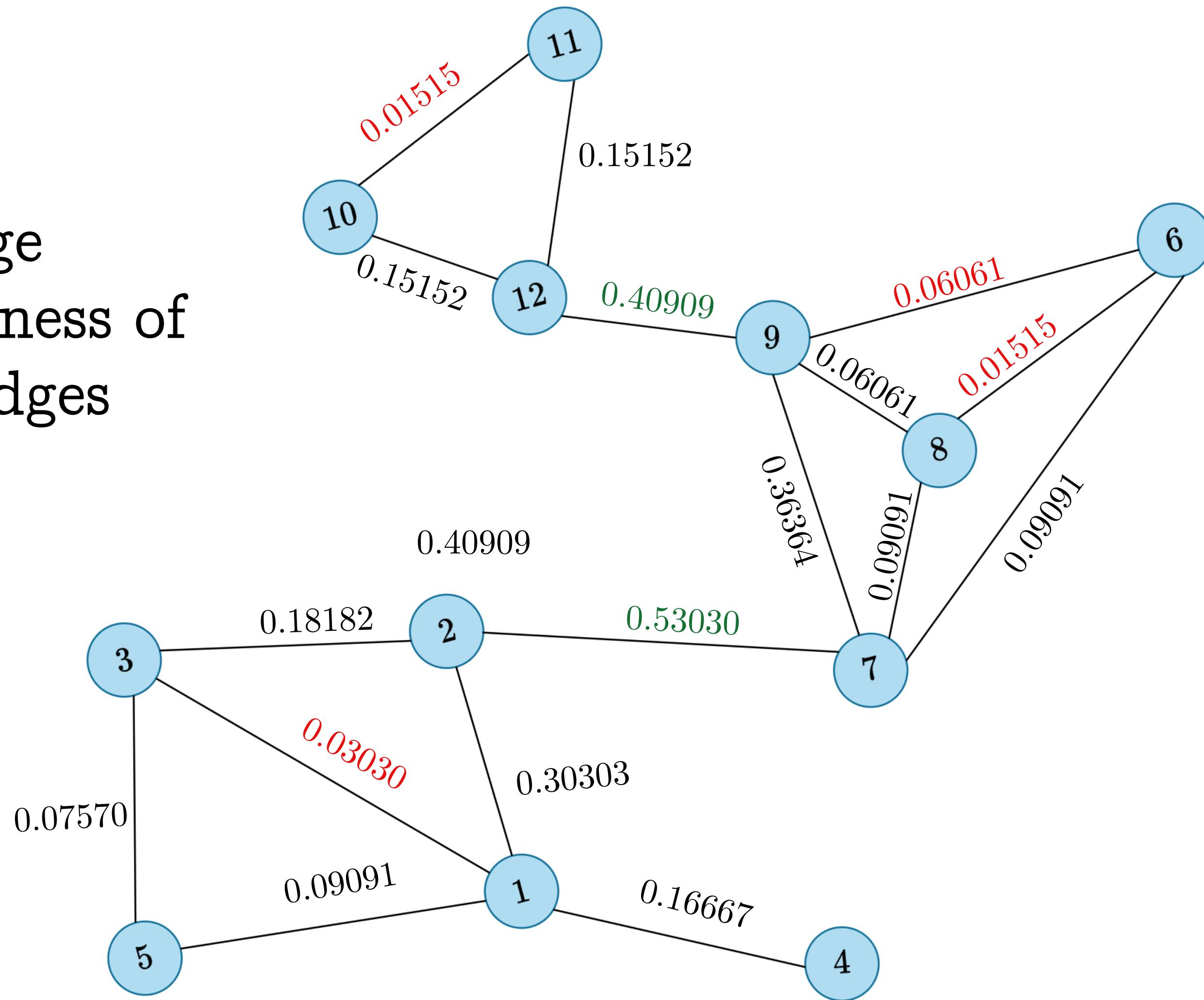


High
Betweenness
Edges

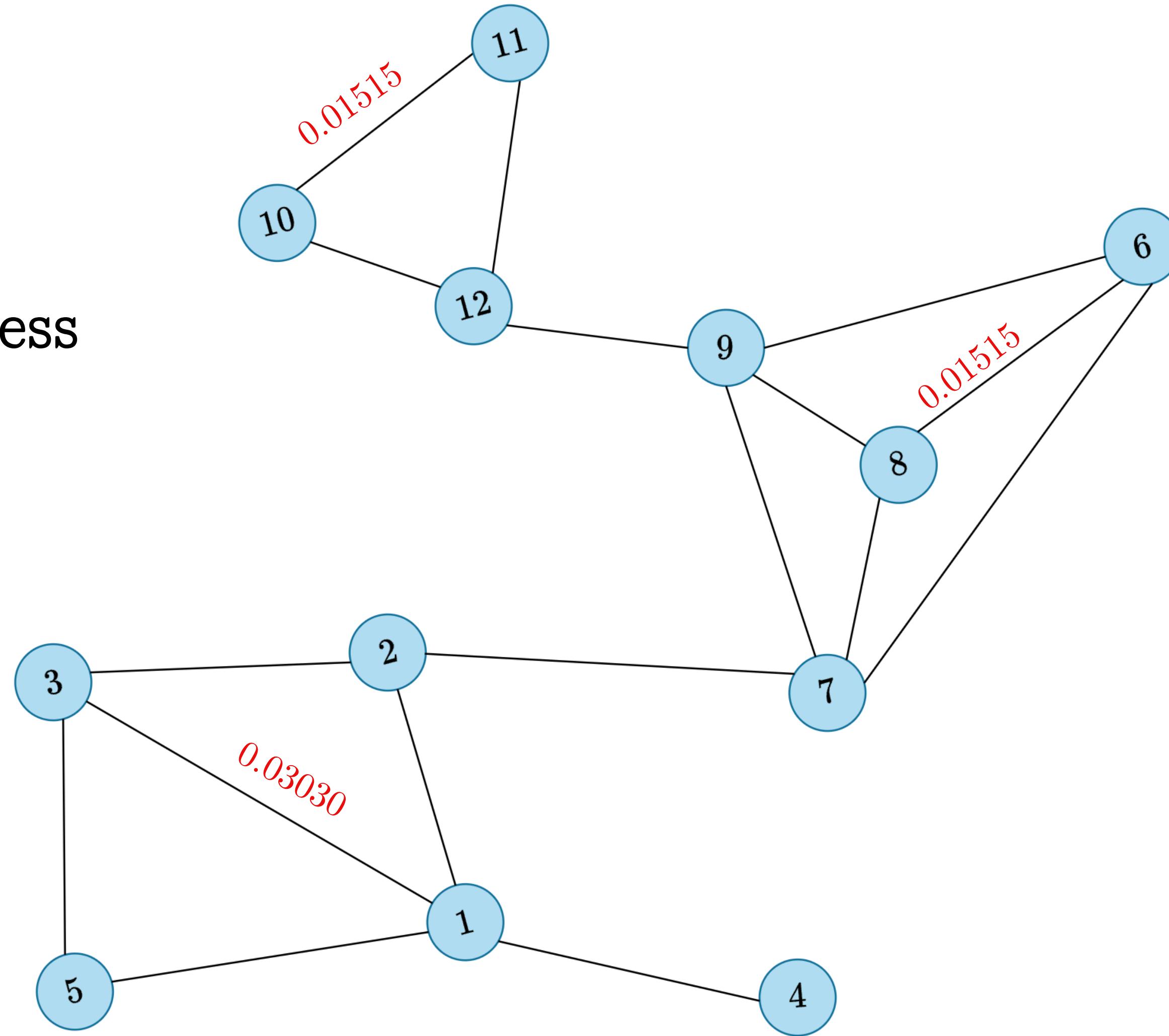




Edge Betweenness of the Edges

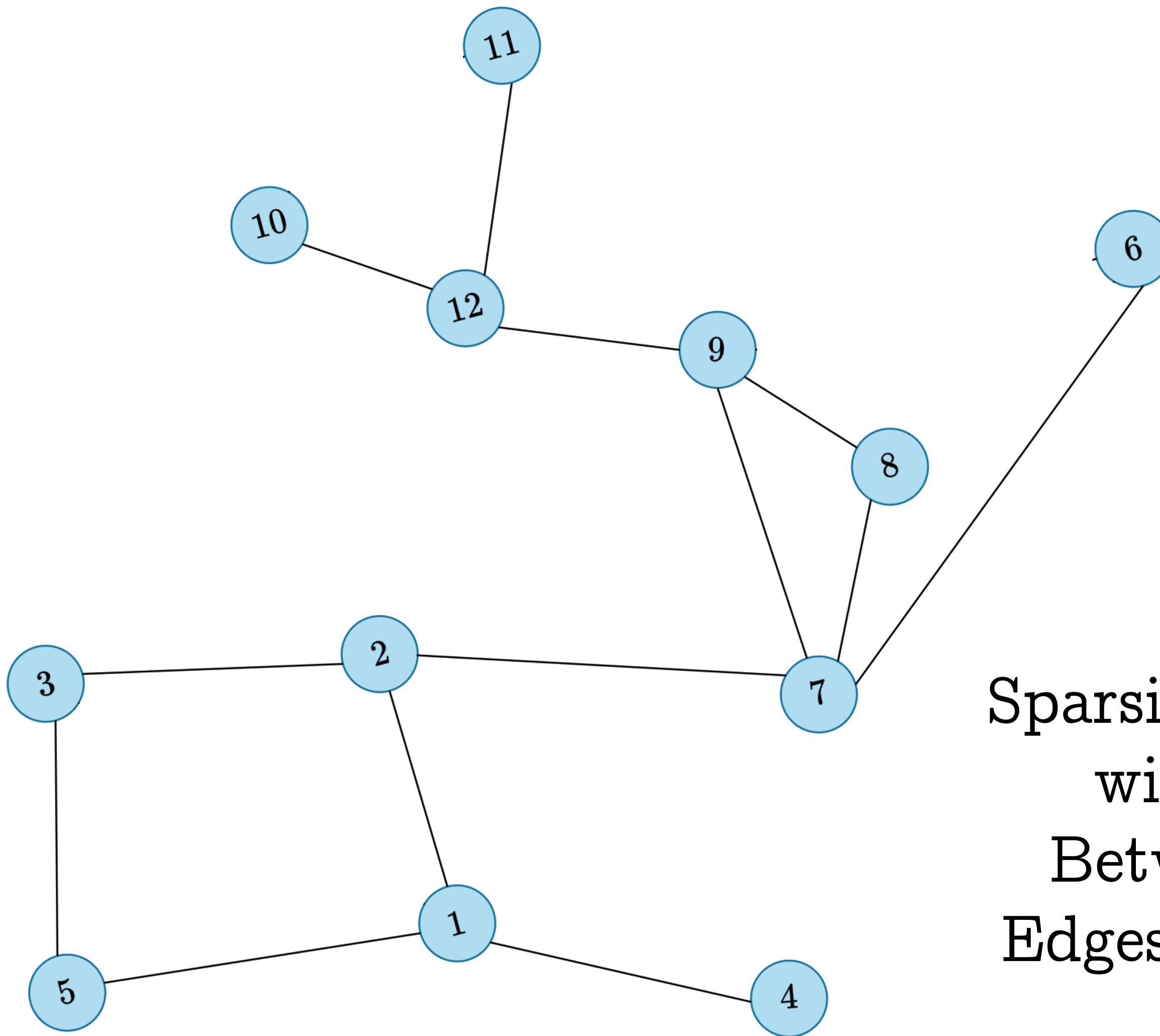


Low
Betweenness
Edges

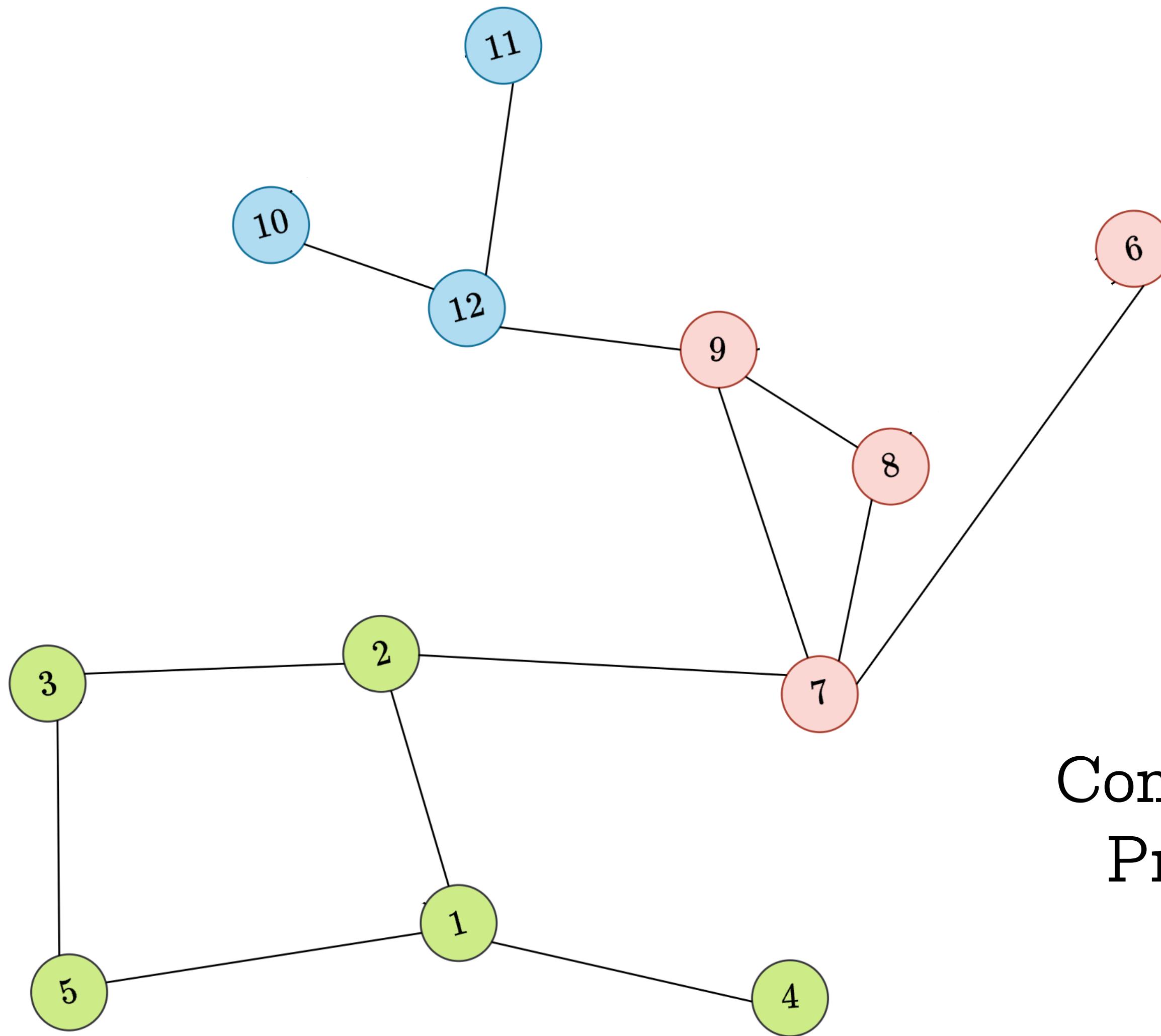


$$|V| = 12$$

$$|E| = 11$$

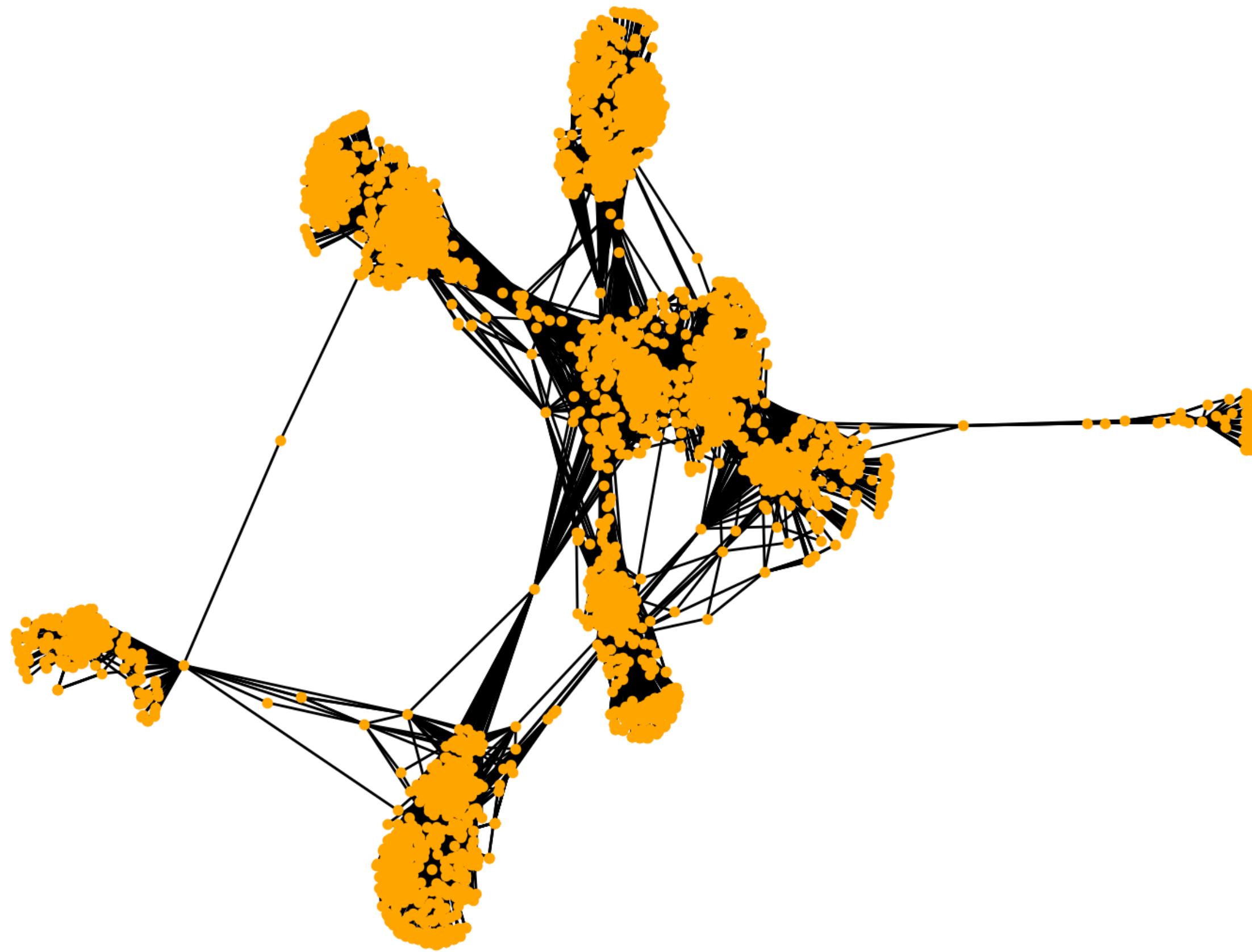


Sparsified Graph
with Low
Betweenness
Edges Removed



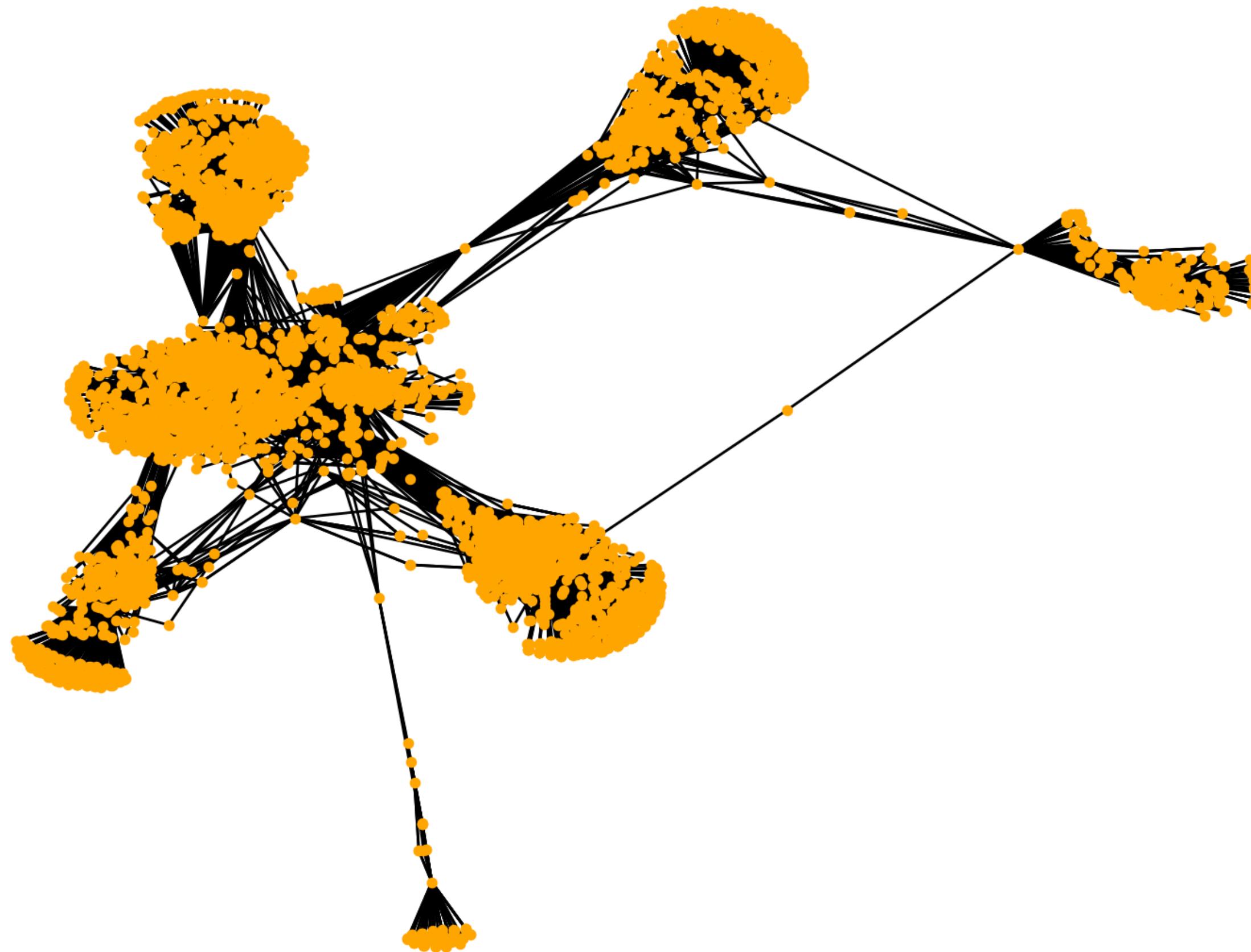
Communities
Preserved

Graph with Top 50.0 % Highest Betweenness Edges Retained
Number of edges: 44117



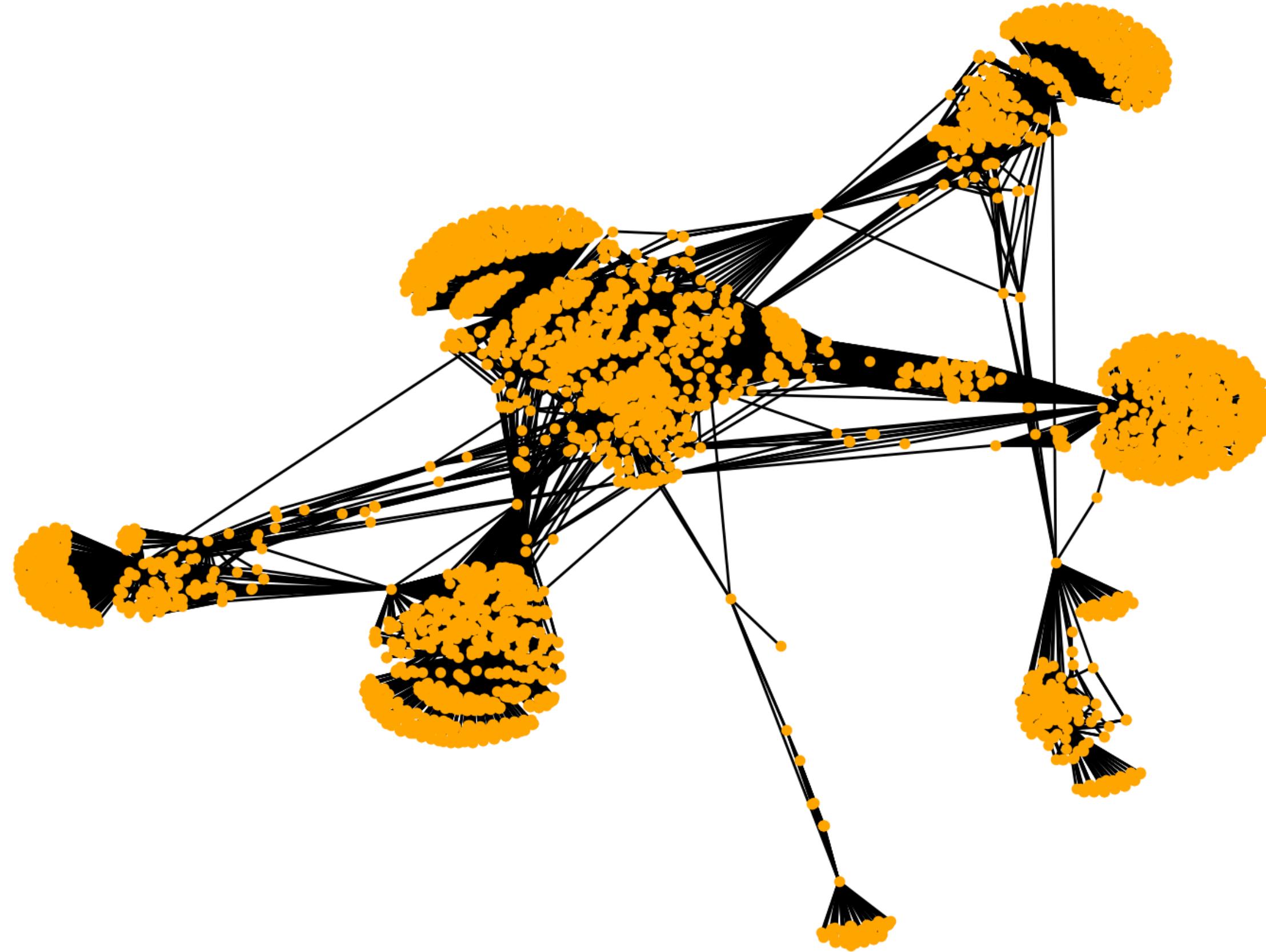
FaceBook

Graph with Top 30.0 % Highest Betweenness Edges Retained
Number of edges: 26471



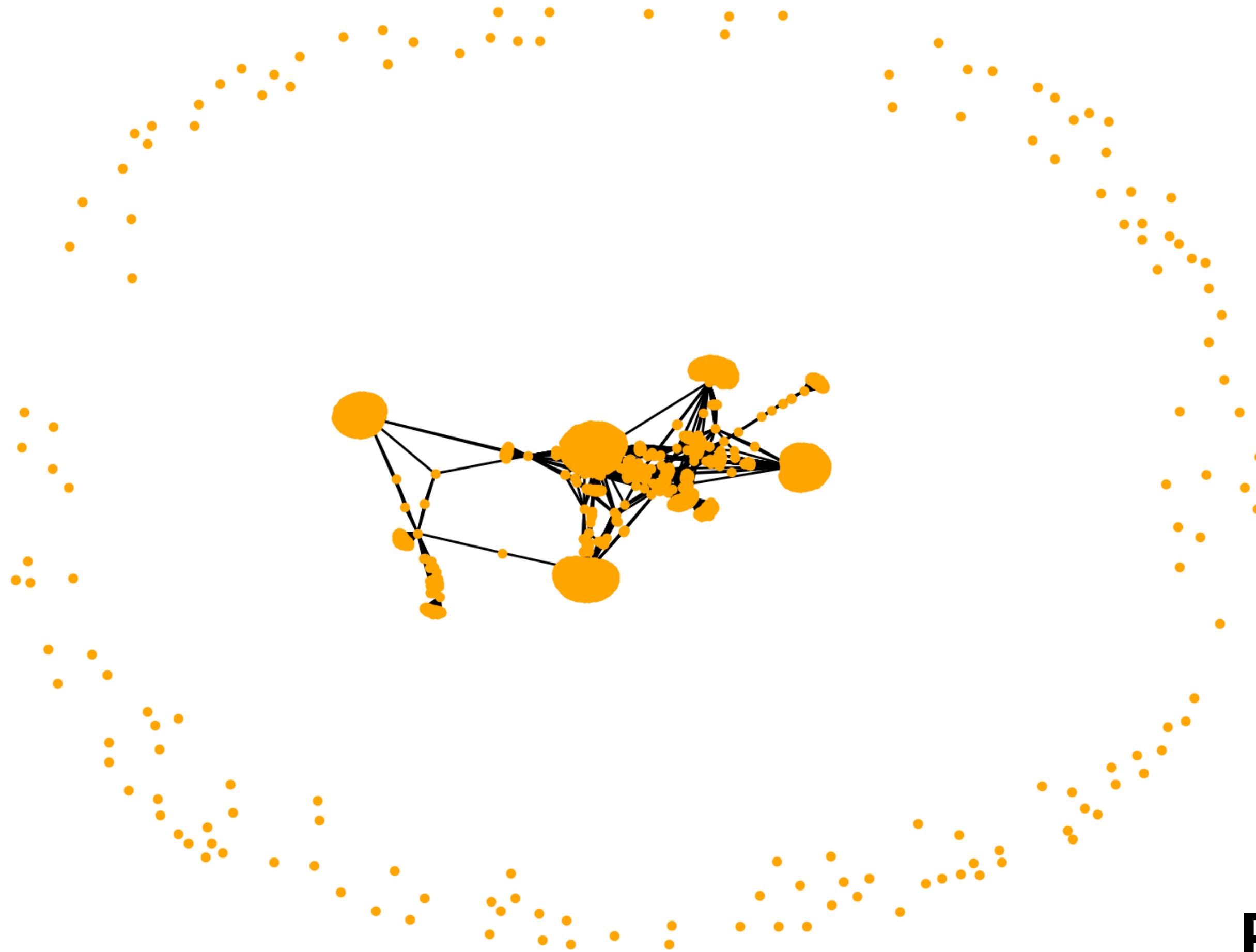
FaceBook

Graph with Top 10.0 % Highest Betweenness Edges Retained
Number of edges: 8824



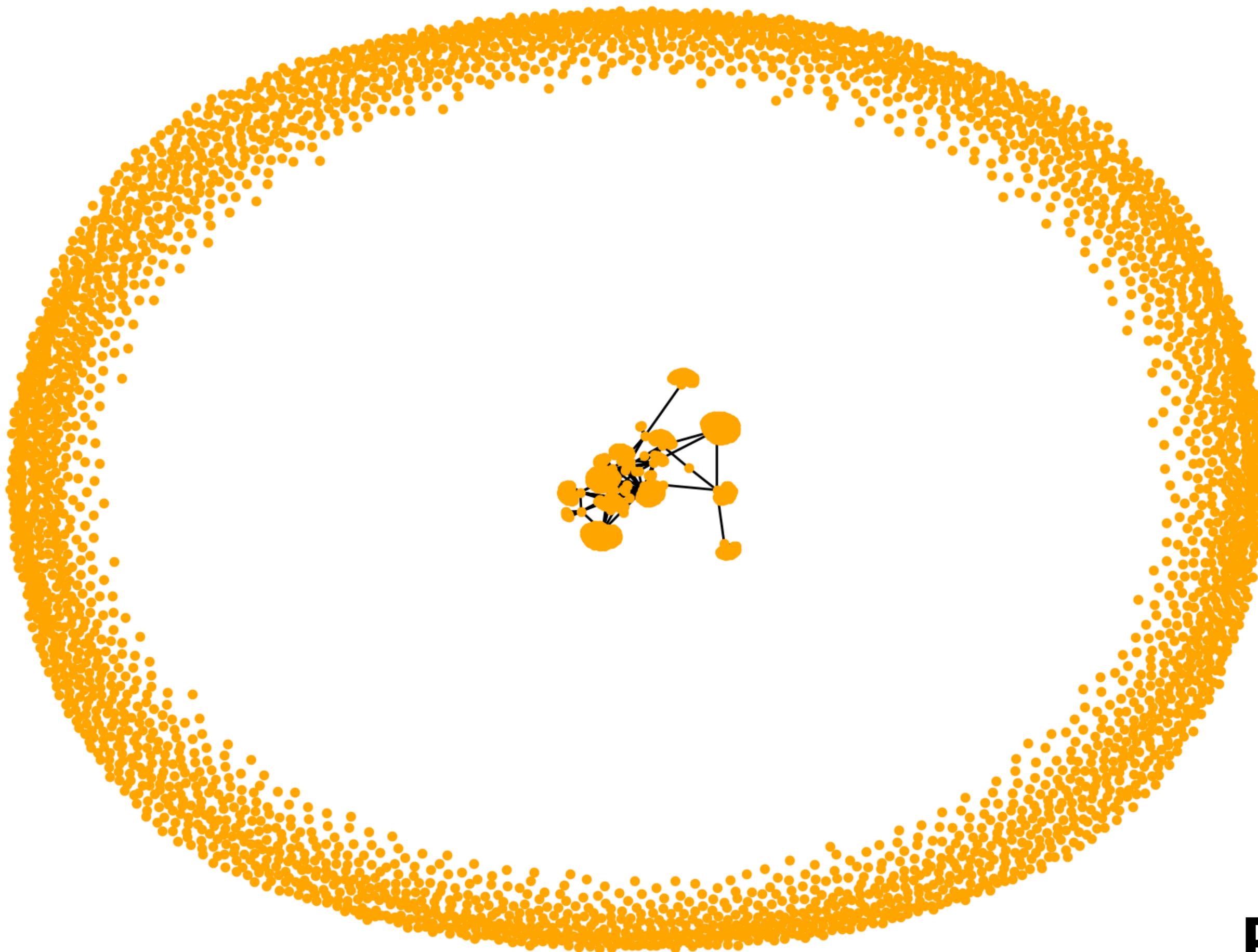
FaceBook

Graph with Top 5.0 % Highest Betweenness Edges Retained
Number of edges: 4412



FaceBook

Graph with Top 1.0 % Highest Betweenness Edges Retained
Number of edges: 883



FaceBook

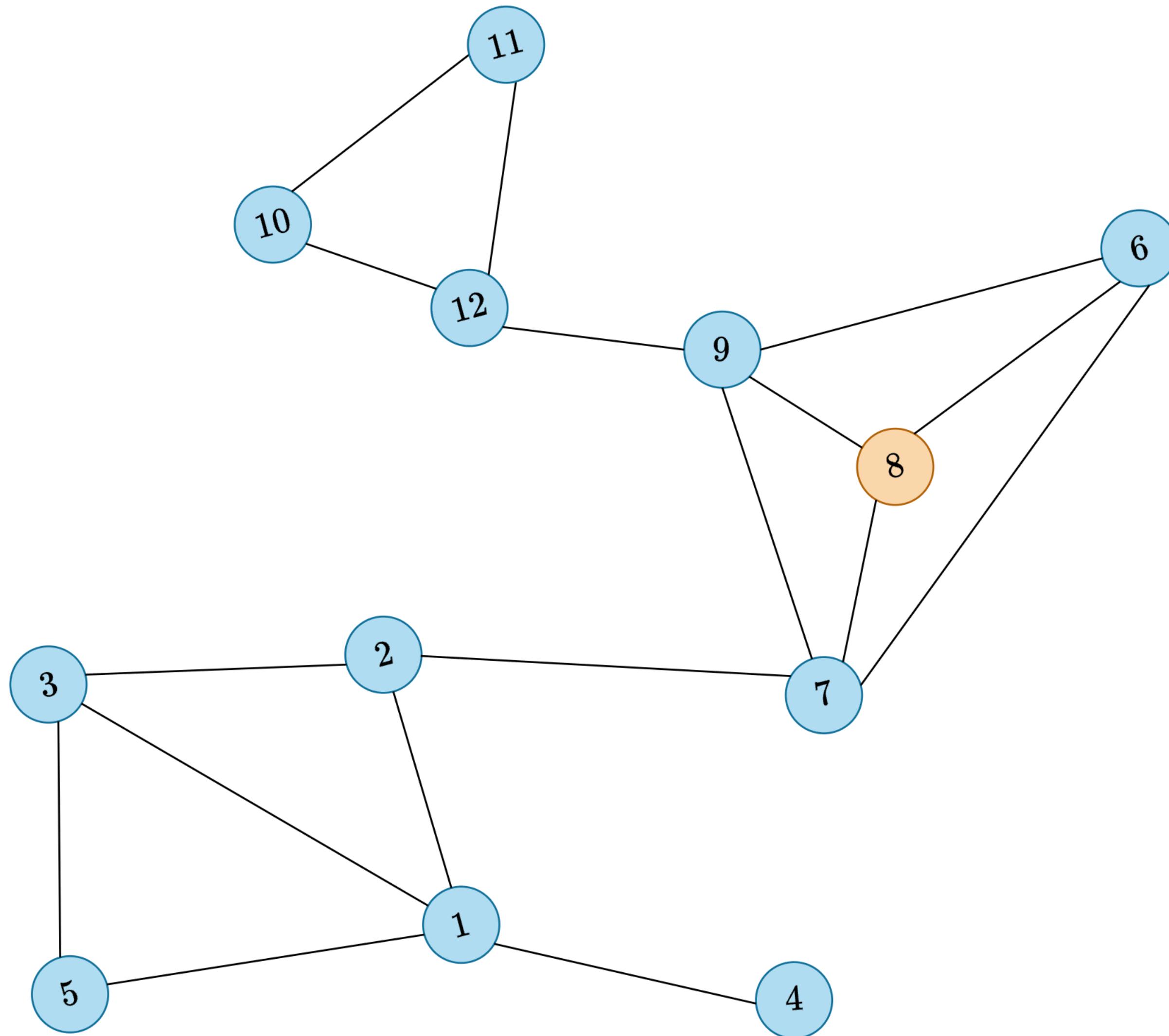
Jaccard Similarity

An edge (i, j) is *likely* to lie within a cluster if the vertices i and j have adjacency lists with *high overlap*.

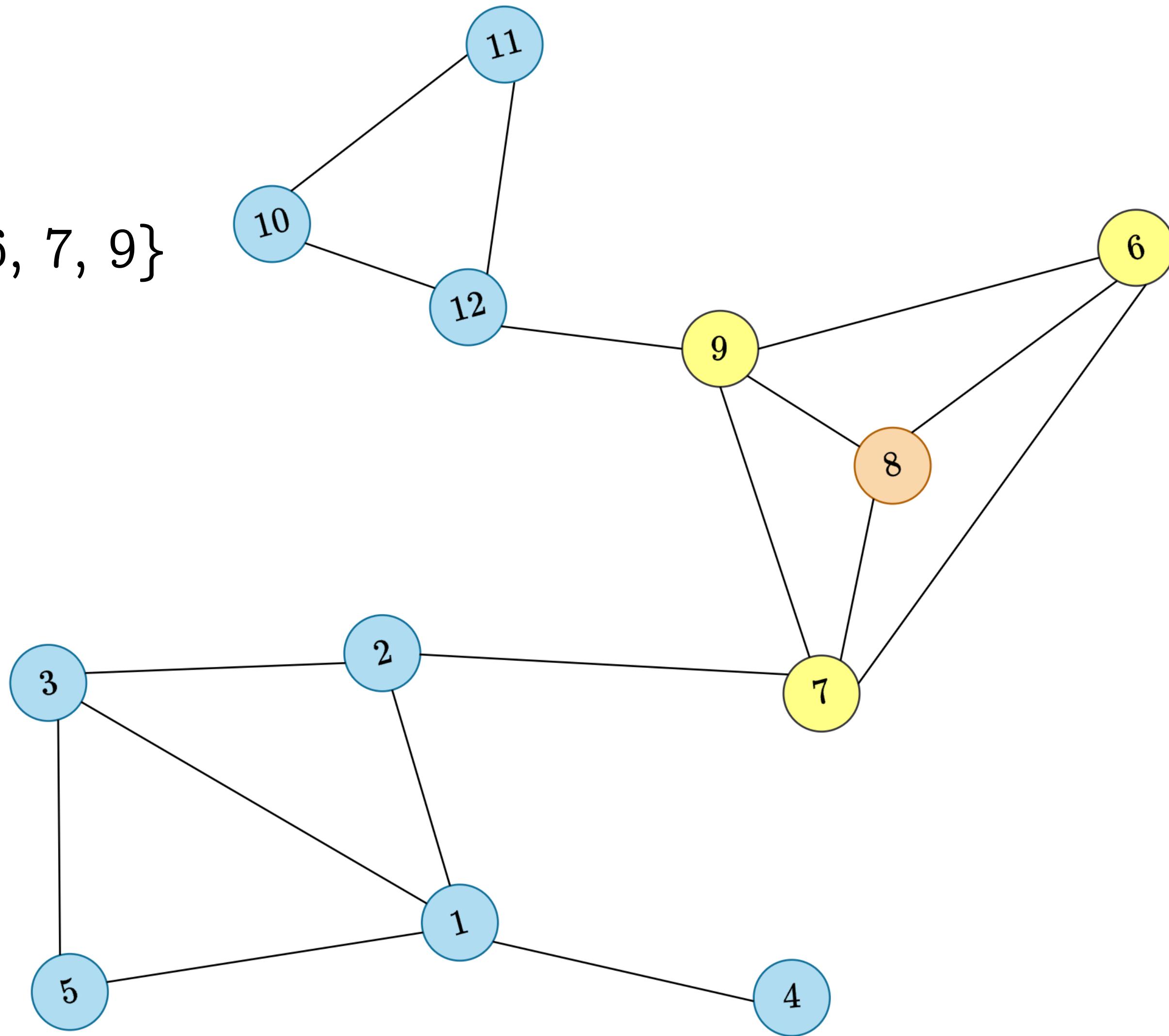
$$Sim(i, j) = \frac{|Adj(i) \cap Adj(j)|}{|Adj(i) \cup Adj(j)|}$$

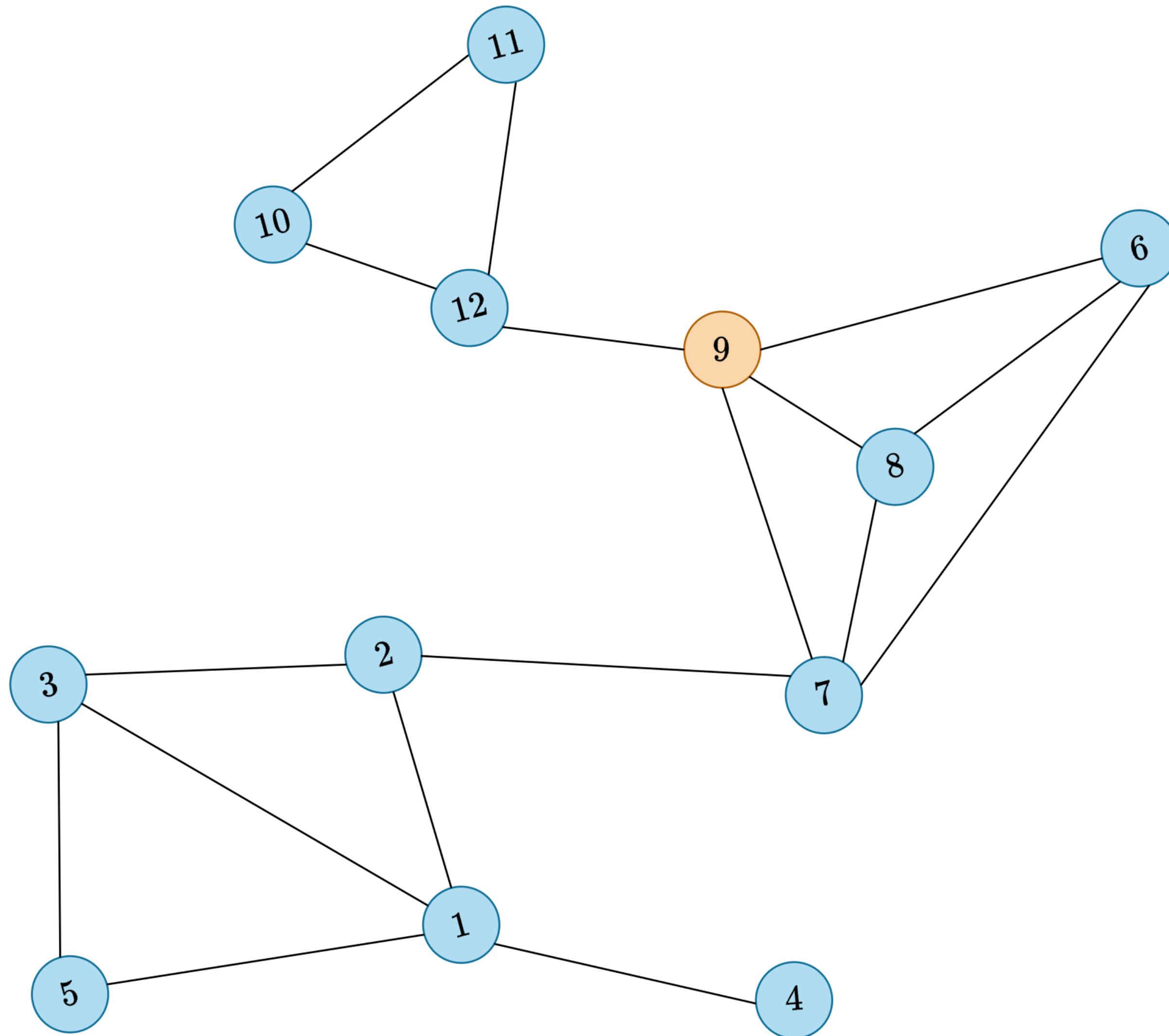


Adjacency List of Node i

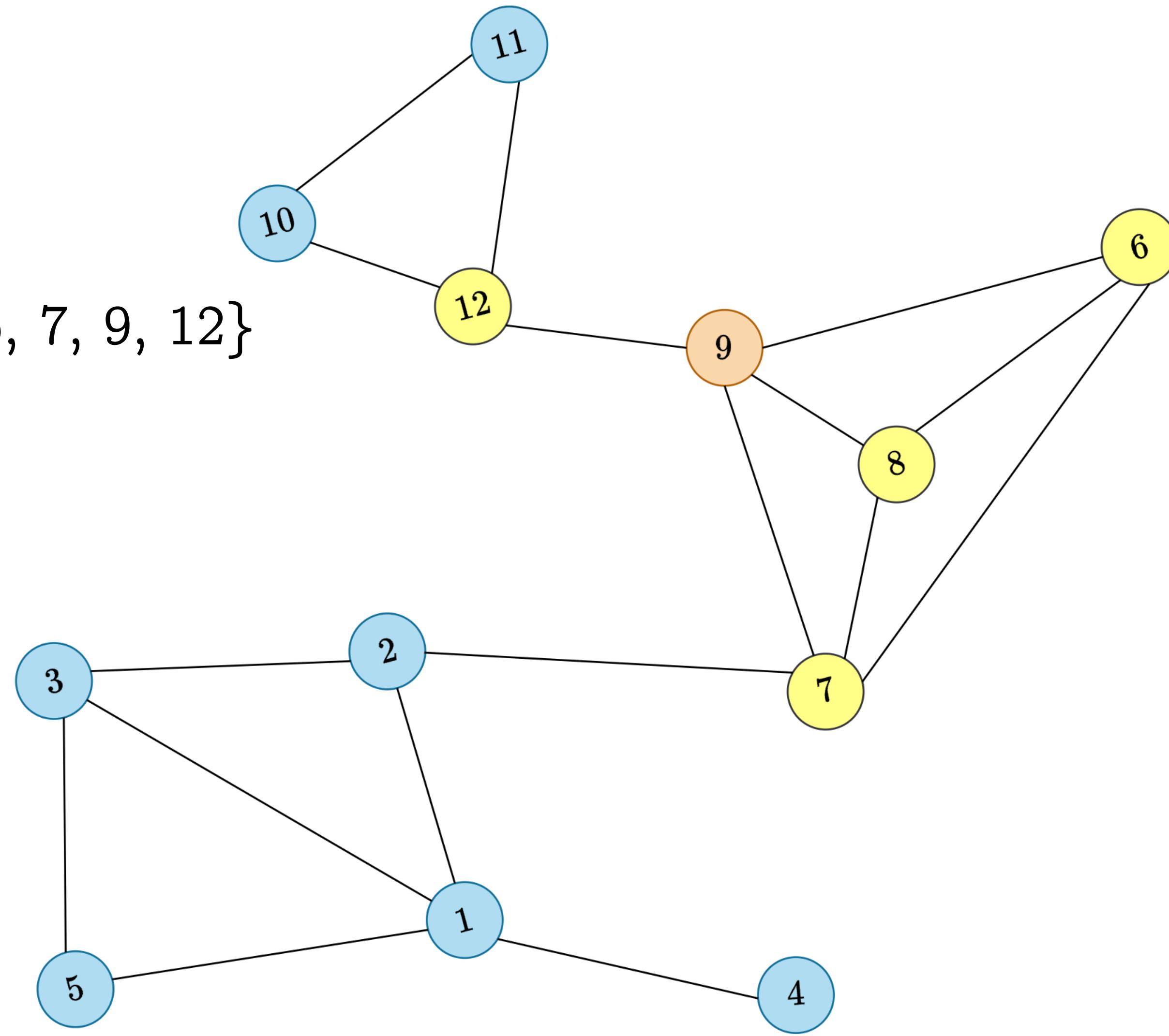


$$\text{Adj}(8) = \{6, 7, 9\}$$





$$\text{Adj}(9) = \{6, 7, 9, 12\}$$



$$\text{Adj}(8) = \{6, 7, 9\}$$

$$\text{Adj}(9) = \{6, 7, 9, 12\}$$

$$\text{Adj}(8) = \{6, 7, 9\}$$

$$\text{Adj}(9) = \{6, 7, 9, 12\}$$

$$\text{Sim}(8, 9) = \frac{|\{6, 7, 9\} \cap \{6, 7, 9, 12\}|}{|\{6, 7, 9\} \cup \{6, 7, 9, 12\}|}$$

$$\text{Adj}(8) = \{6, 7, 9\}$$

$$\text{Adj}(9) = \{6, 7, 9, 12\}$$

$$\text{Sim}(8, 9) = \frac{3}{4}$$

Algorithm 1 Global Sparsification Algorithm

Input: Graph $G = (V, E)$, Sparsification ratio s

```
 $G_{sparse} \leftarrow \emptyset$ 
for each edge  $e=(i,j)$  in  $E$  do
     $e.sim = Sim(i, j)$  according to Eqn 1
end for
```

Sort all edges in E by $e.sim$

Add the top $s\%$ edges to G_{sparse}

return G_{sparse}

Time Complexity : $O(|E|(m+\log|E|))$
 m is size of largest adjacency list.

Space Complexity : $O(|V|+|E|)$

Algorithm 2 Local Sparsification Algorithm

Input: Graph $G = (V, E)$, Local Sparsification exponent e

Output: Sparsified graph G_{sparse}

```
 $G_{sparse} \leftarrow \emptyset$ 
for each node  $i$  in  $V$  do
    Let  $d_i$  be the degree of  $i$ 
    Let  $E_i$  be the set of edges incident to  $i$ 
    for each edge  $e=(i,j)$  in  $E_i$  do
         $e.sim = Sim(i, j)$  according to Eqn. 1
    end for
    Sort all edges in  $E_i$  by  $e.sim$ 
    Add top  $d_i^e$  edges to  $G_{sparse}$ 
end for

return  $G_{sparse}$ 
```

Time Complexity : $O(|E|/(m+\log|E|))$
 m is size of largest adjacency list.

Space Complexity : $O(|V|+|E|)$

Modularity

$$Q(\mathcal{C}) = \frac{1}{2m} \sum_{C \in \mathcal{C}} \sum_{u \in C, v \in C} \left(A_{u,v} - \frac{d_u d_v}{2m} \right)$$

↑
Expected Edges b/w u, v

↓
Edges b/w u, v

Clustering Coefficients

$$c_u = \frac{\# \text{ triangles through node } u}{deg(u)(deg(u)) - 1}$$

The diagram illustrates the components of the clustering coefficient formula. An upward arrow points from the term $2T(u)$ to the text "# triangles through node u". A downward arrow points from the variable c_u to the symbol \in , which is followed by the letter E .

Networks Used

1.

Amazon Co-Purchase

2.

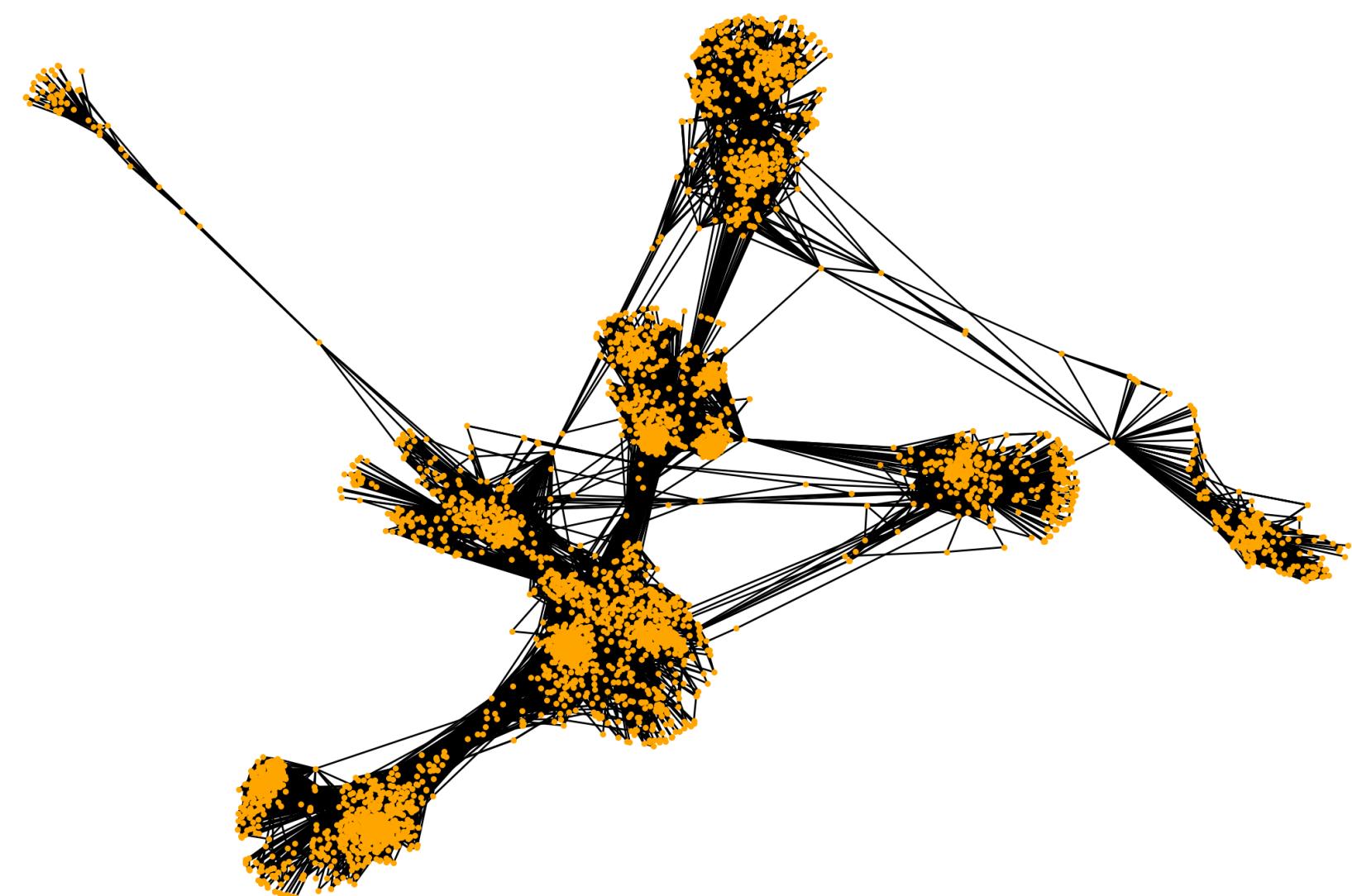
DBLP

3.

Facebook Social

4.

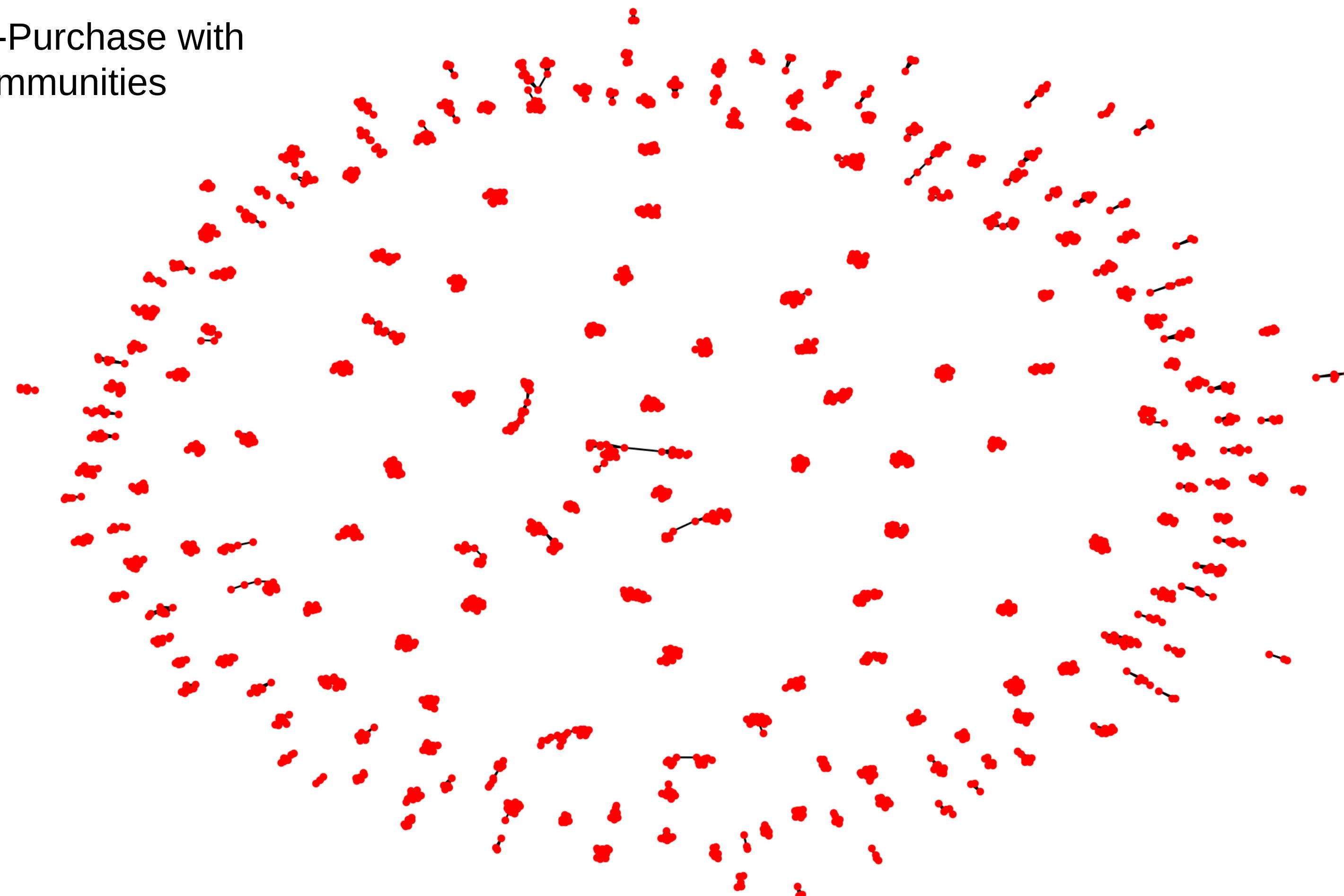
Email EU Core



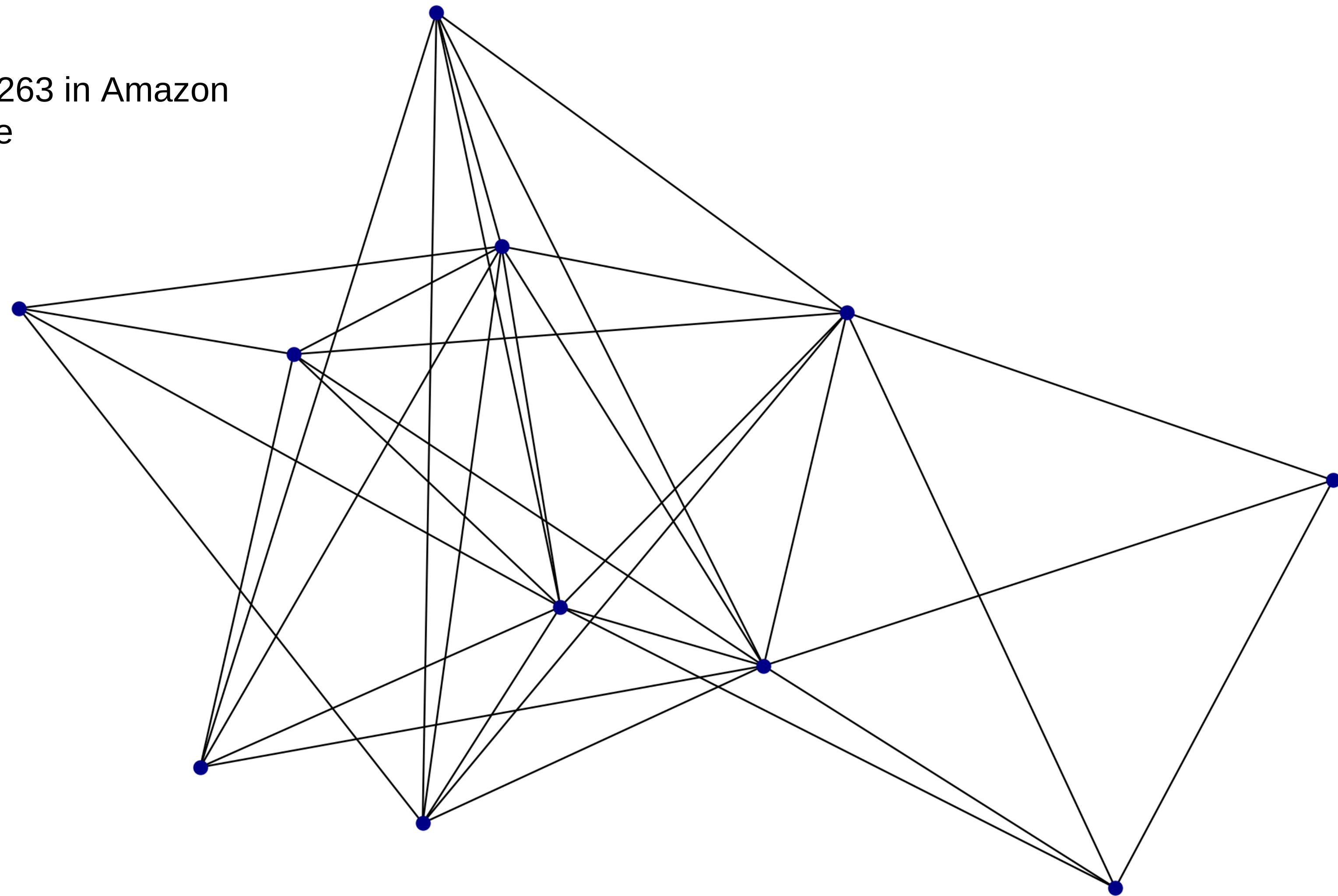
FaceBook Network

Network	 V 	 E 	 C 	 V _{ind}	 E _{ind}
DBLP	317080	1049866	150	1420	4609
Amazon	334863	925872	300	2008	5960
EU	1005	16064	42	1005	16064
Facebook	4039	88234	NA	4039	88234

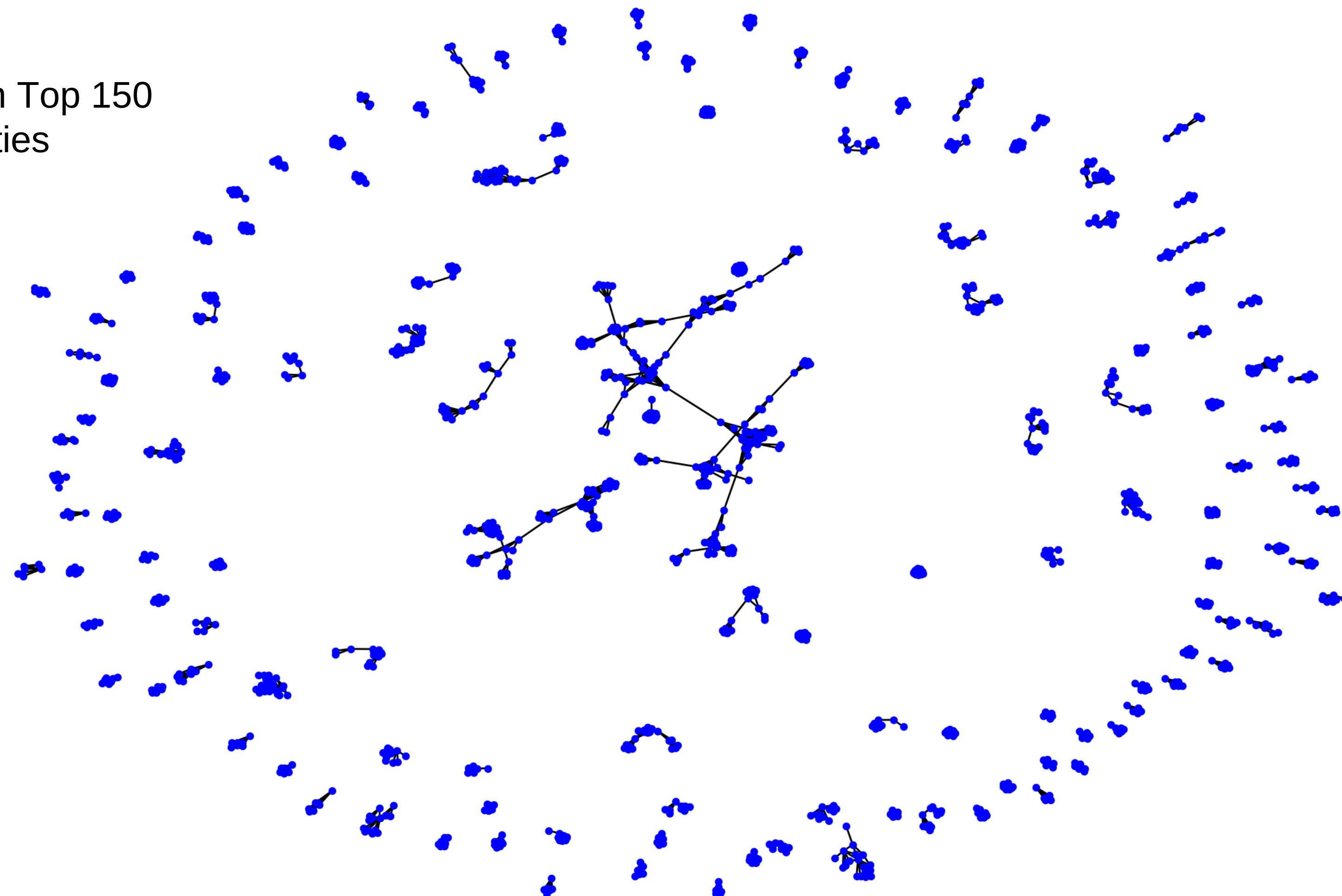
Amazon Co-Purchase with Top 300 Communities



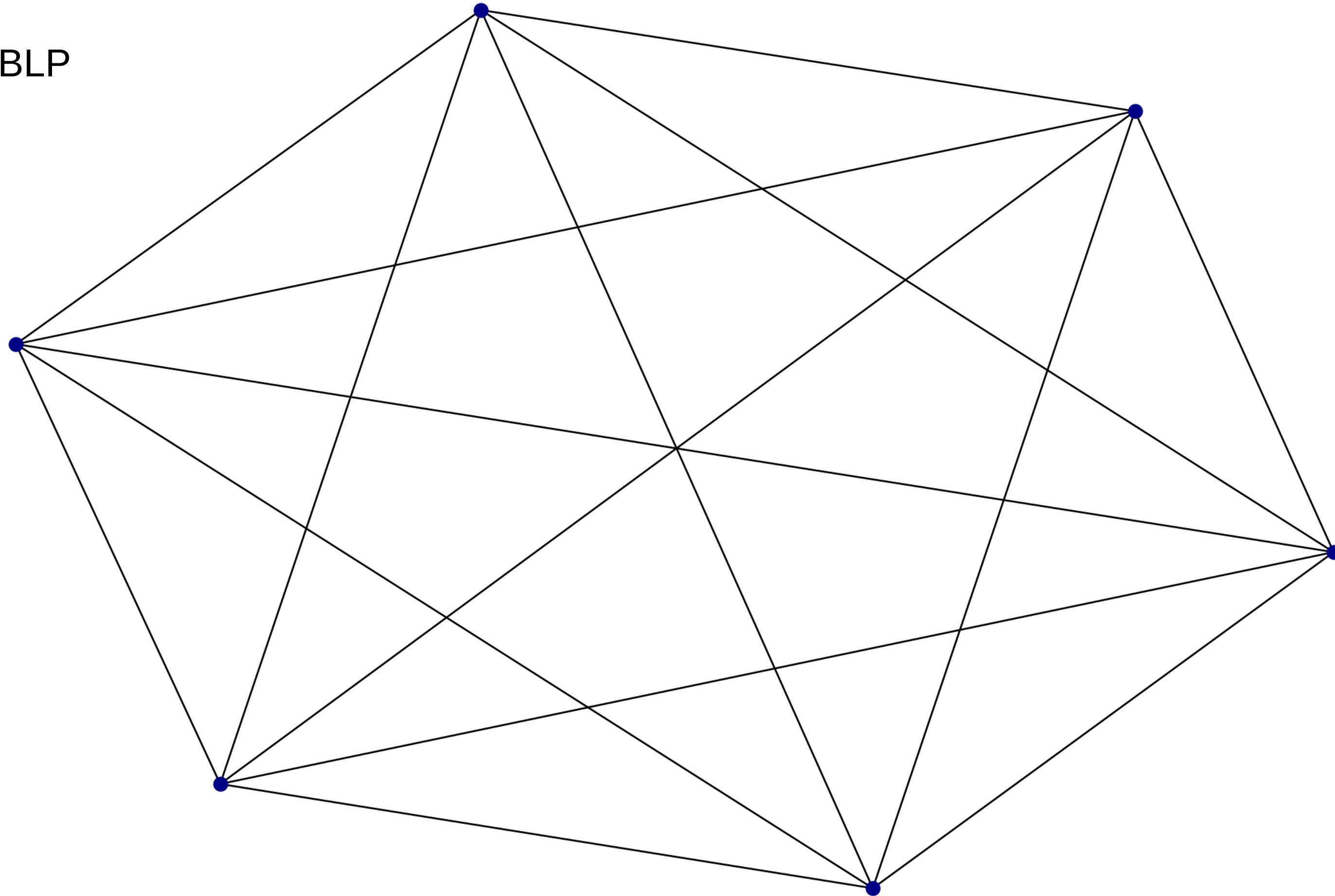
Community 263 in Amazon
Co-Purchase



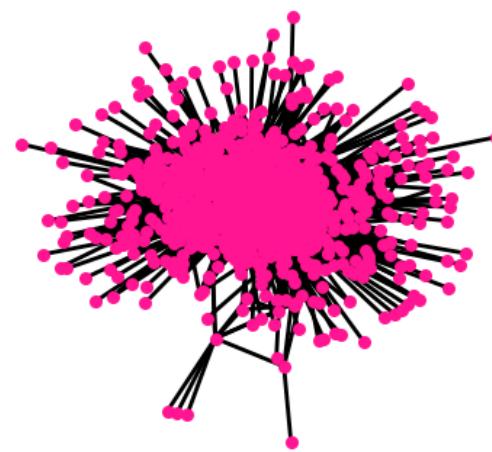
DBLP with Top 150
Communities



Community 4 in DBLP



Email EU Core Network



Community Detection Algorithms Used

1.

Louvain Algorithm

A greedy algorithm to extract the community structure of a network based on modularity optimization.

2.

Label Propogation Algorithm

The algorithm detects communities using network structure alone. At every iteration of propagation, each node updates its label to the one that the maximum numbers of its neighbours belongs to

3.

InfoMap Algorithm

The algorithm uses the probability flow of random walks on a network as a proxy for information flows in the real system and decomposes the network into modules by compressing a description of the probability flow

Metrics for comparing Community Quality

1.

Adjusted Rand Index

$$ARI = \frac{\sum_{ij} \binom{n_{ij}}{2} - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}] / \binom{n}{2}}{\frac{1}{2} [\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}] - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}] / \binom{n}{2}}$$

2.

Normalized Mutual Information

$$NMI(X, Y) = \frac{-2 \sum_{ij} p_{ij} \log(p_{ij}/p_{i+}p_{+j})}{\sum_i p_{i+} \log(p_{i+}) + \sum_j p_{+j} \log(p_{+j})}$$

3.

Modularity

$$Q(\mathcal{C}) = \frac{1}{2m} \sum_{C \in \mathcal{C}} \sum_{u \in C, v \in C} \left(A_{u,v} - \frac{d_u d_v}{2m} \right)$$

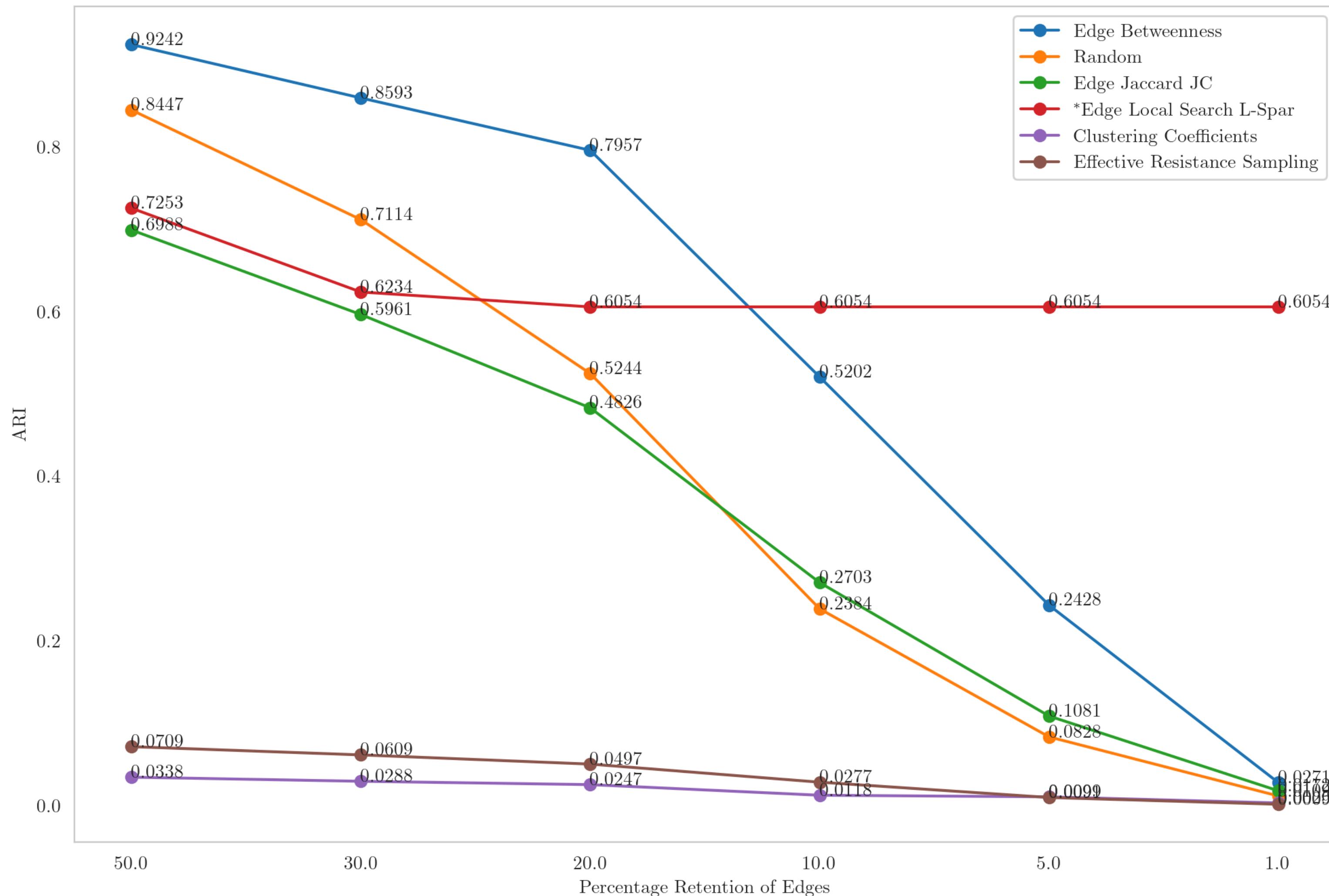
4.

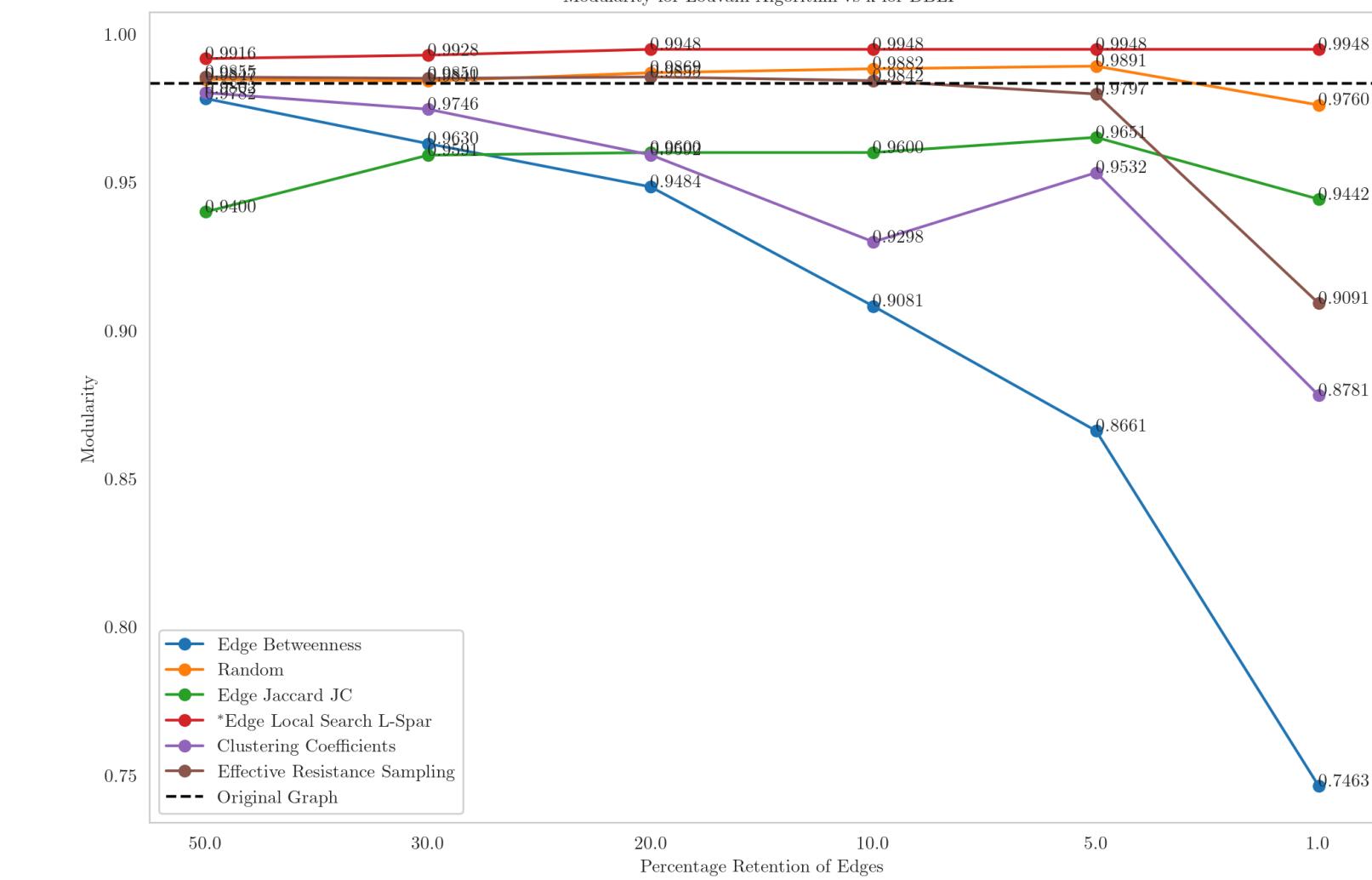
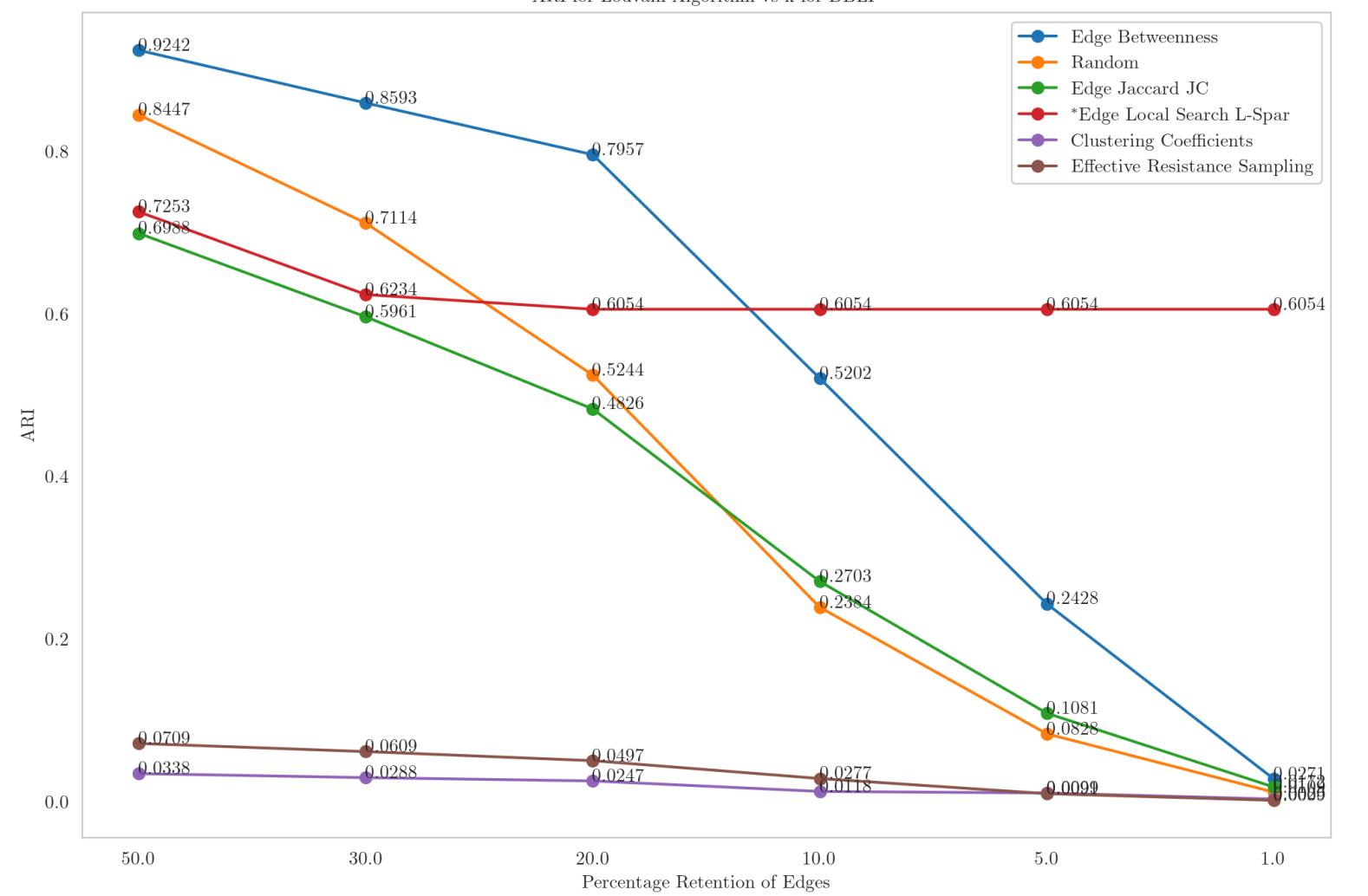
Clustering Coefficient

$$c_u = \frac{2T(u)}{deg(u)(deg(u)) - 1}$$

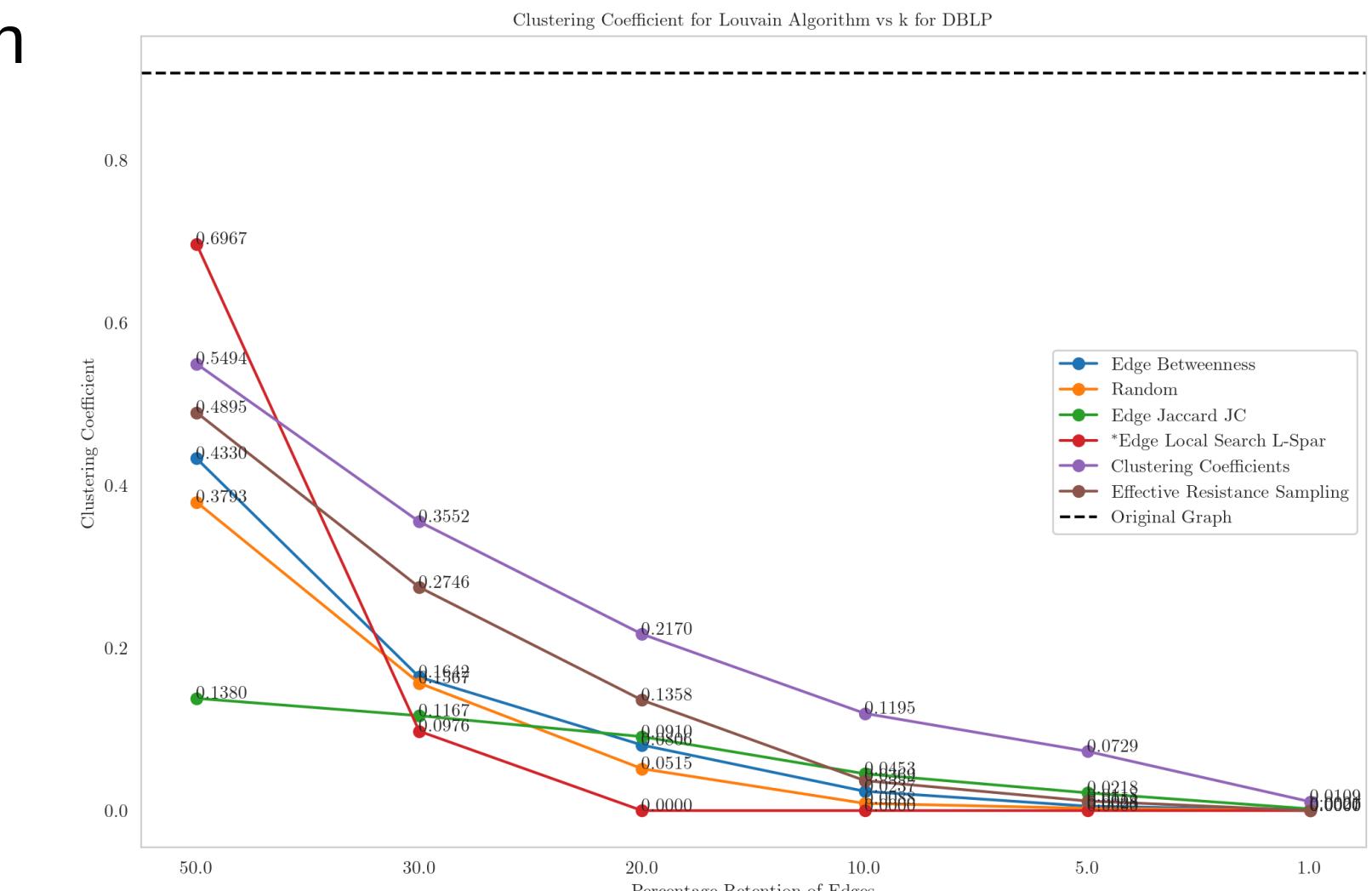
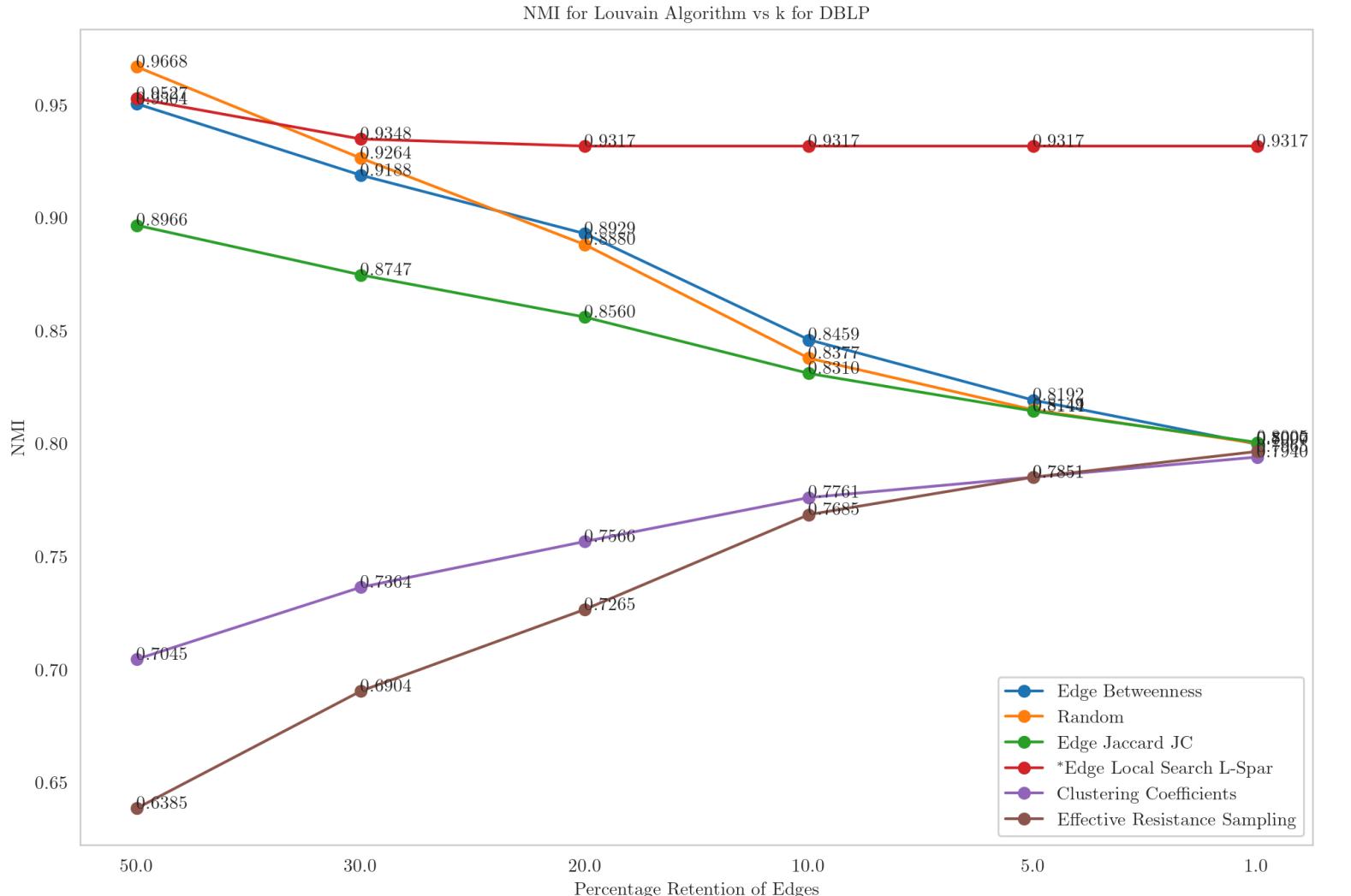
Results for DBLP Network

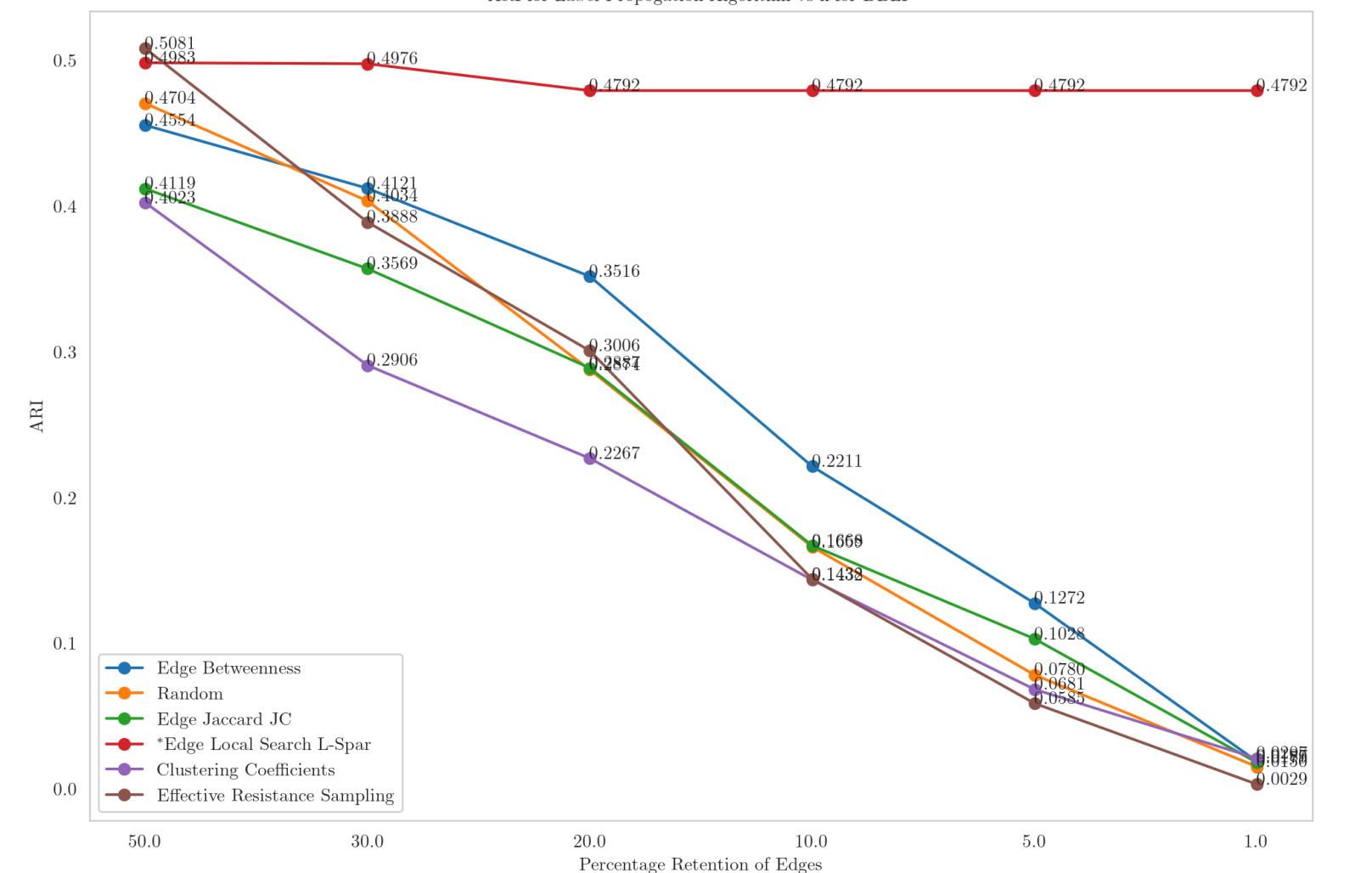
ARI for Louvain Algorithm vs k for DBLP



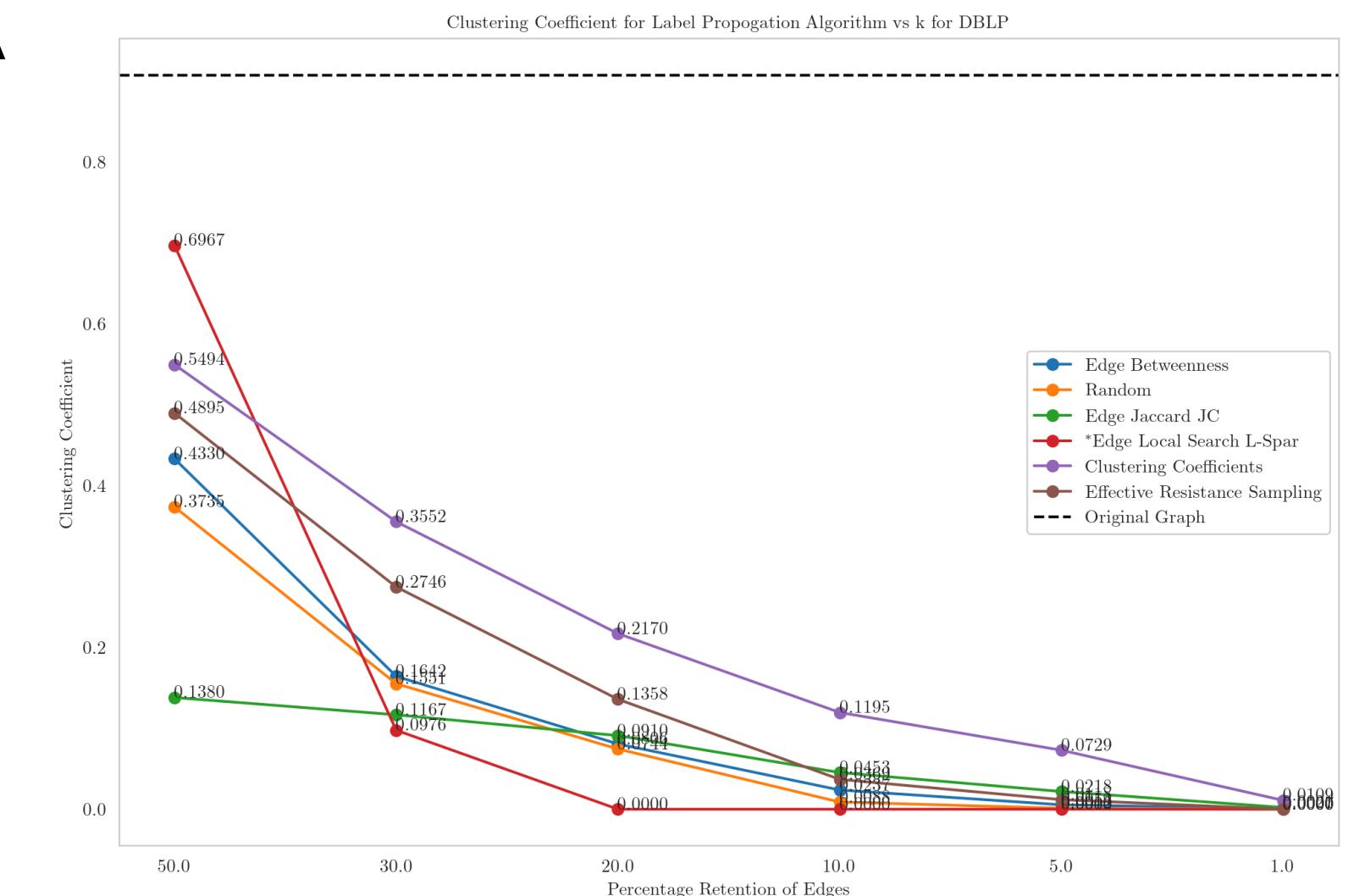
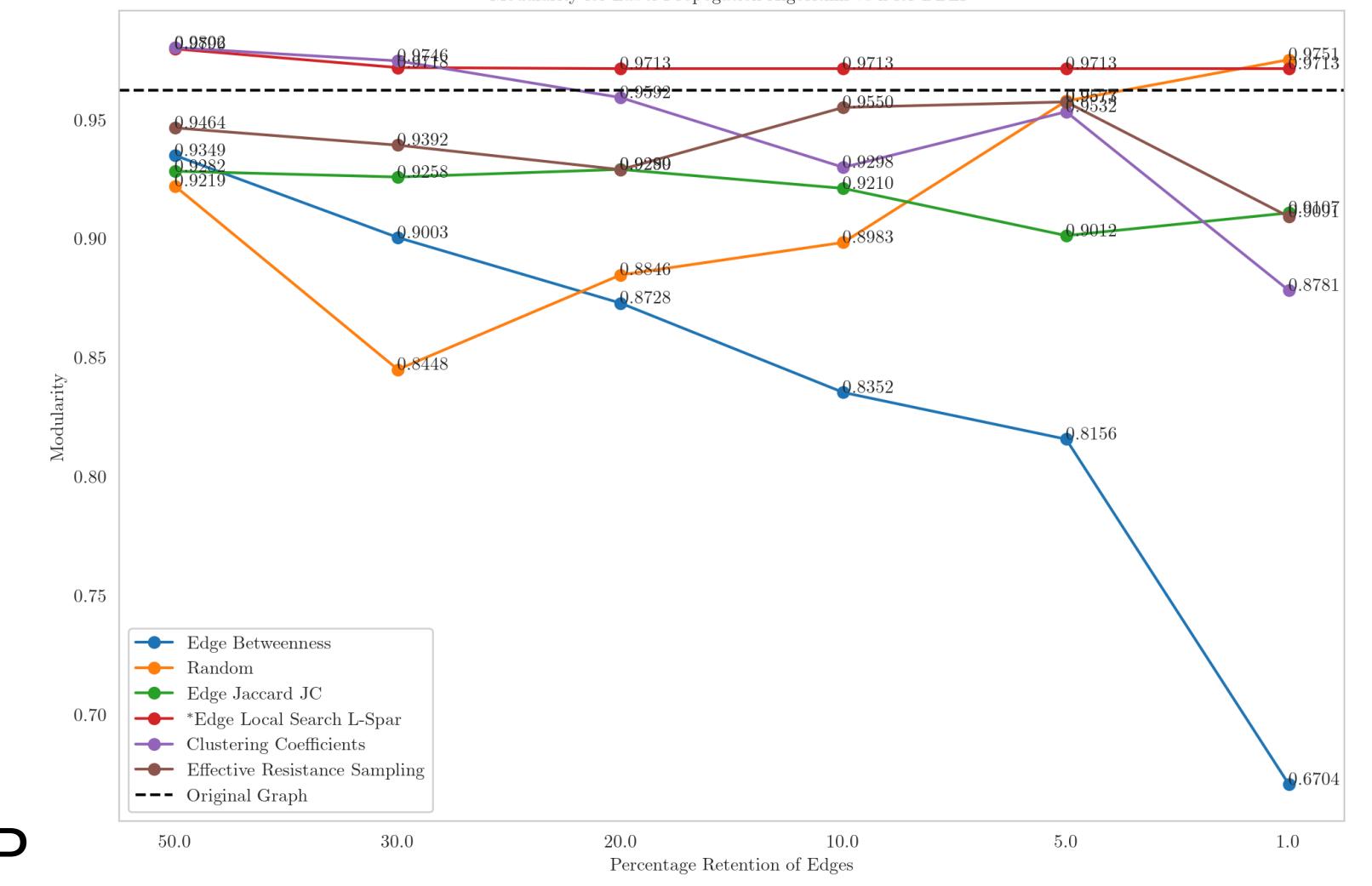
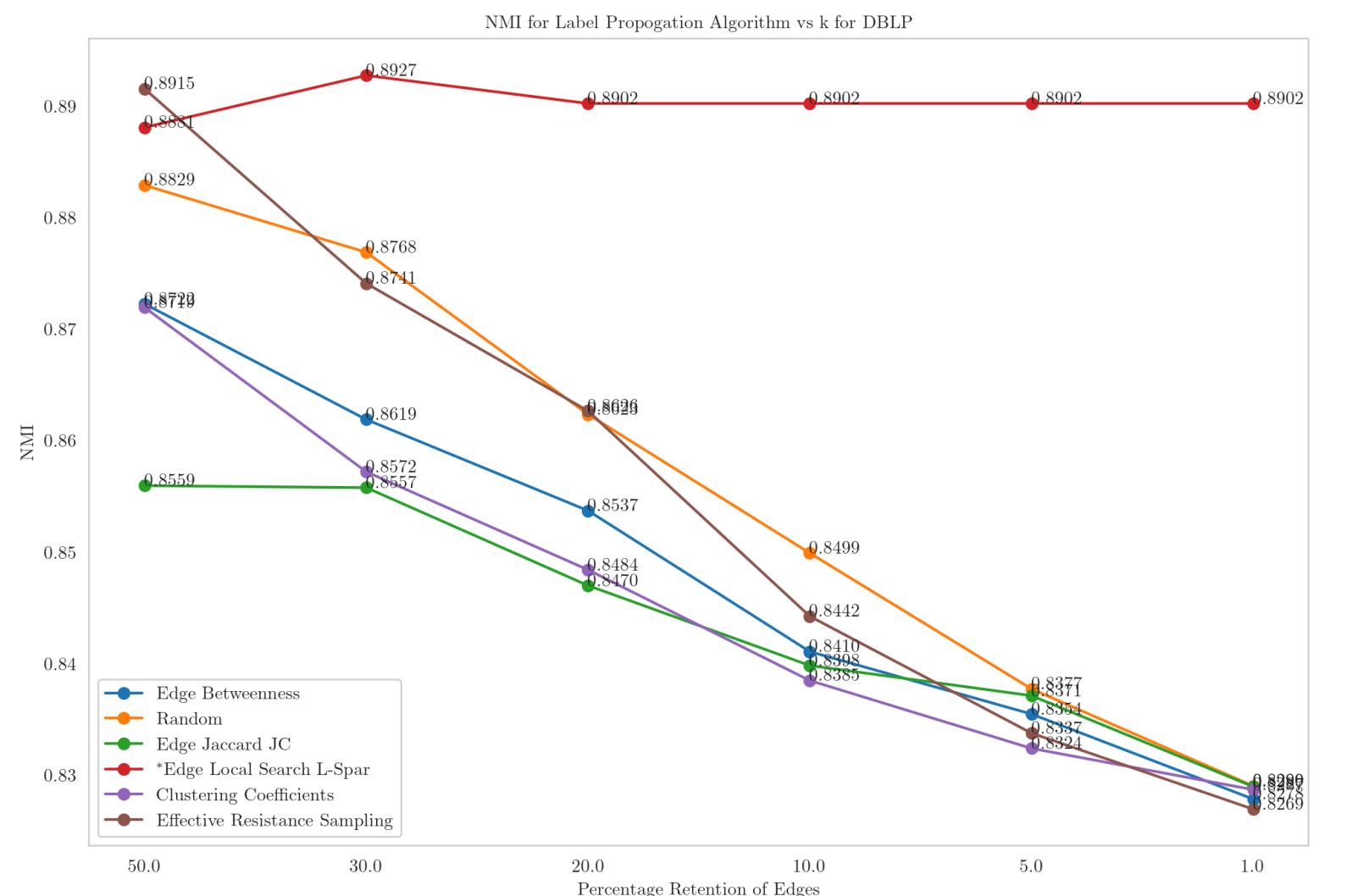


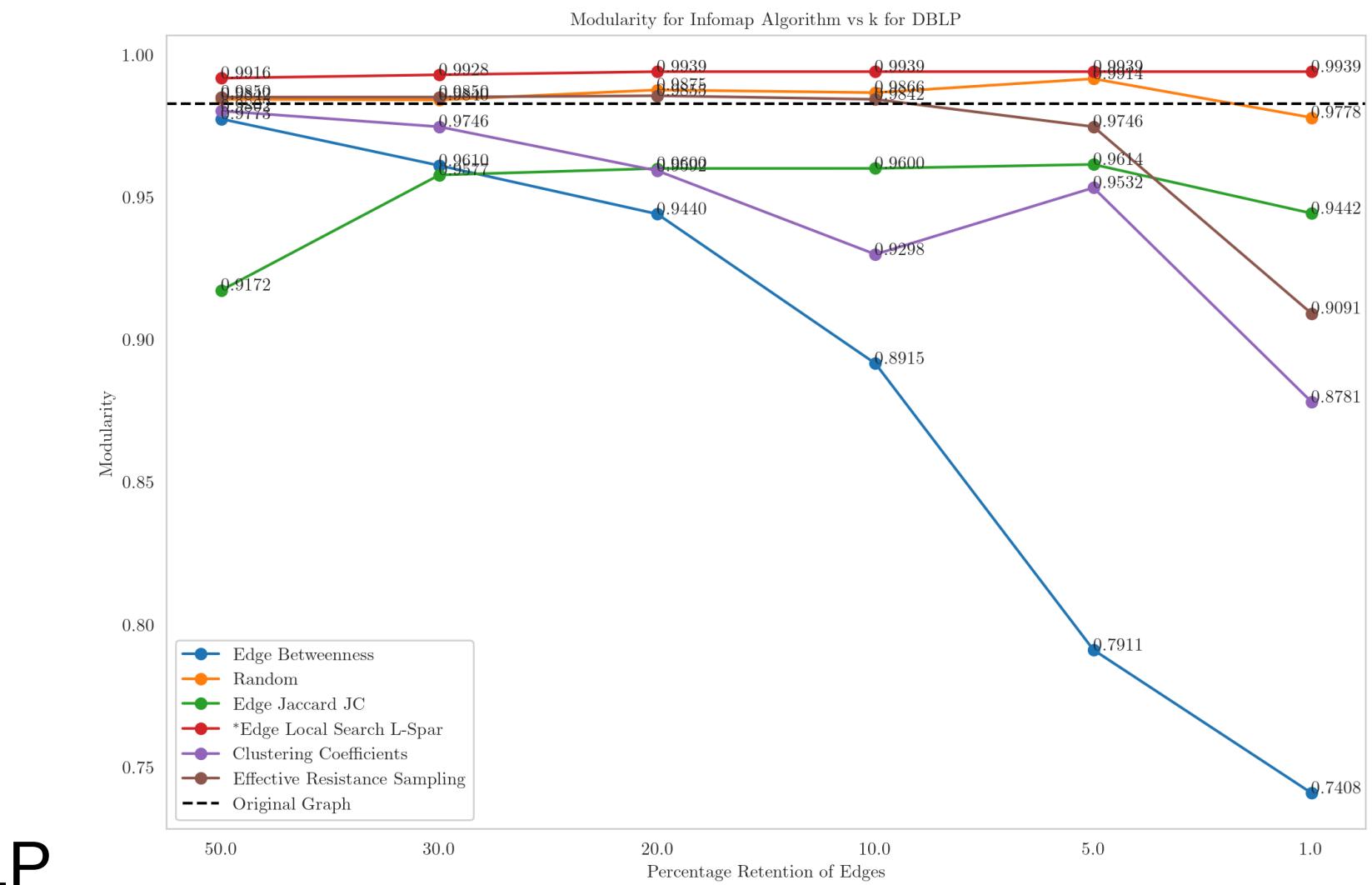
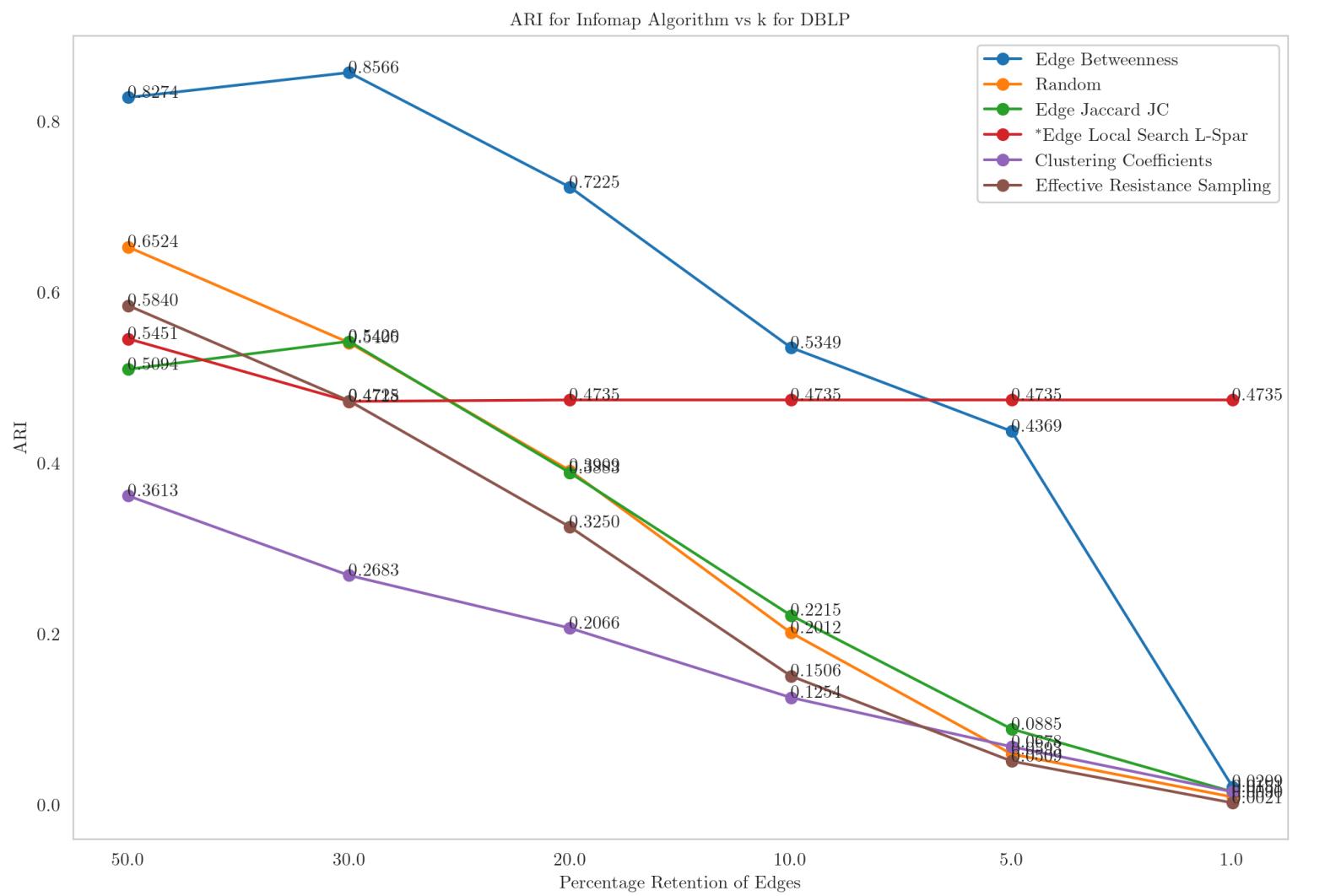
DBLP
Louvain



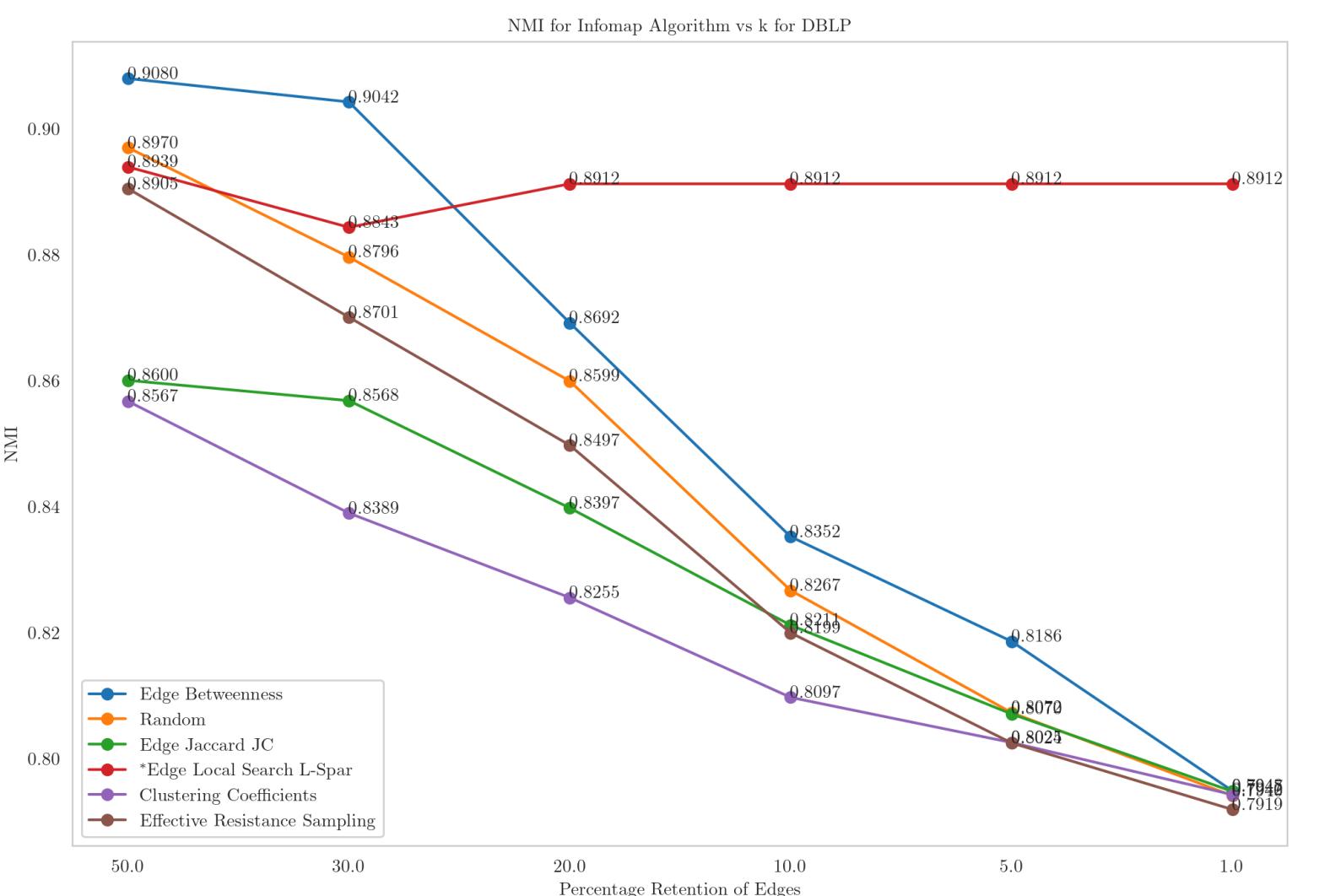
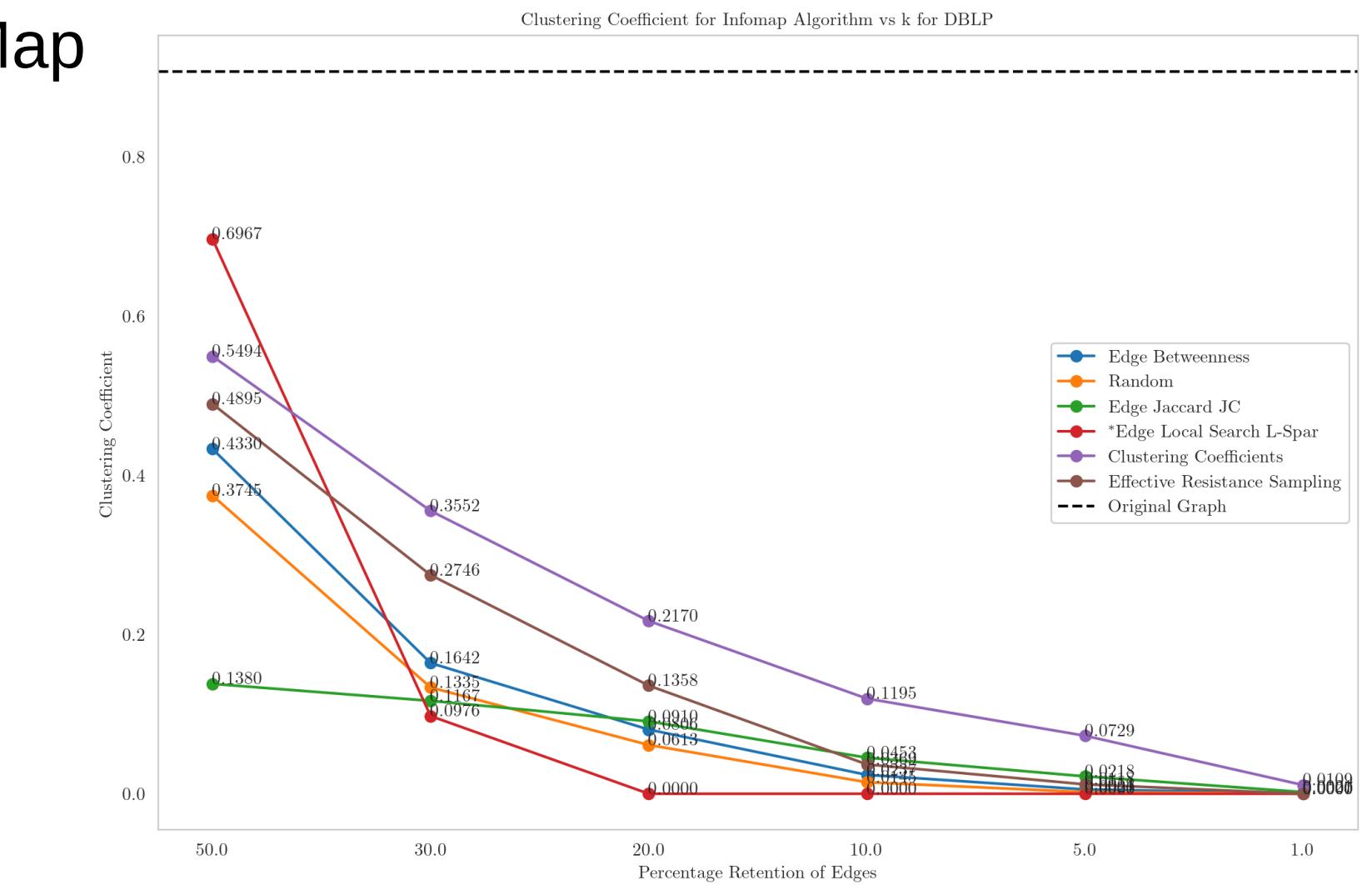


DBL
LPA



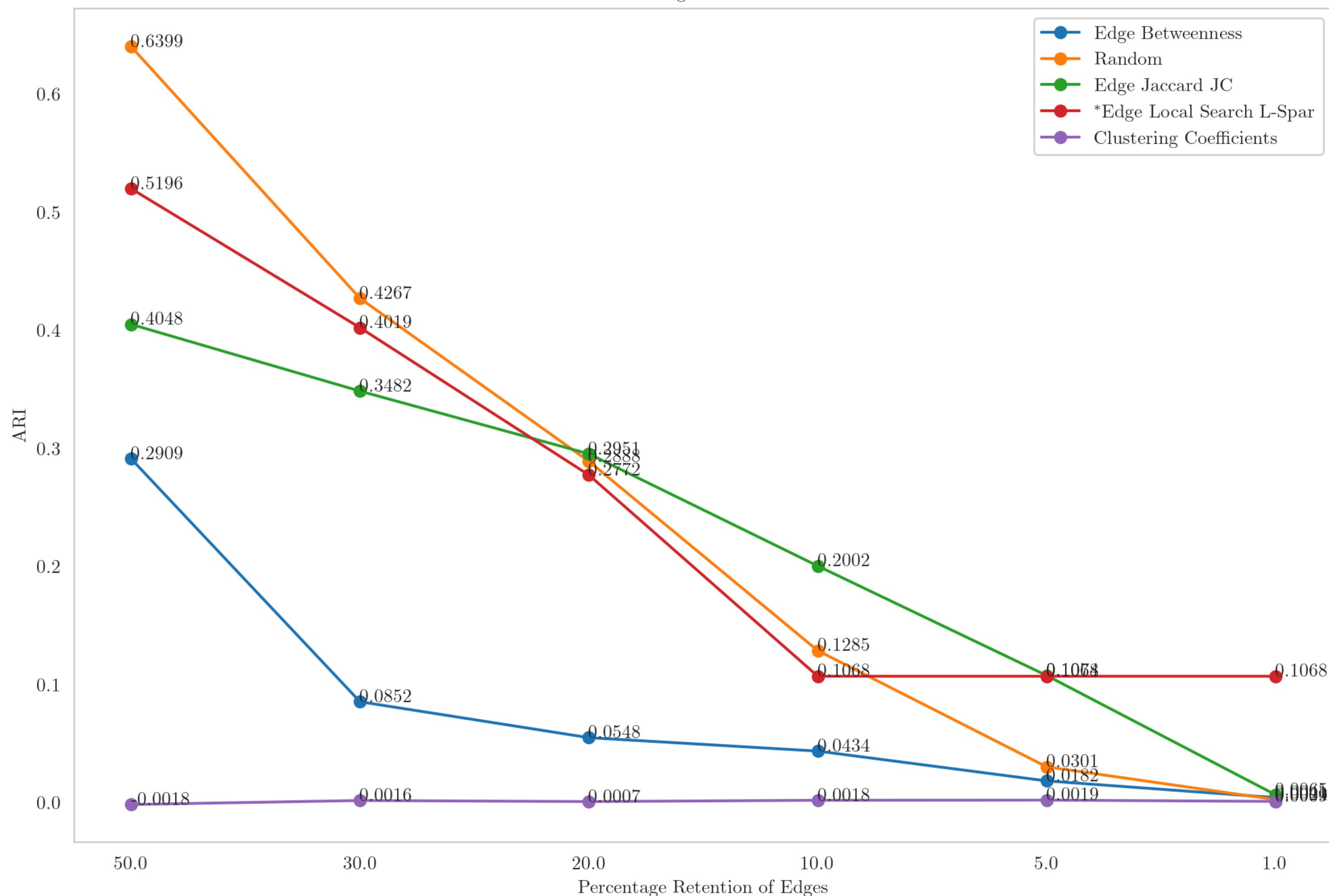


DBLP InfoMap

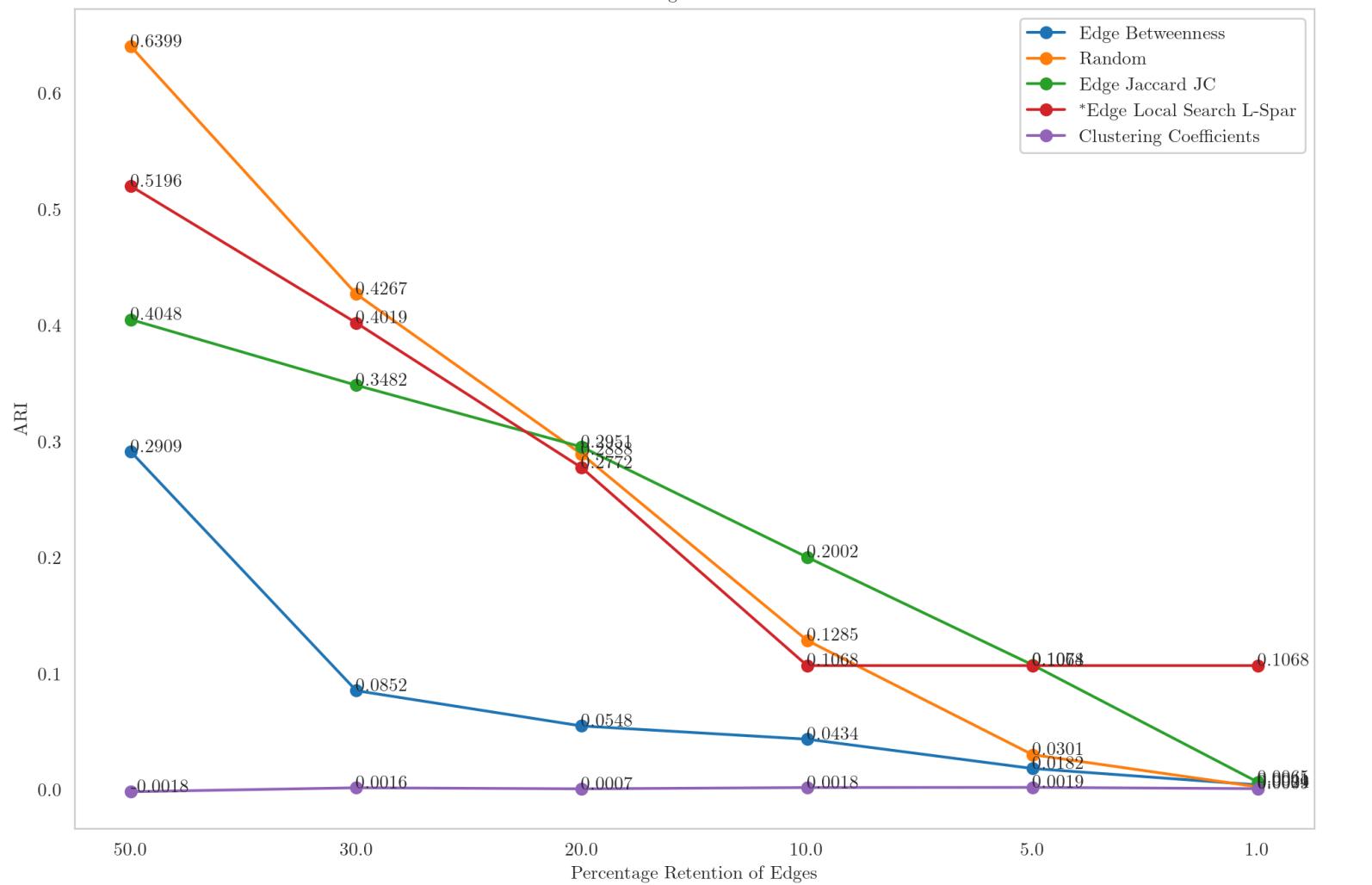


Results for Email EU Core Network

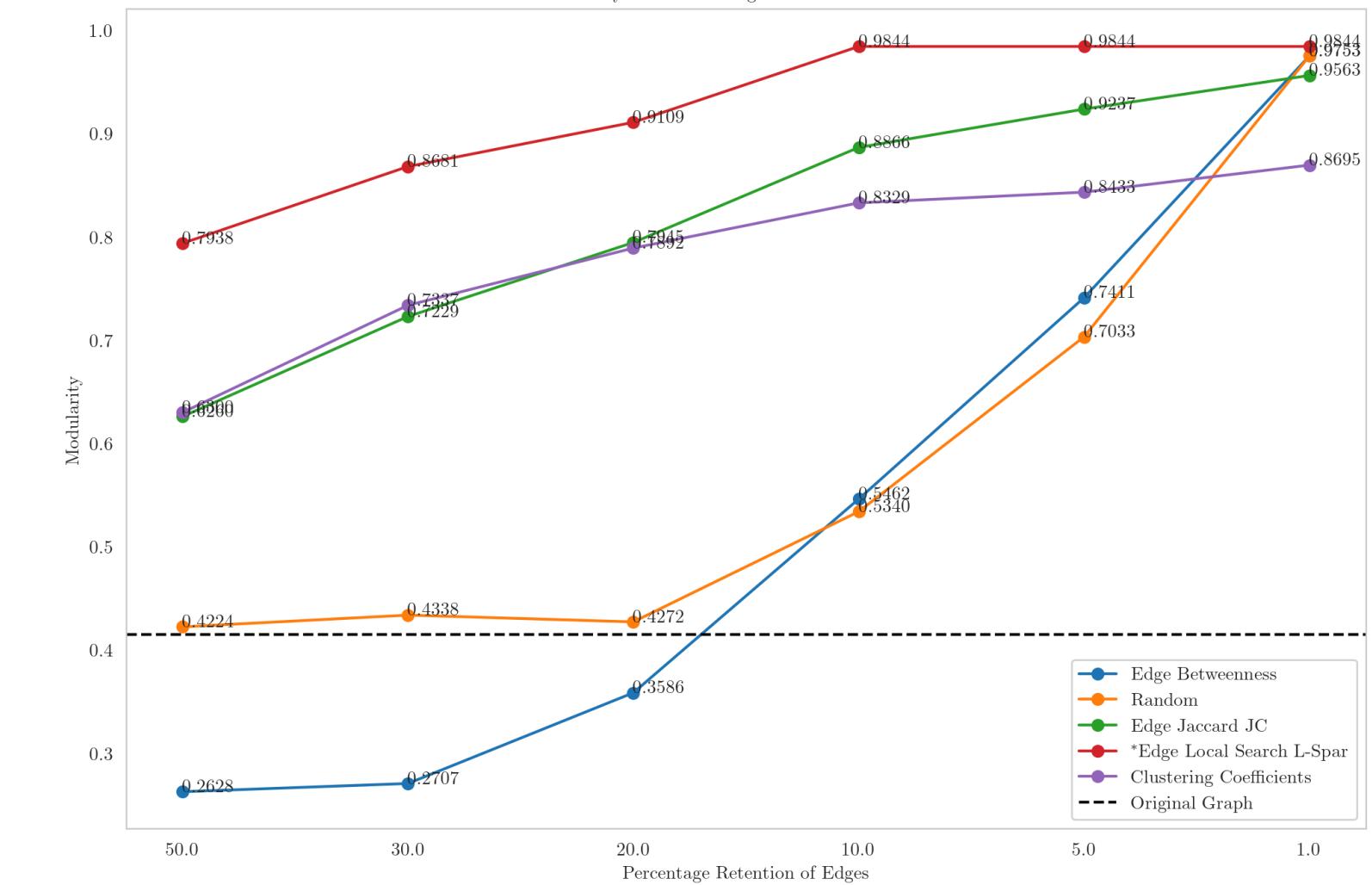
ARI for Louvain Algorithm vs k for EmailEU



ARI for Louvain Algorithm vs k for EmailEU

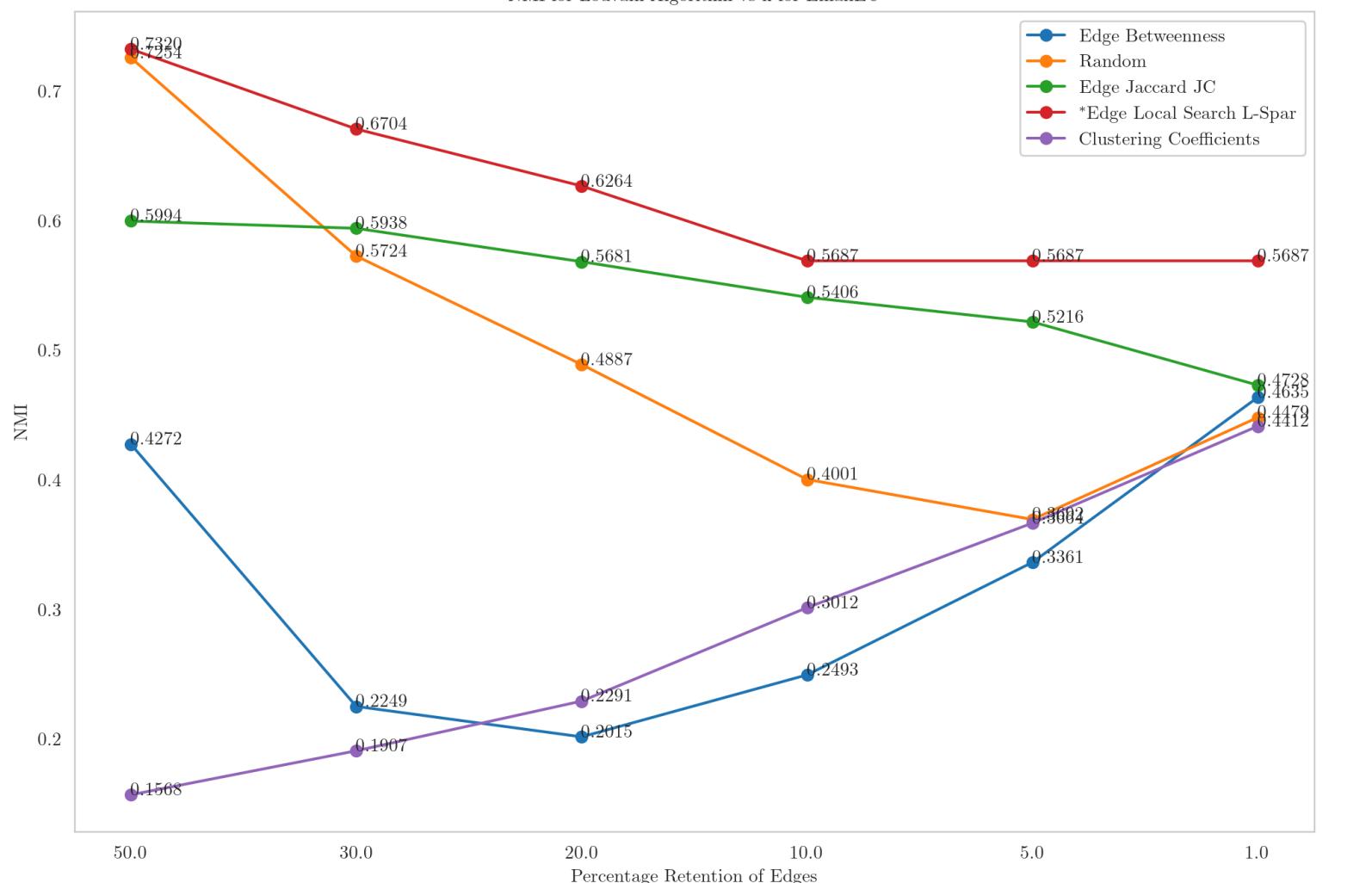


Modularity for Louvain Algorithm vs k for EmailEU

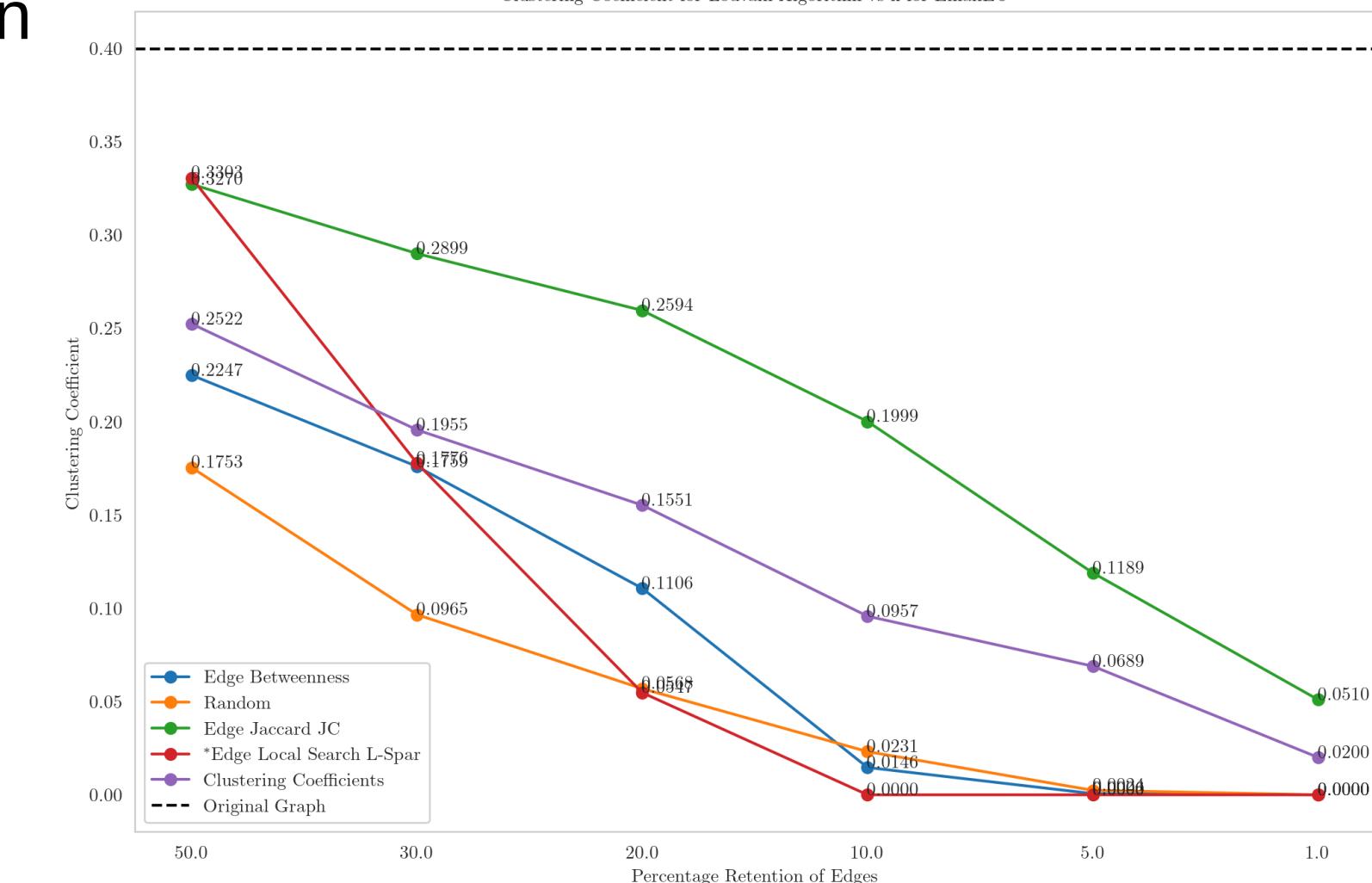


Email Louvain

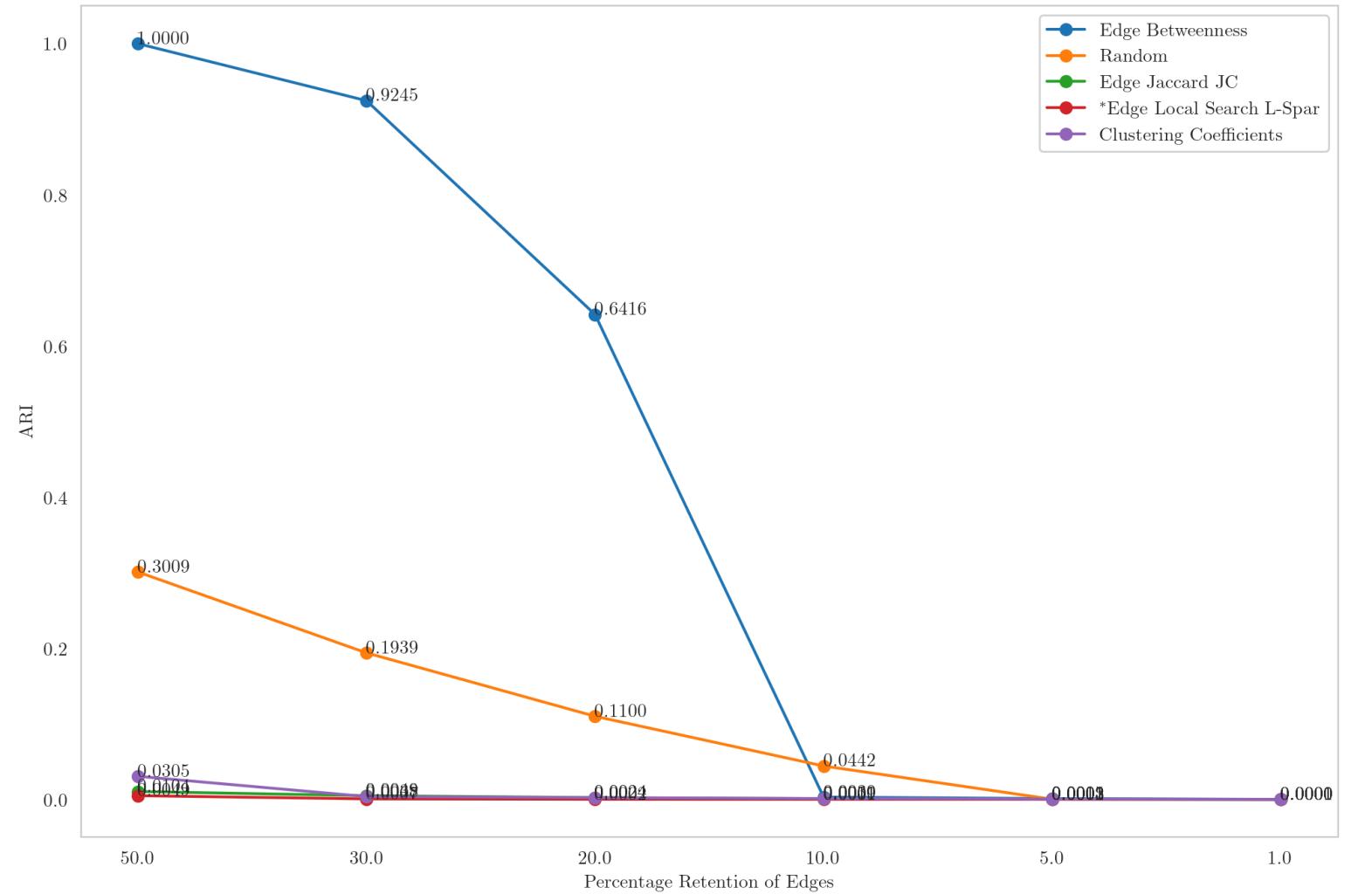
NMI for Louvain Algorithm vs k for EmailEU



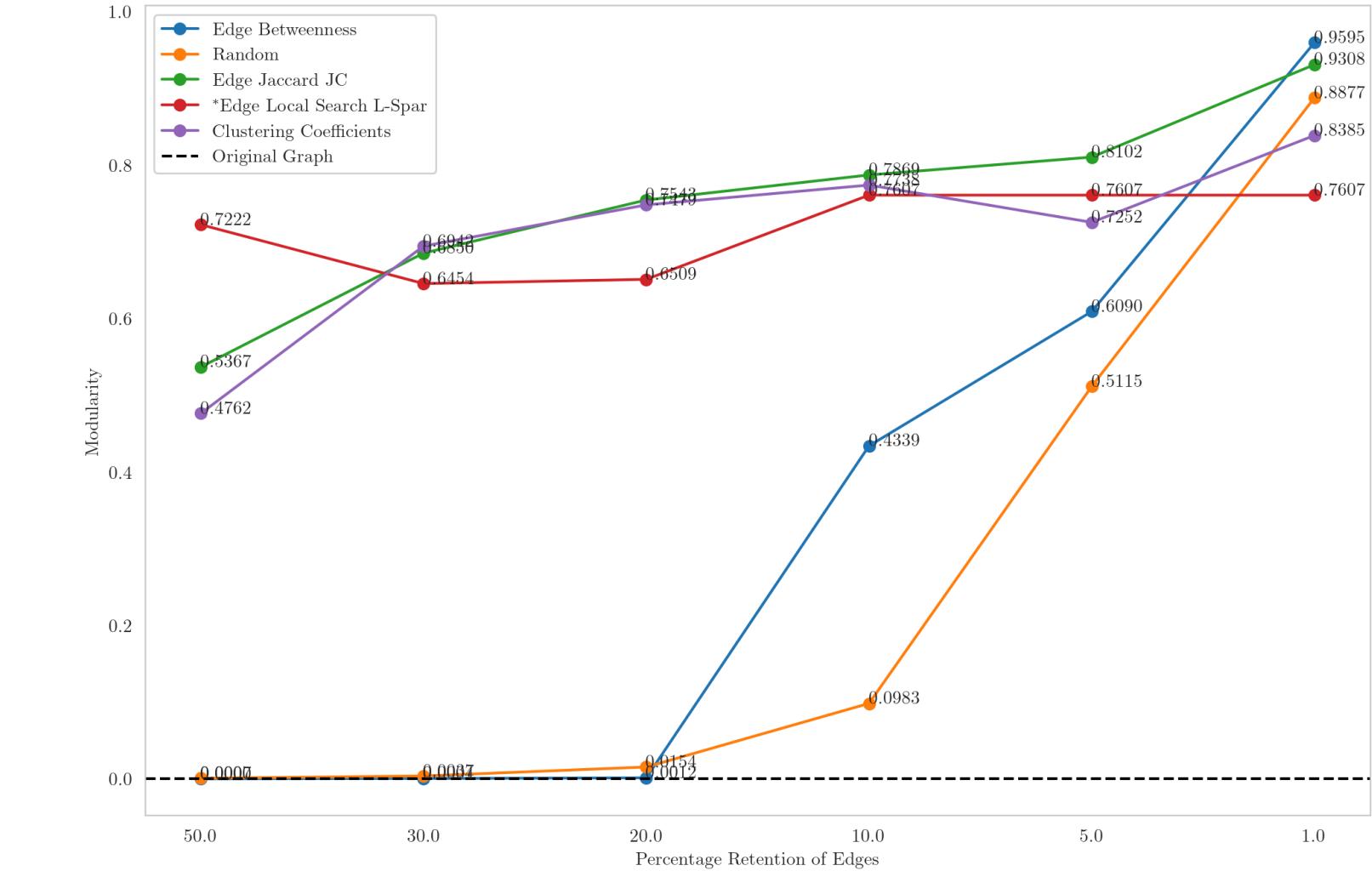
Clustering Coefficient for Louvain Algorithm vs k for EmailEU



ARI for Label Propogation Algorithm vs k for EmailEU

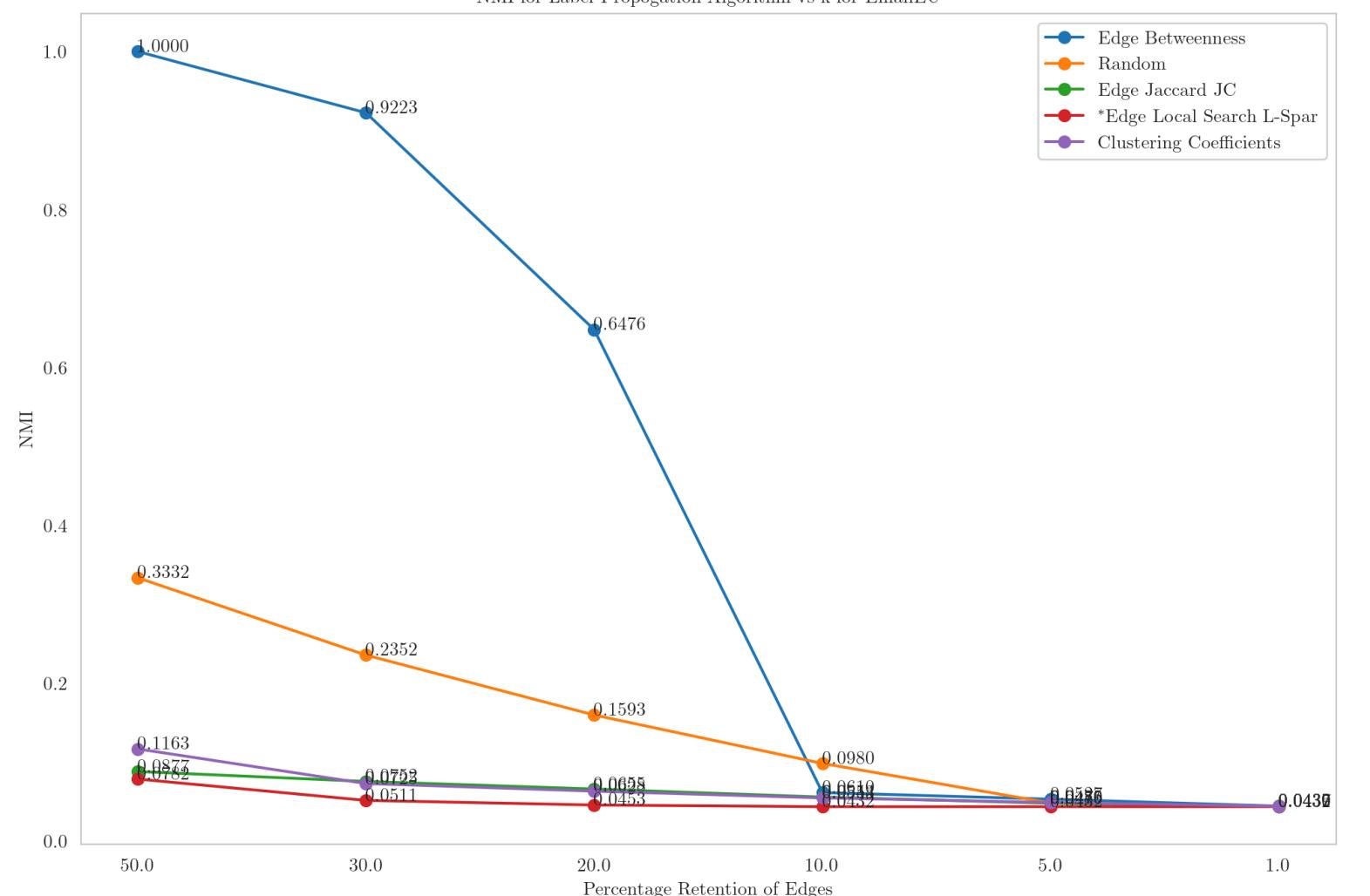


Modularity for Label Propogation Algorithm vs k for EmailEU

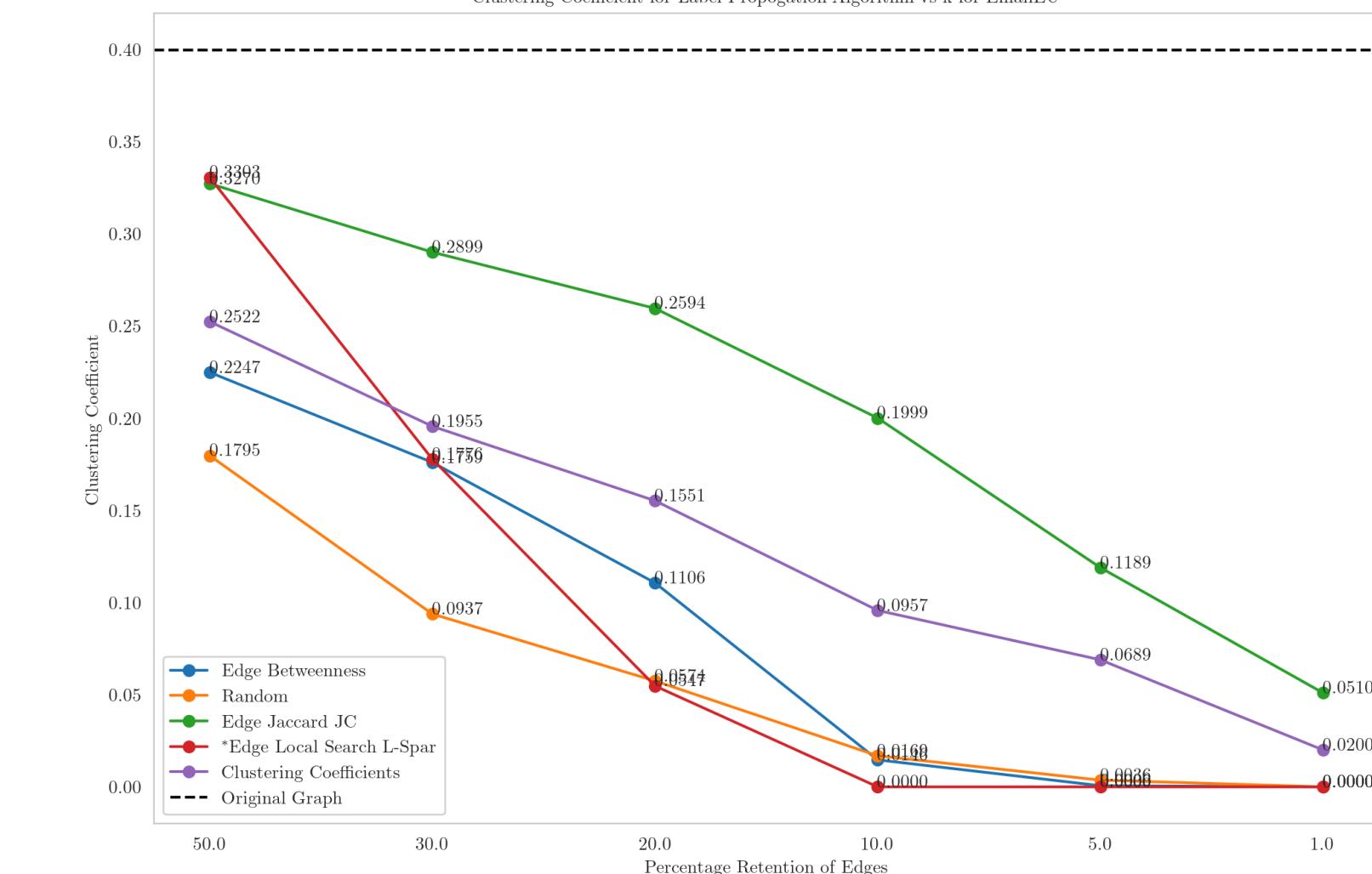


Email LPA

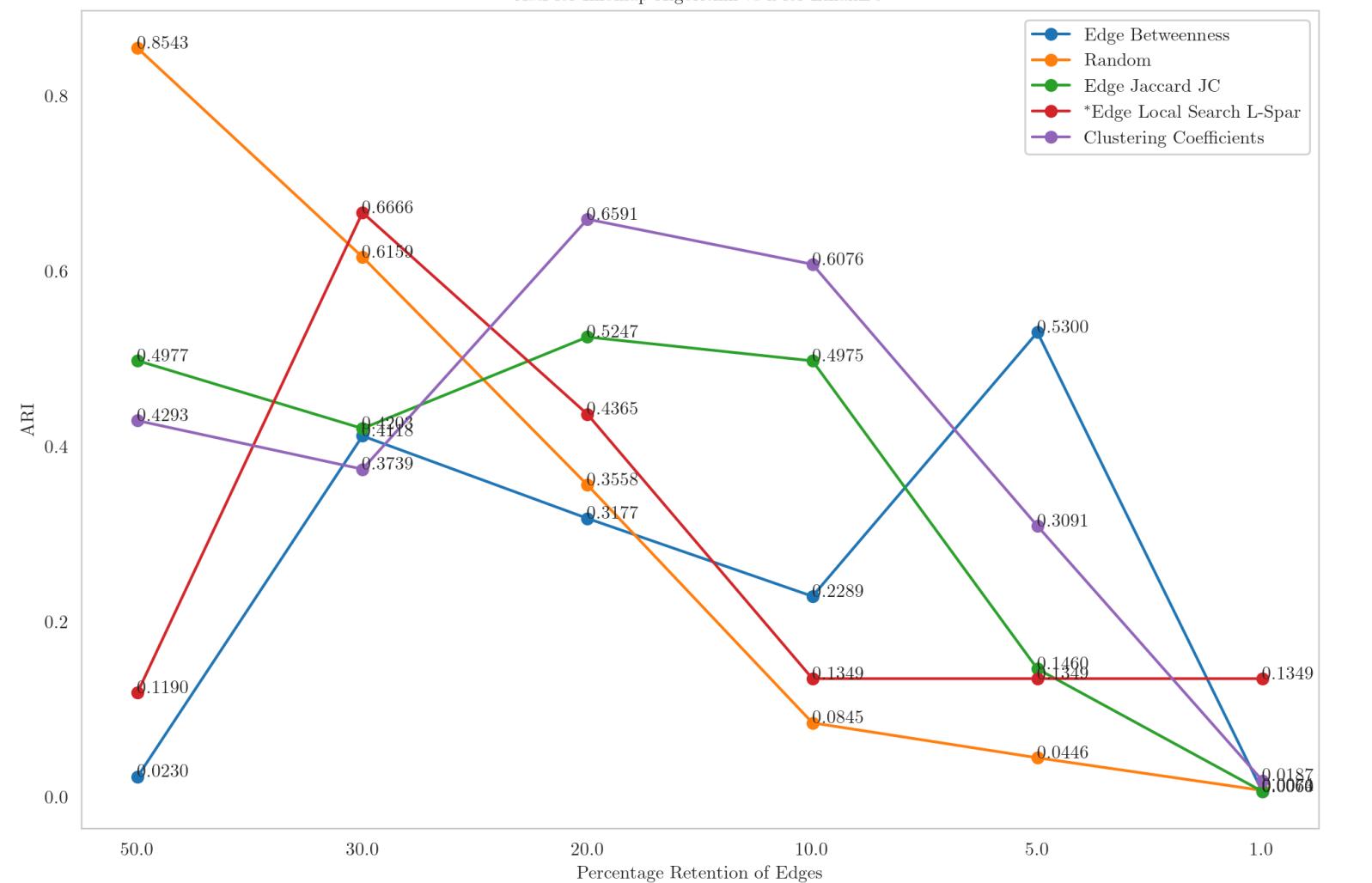
NMI for Label Propogation Algorithm vs k for EmailEU



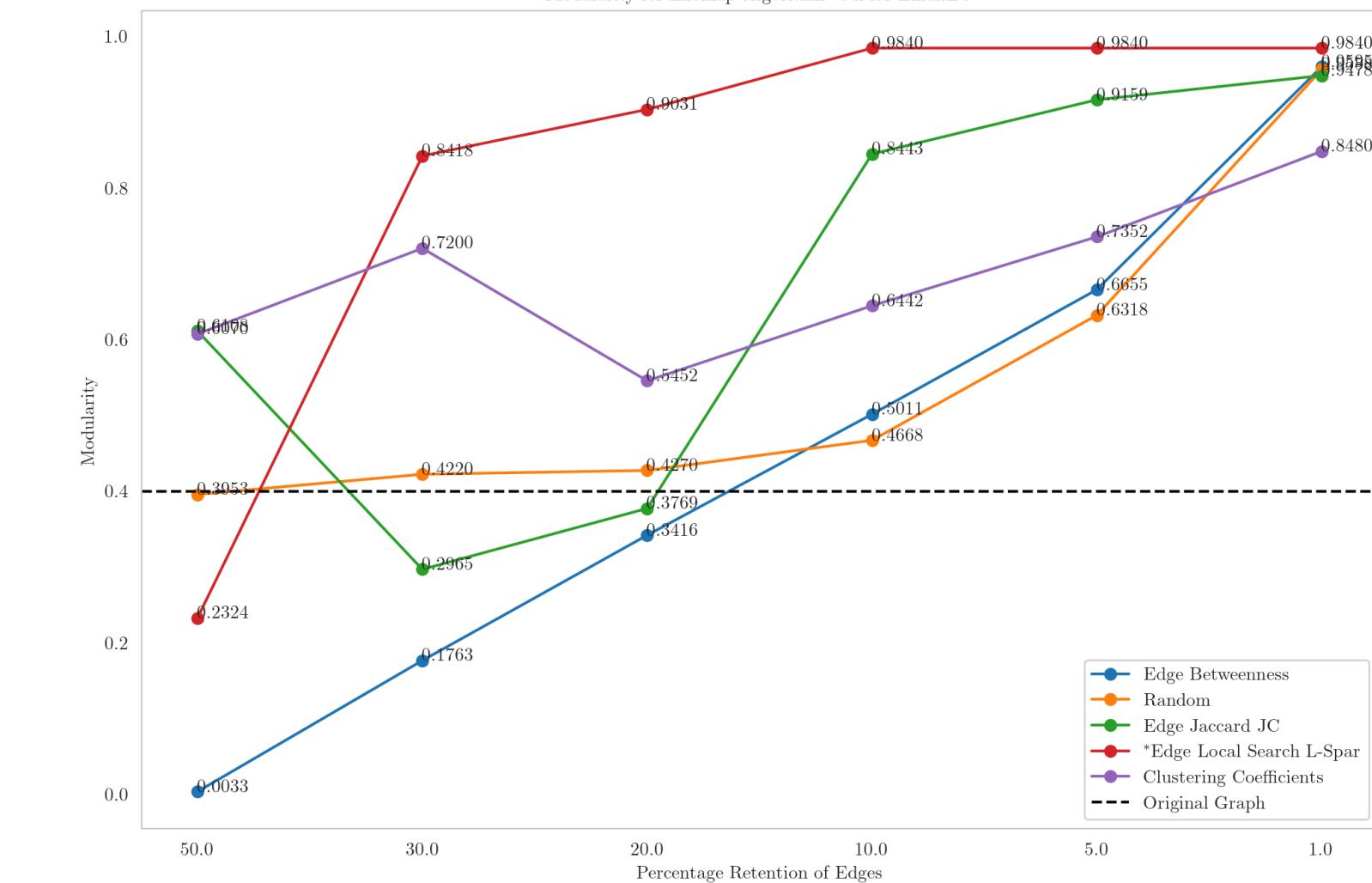
Clustering Coefficient for Label Propogation Algorithm vs k for EmailEU



ARI for Infomap Algorithm vs k for EmailEU

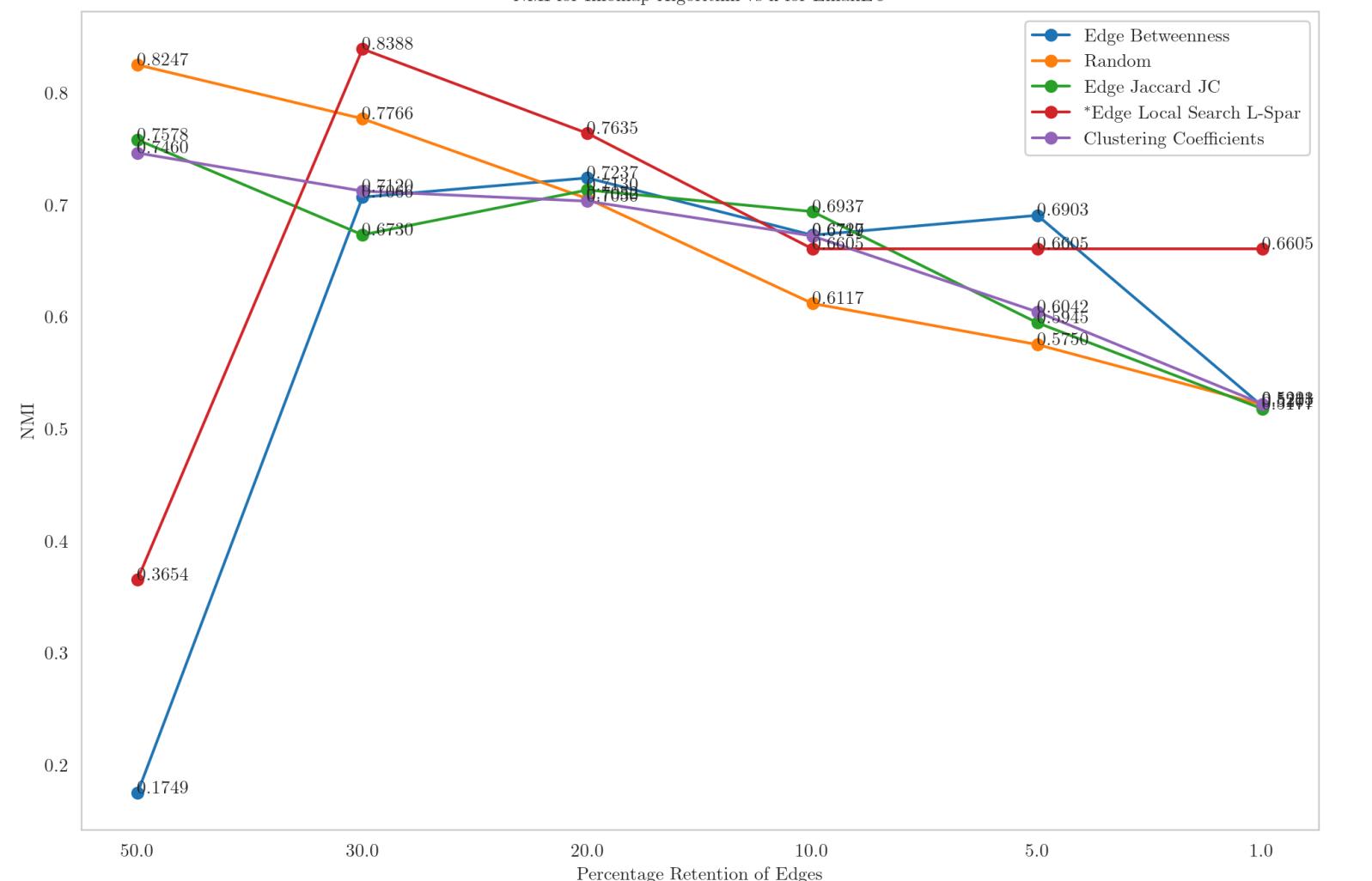


Modularity for Infomap Algorithm vs k for EmailEU

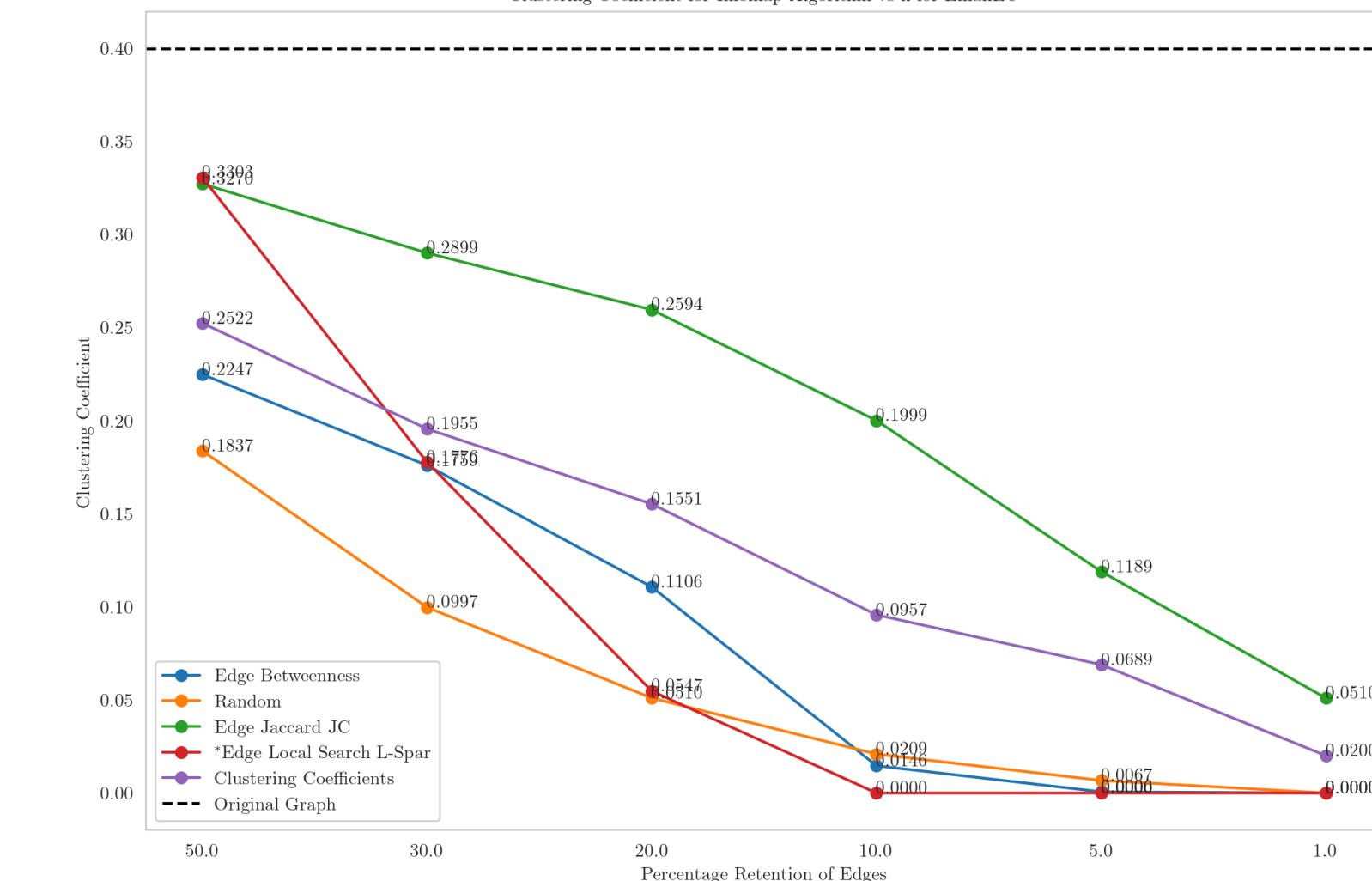


Email InfoMap

NMI for Infomap Algorithm vs k for EmailEU

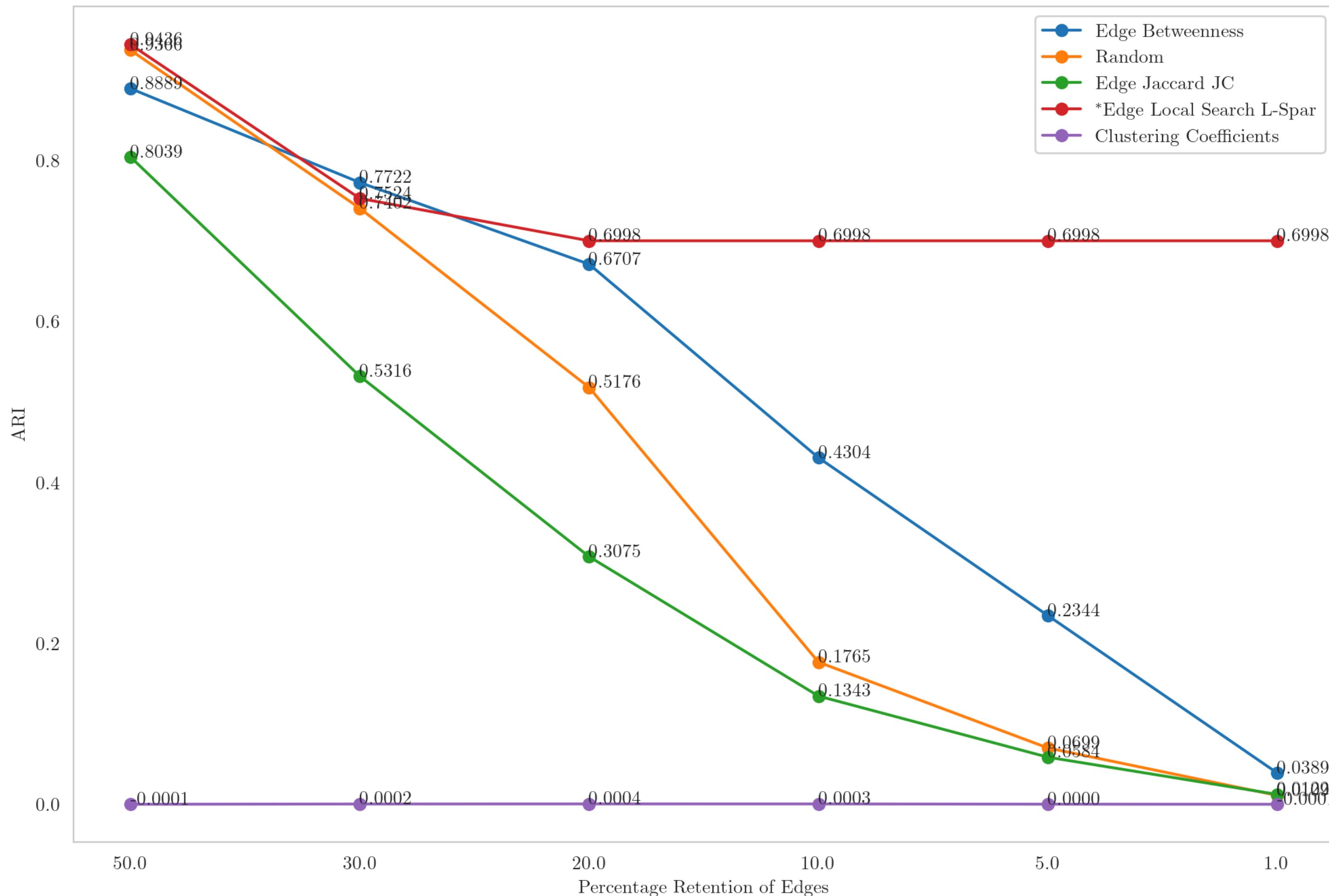


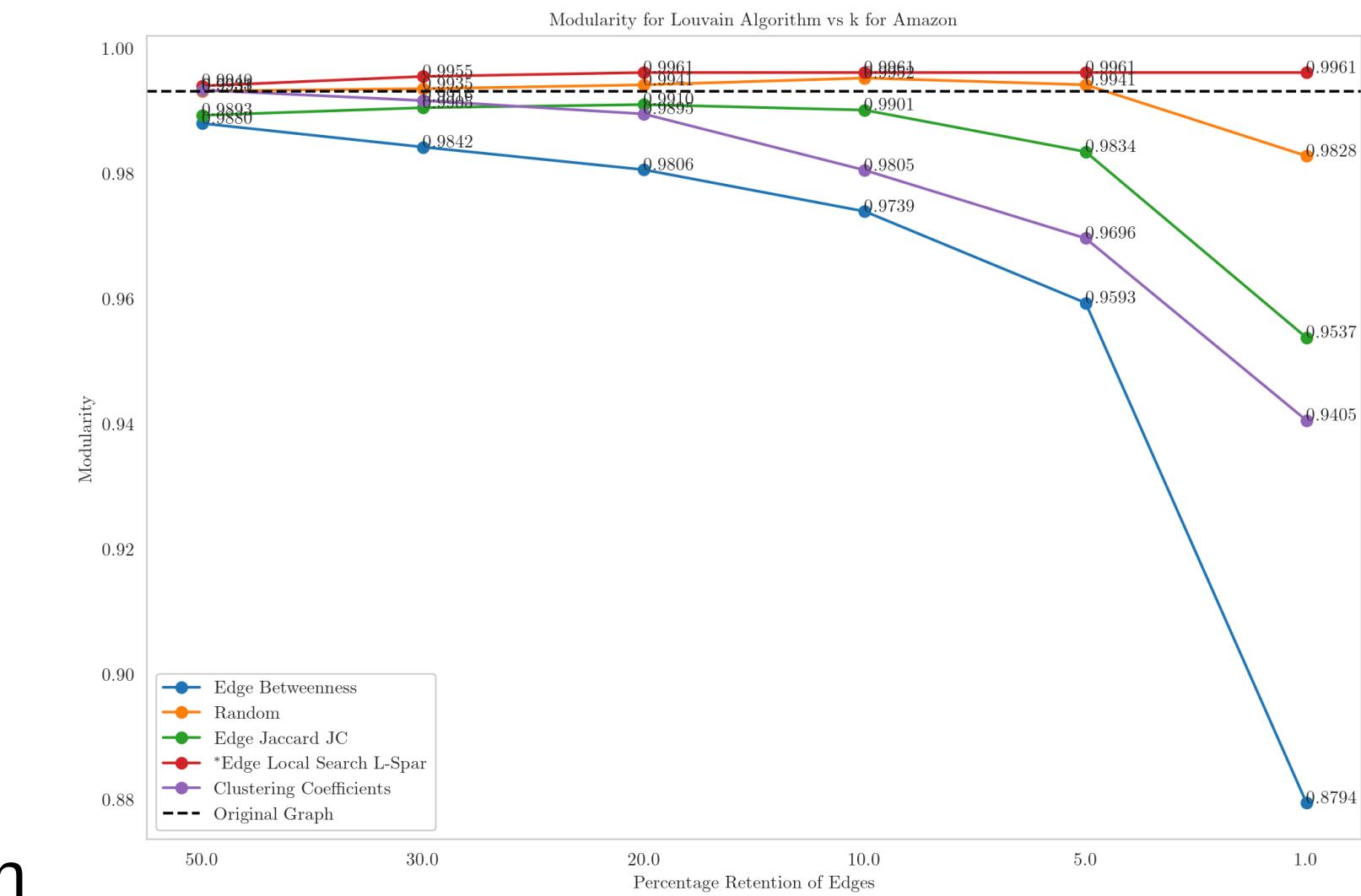
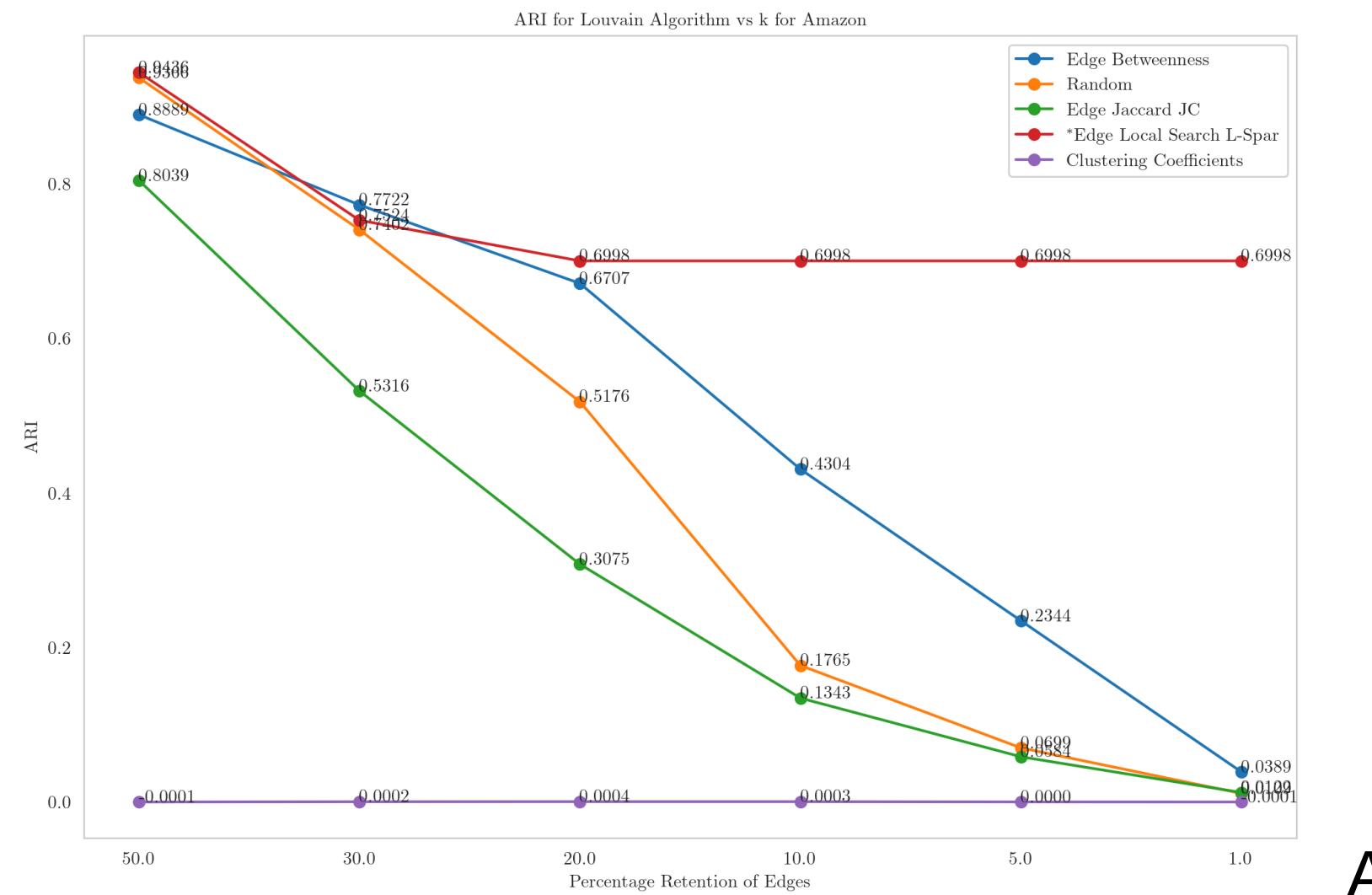
Clustering Coefficient for Infomap Algorithm vs k for EmailEU



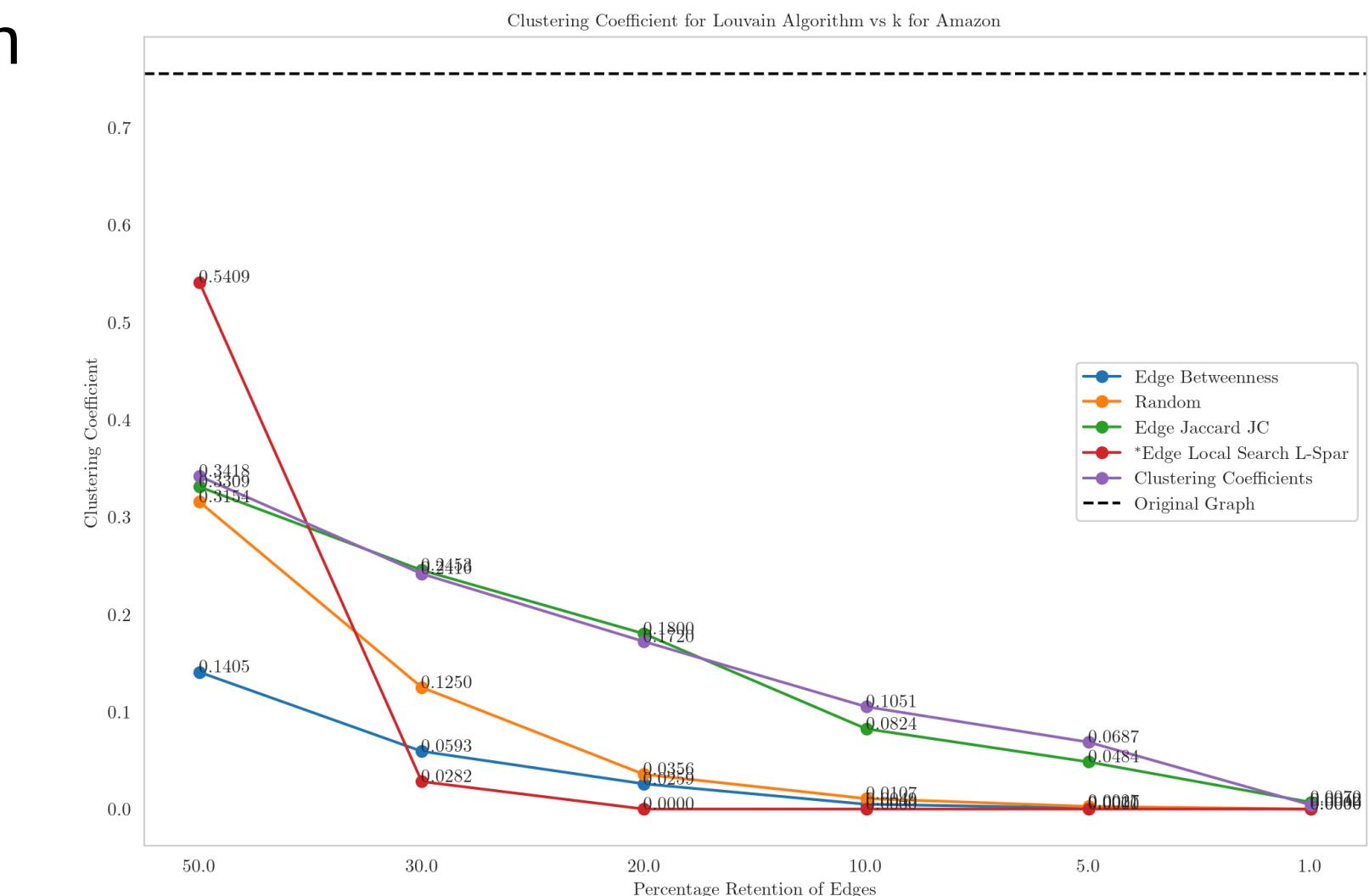
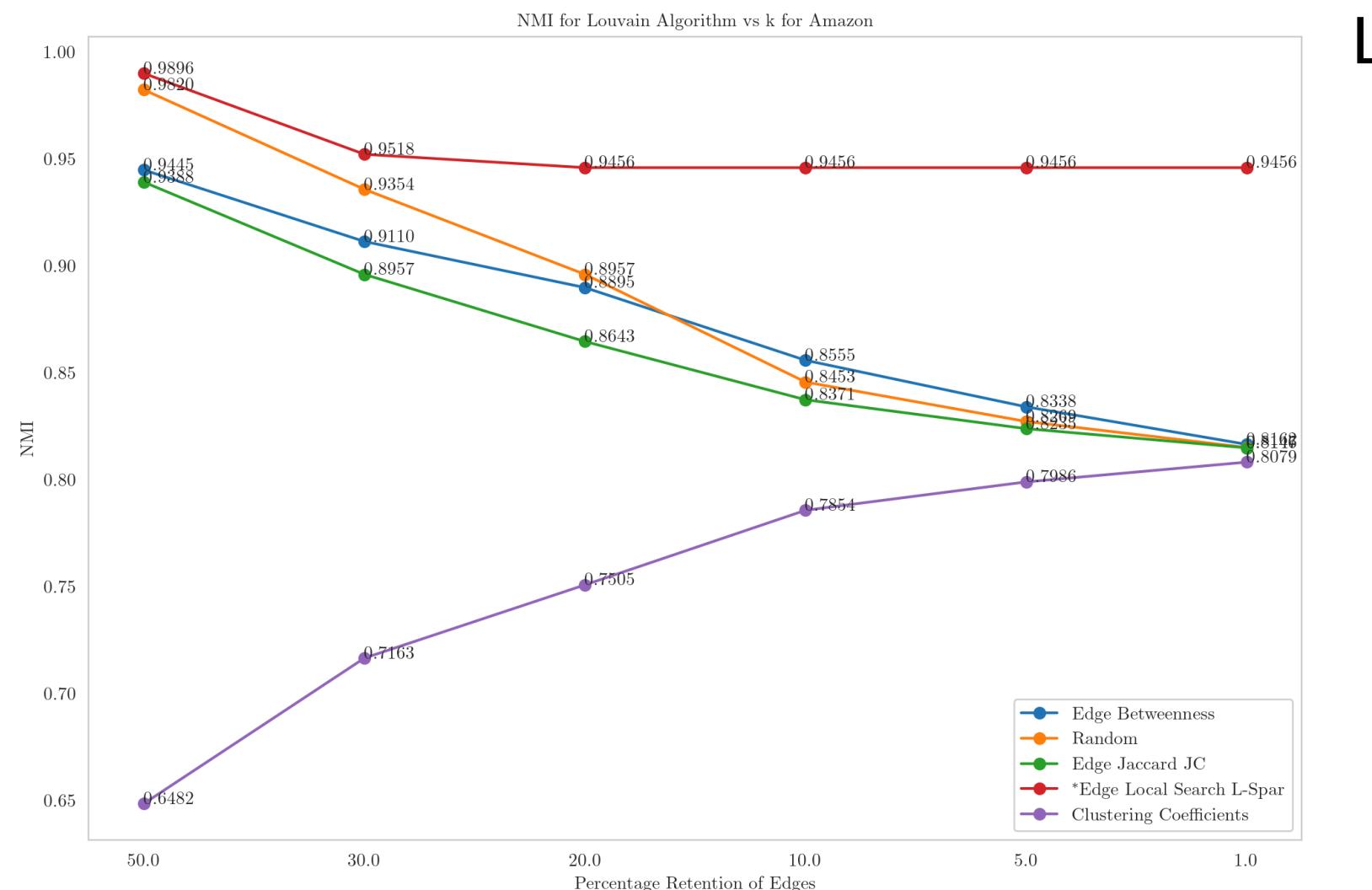
Results for Amazon Co-Purchase Network

ARI for Louvain Algorithm vs k for Amazon

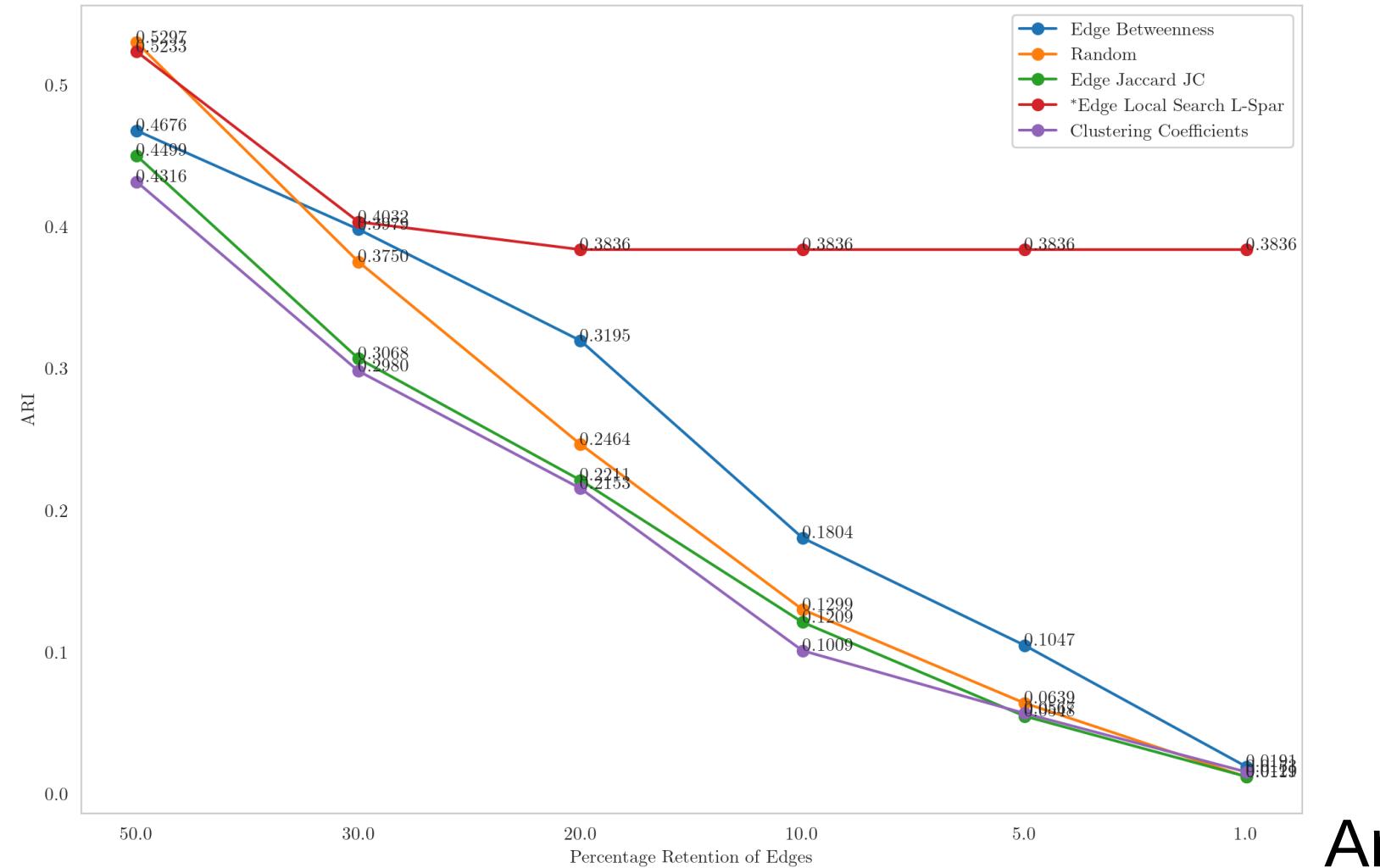




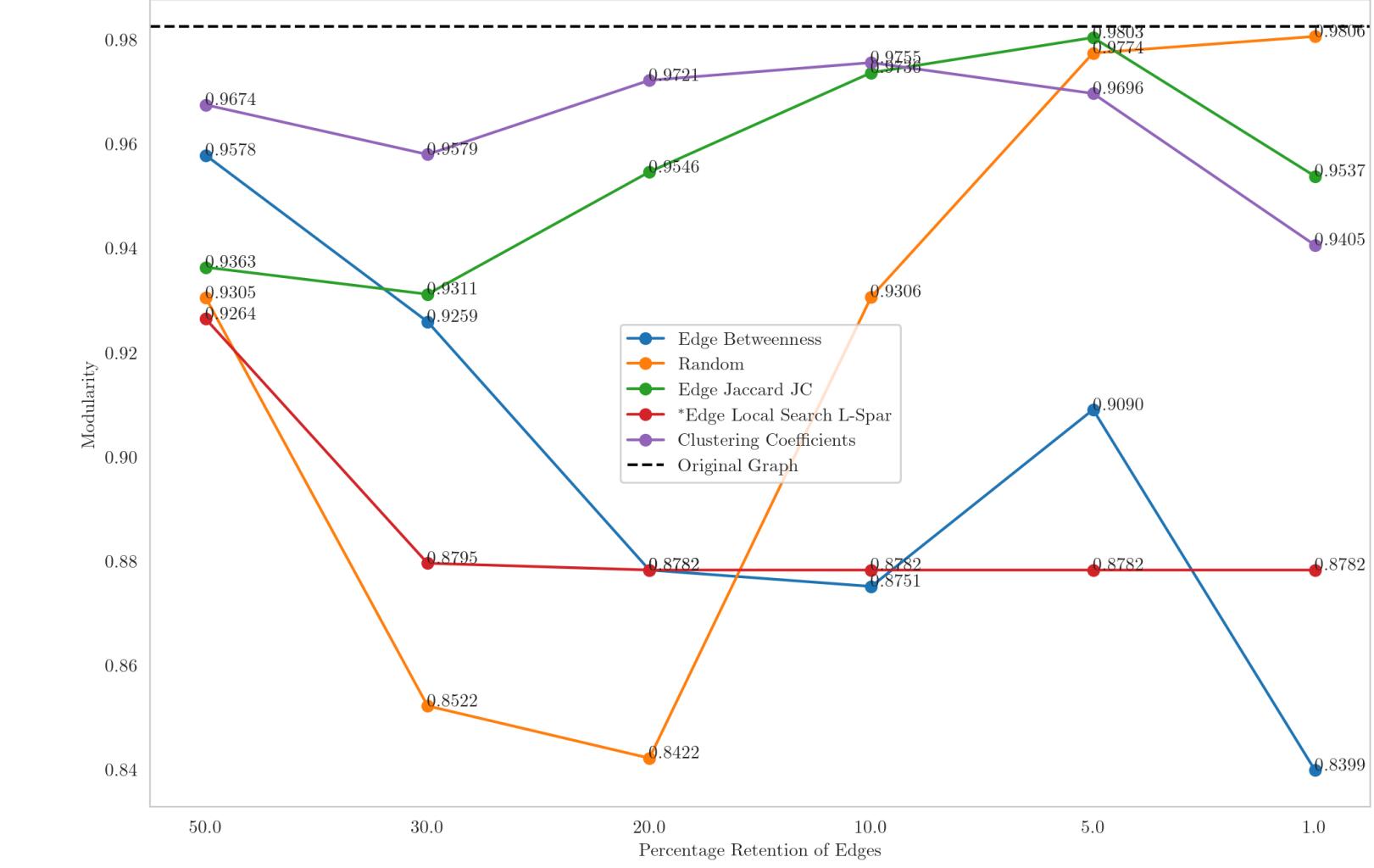
Amazon Louvain



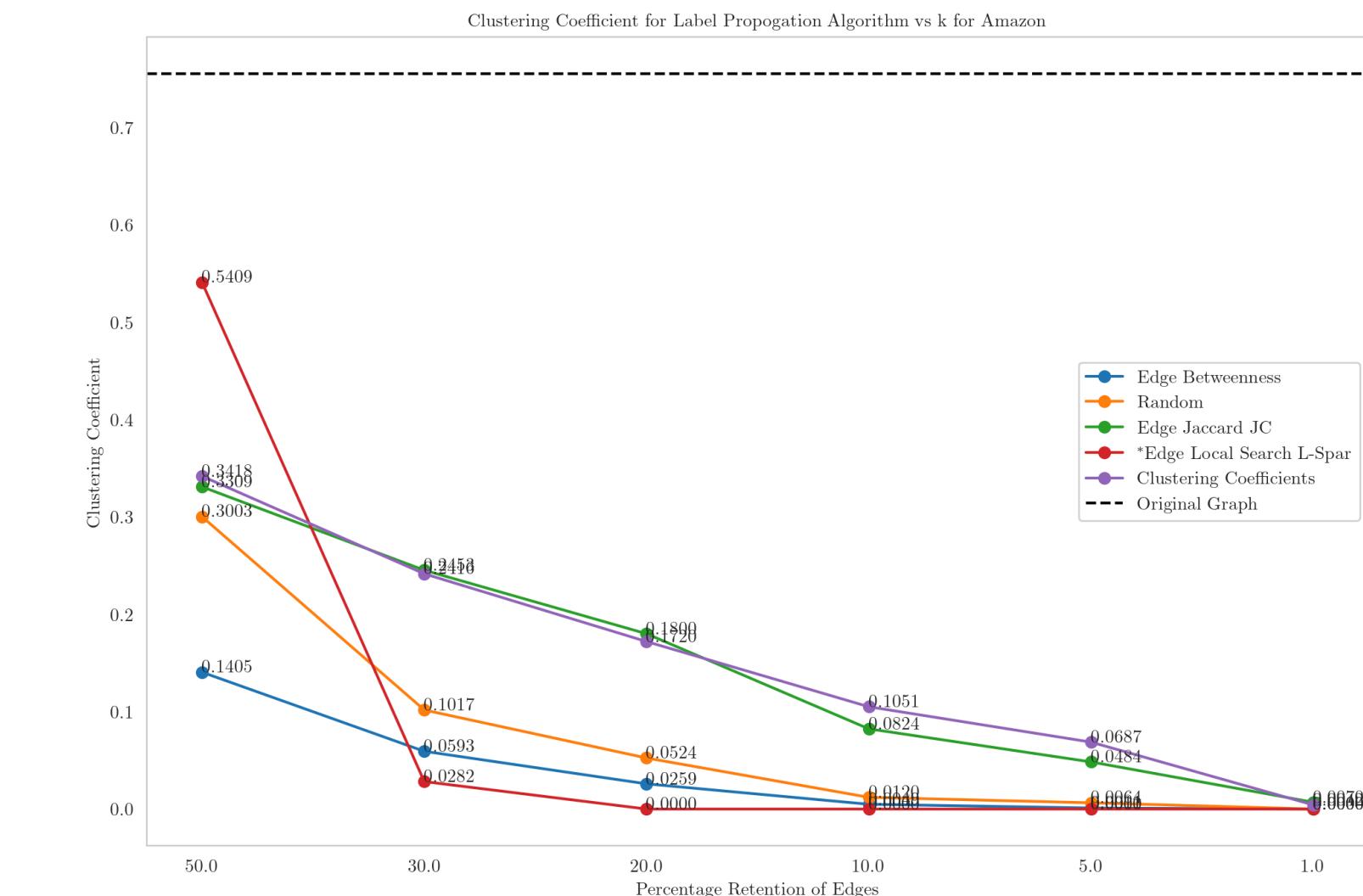
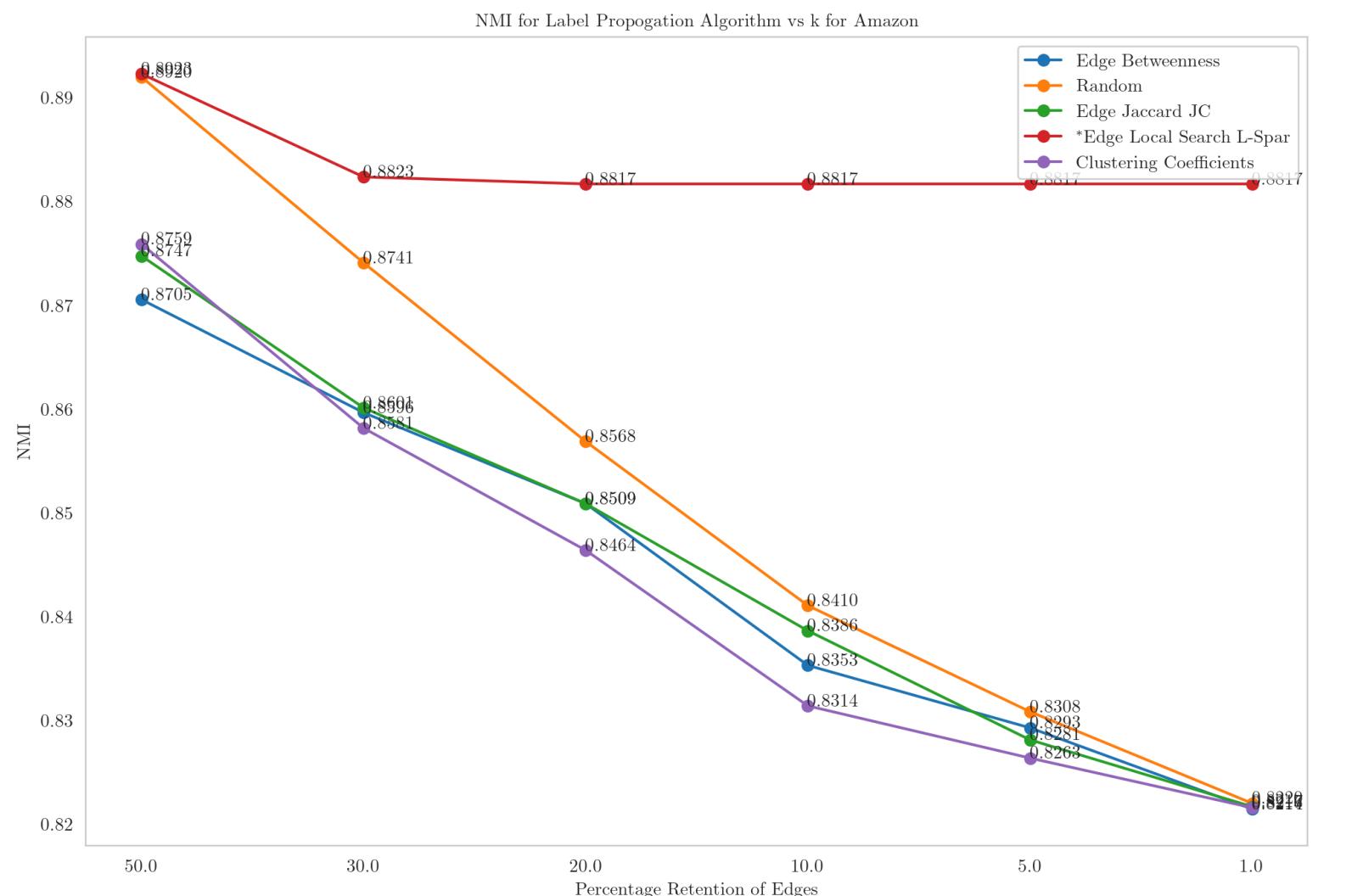
ARI for Label Propogation Algorithm vs k for Amazon

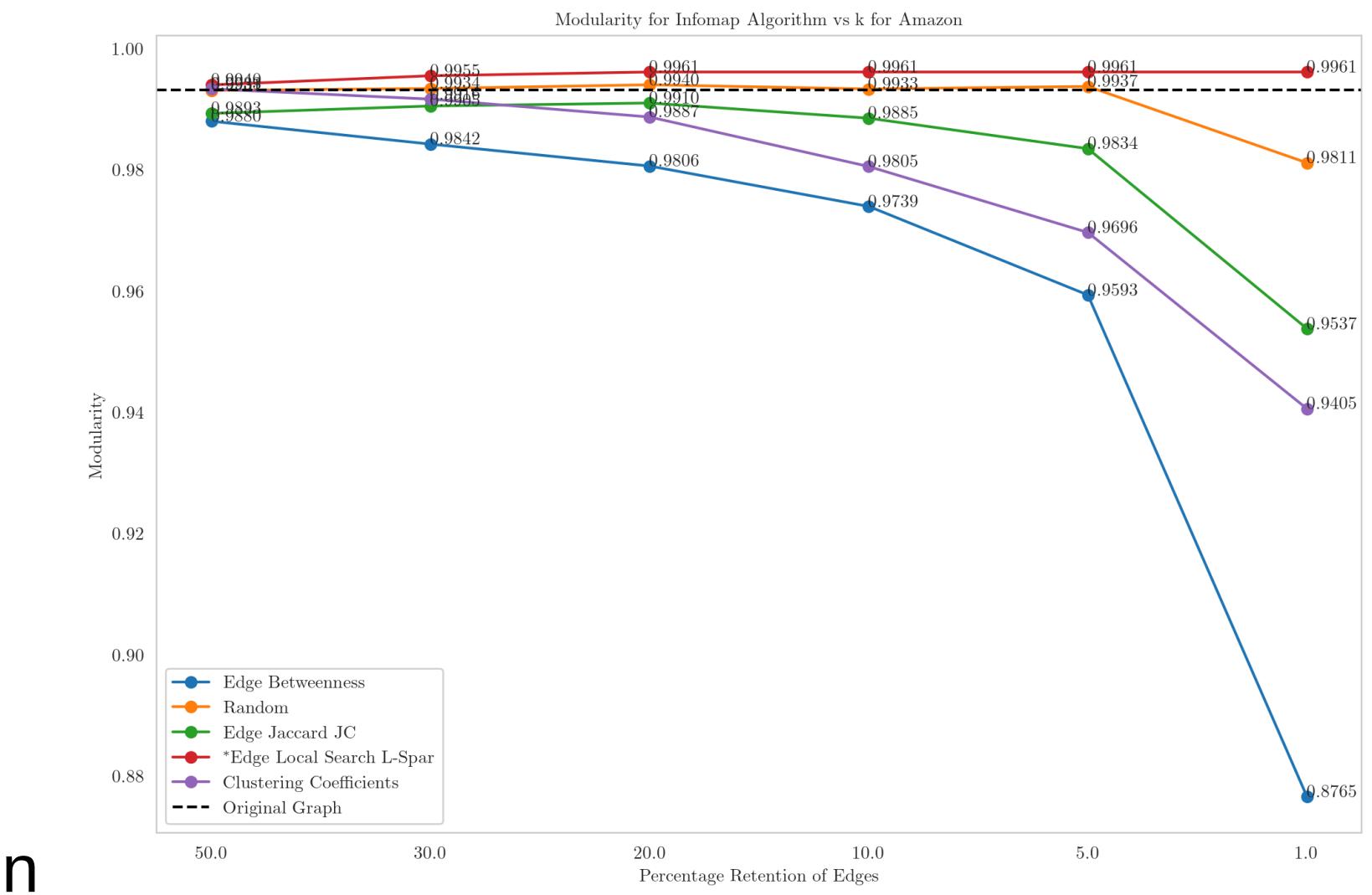
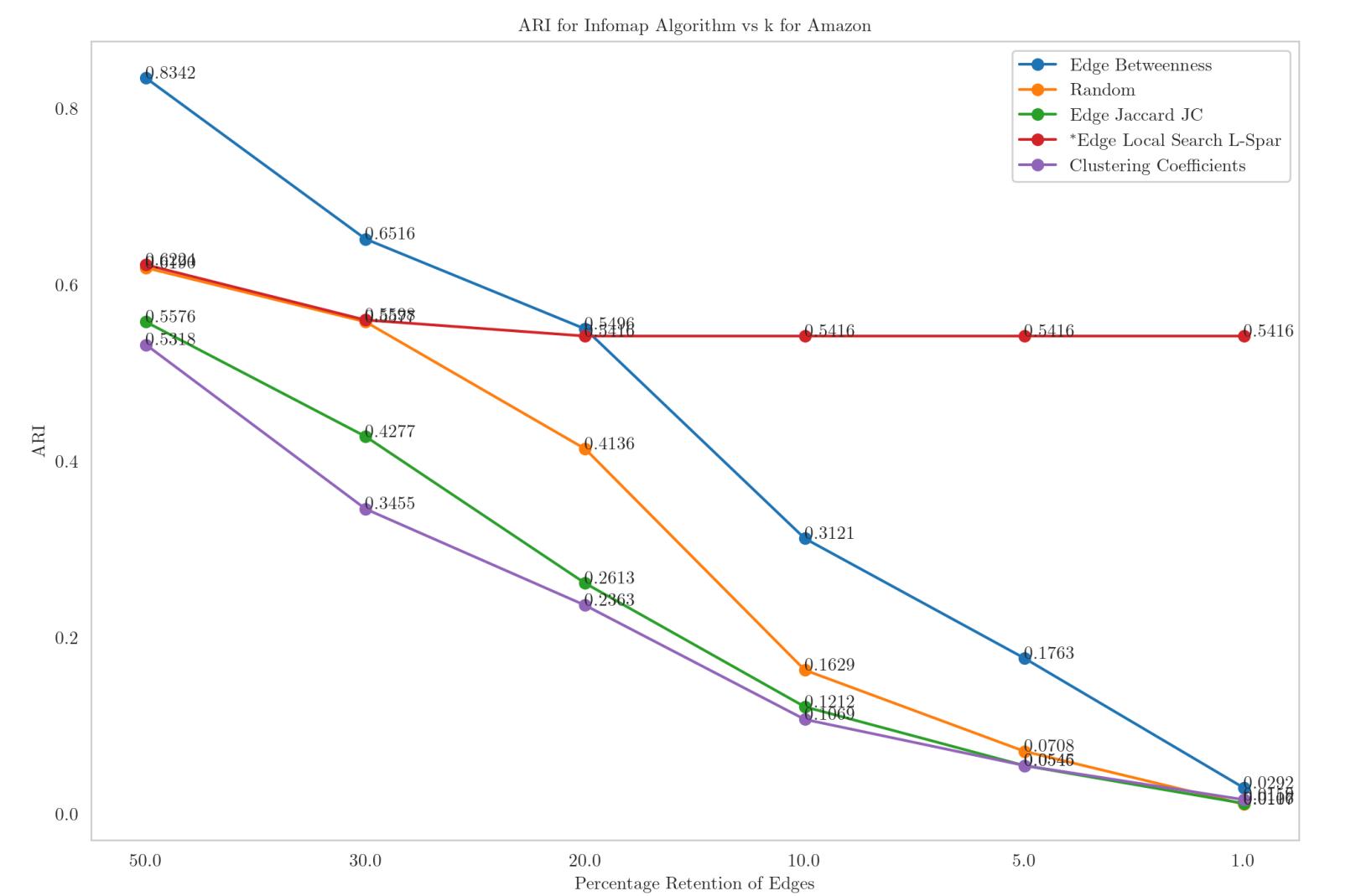


Modularity for Label Propogation Algorithm vs k for Amazon

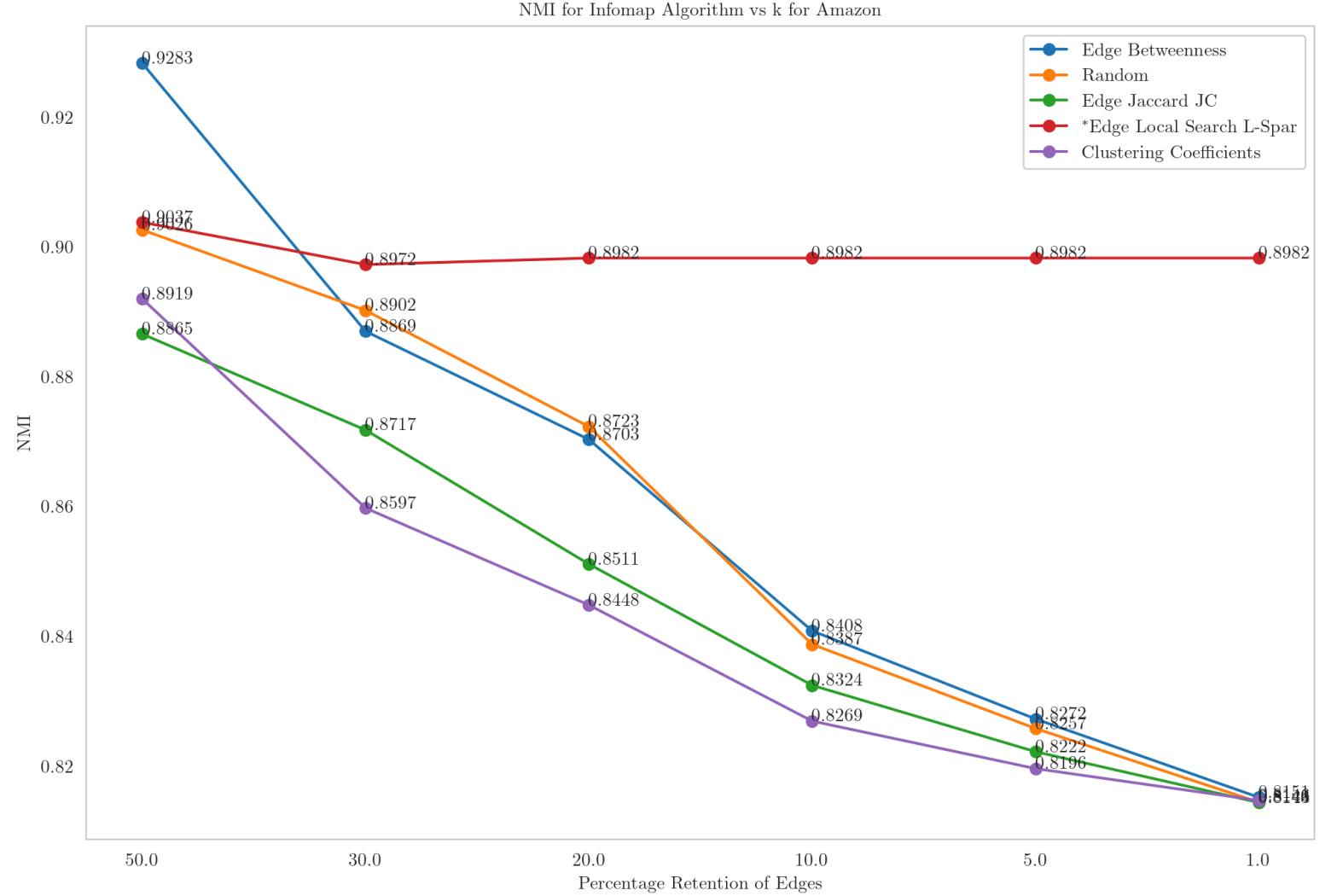
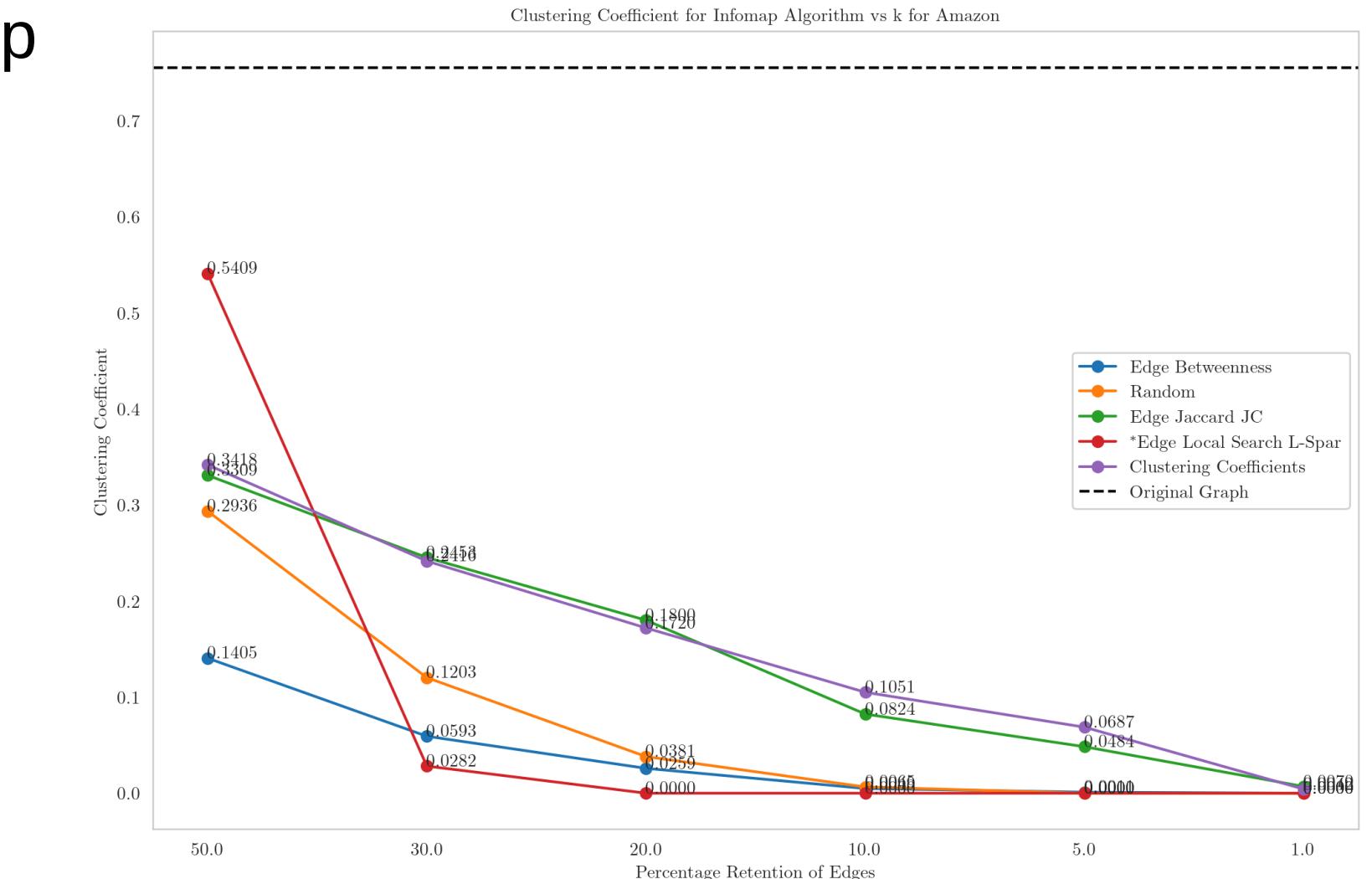


Amazon LPA



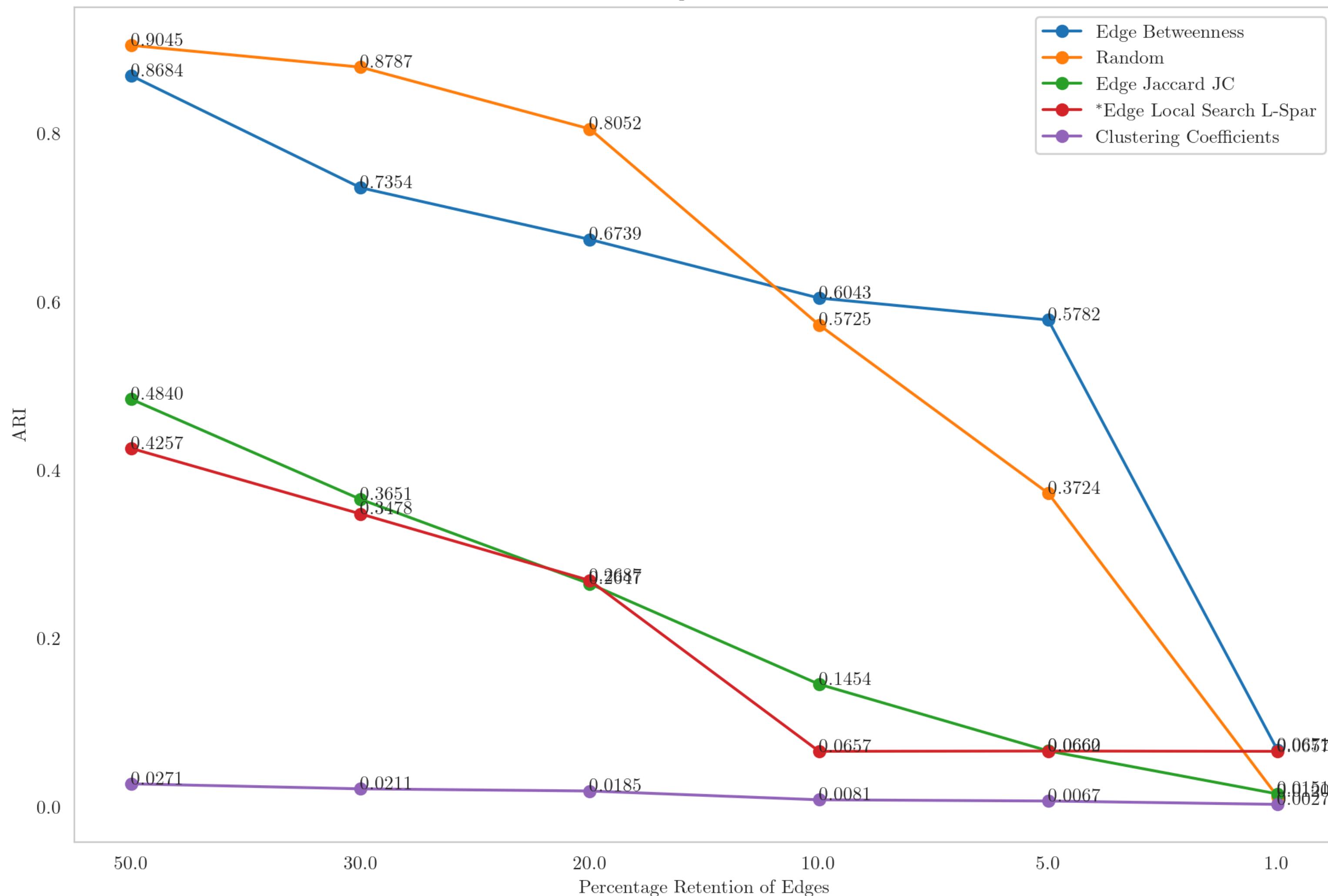


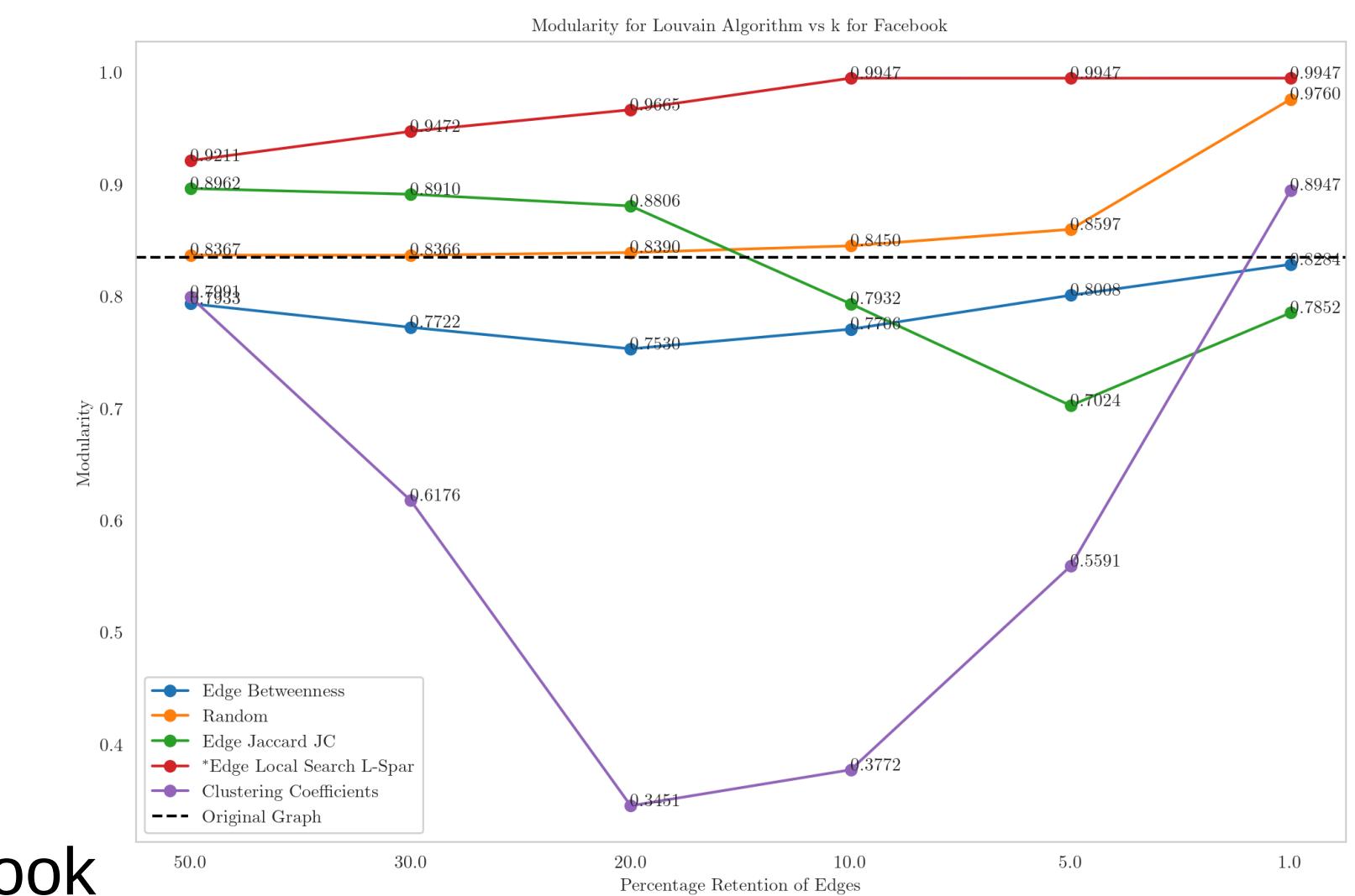
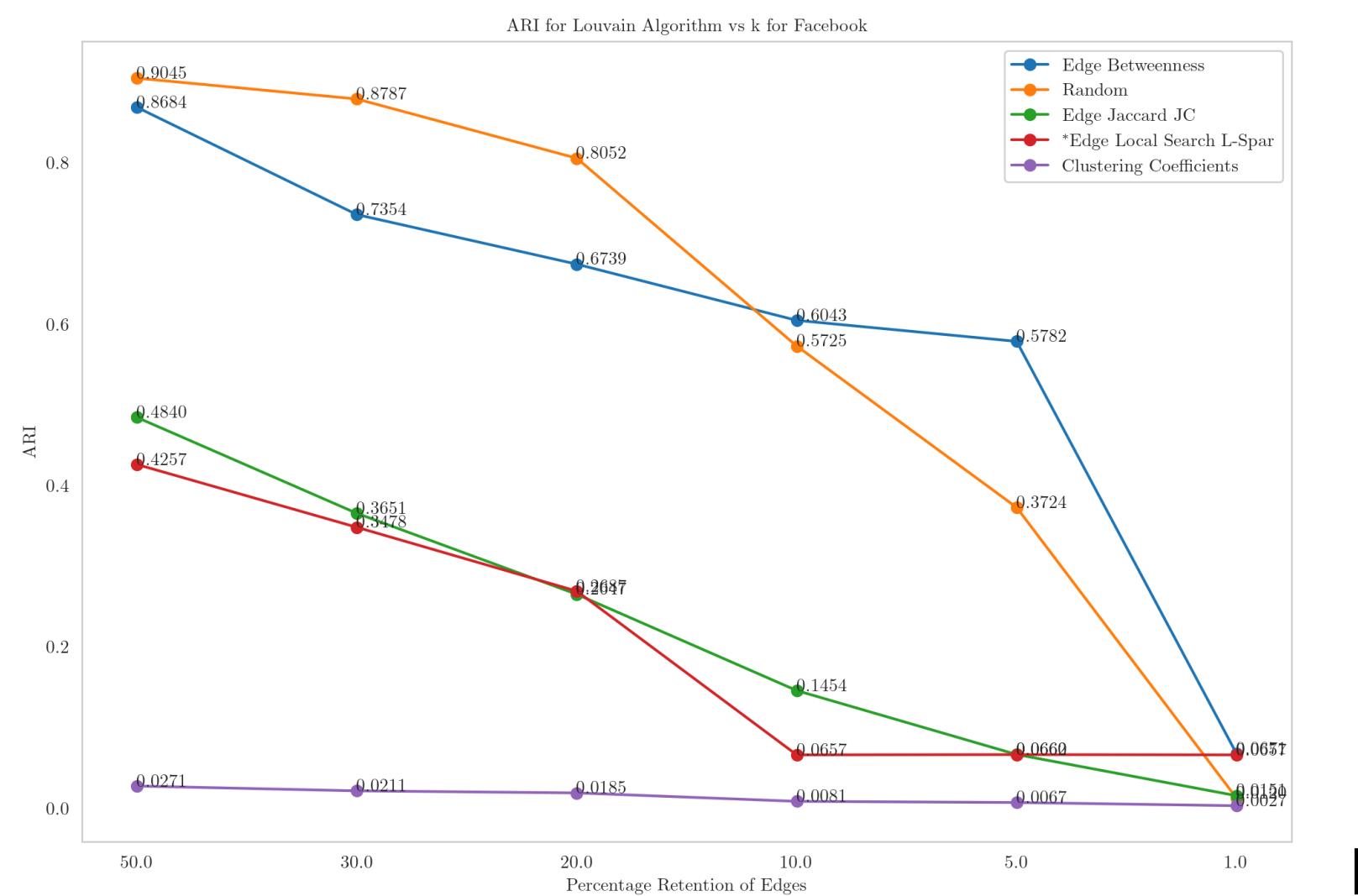
Amazon InfoMap



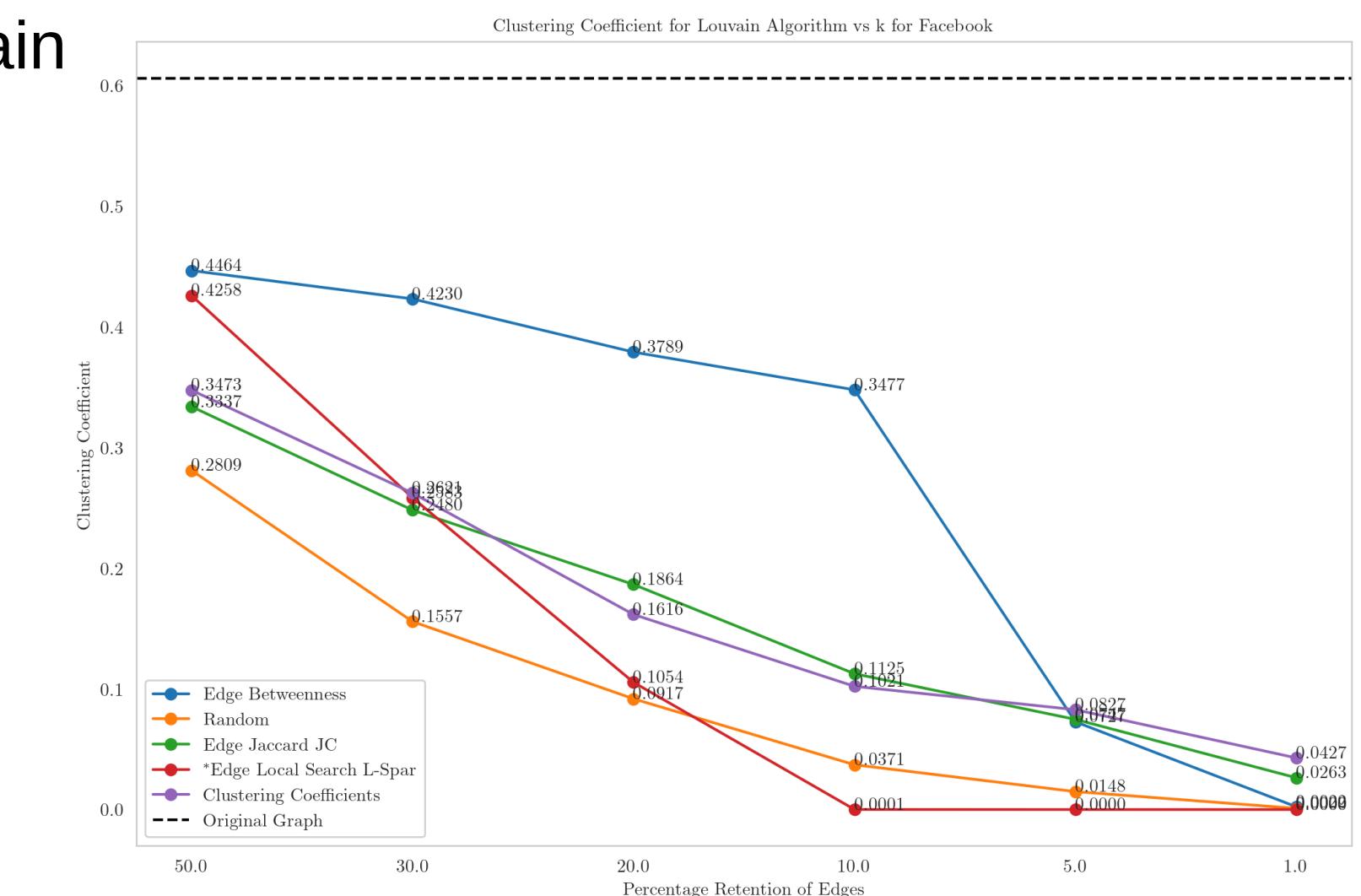
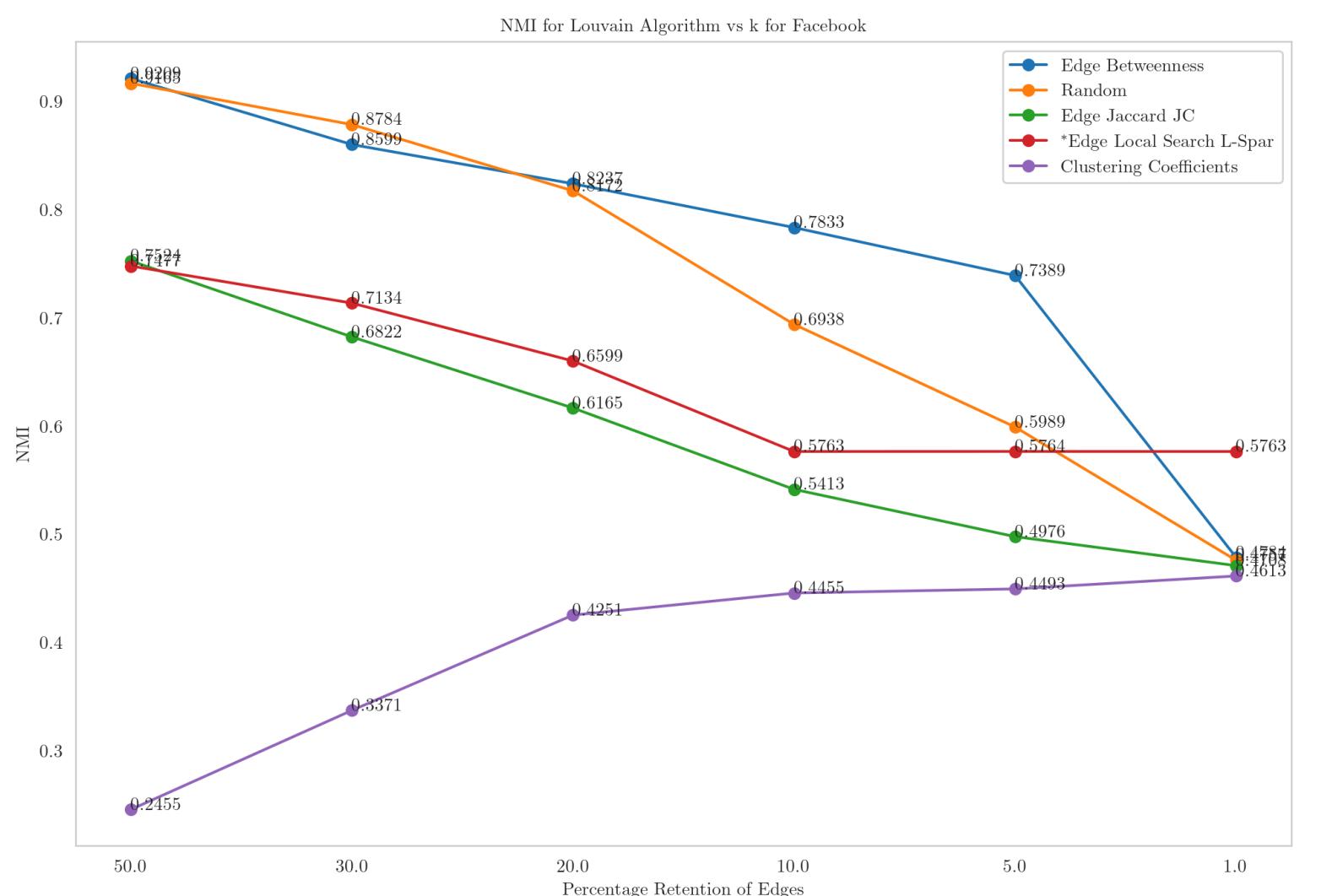
Results for Facebook Social Network

ARI for Louvain Algorithm vs k for Facebook

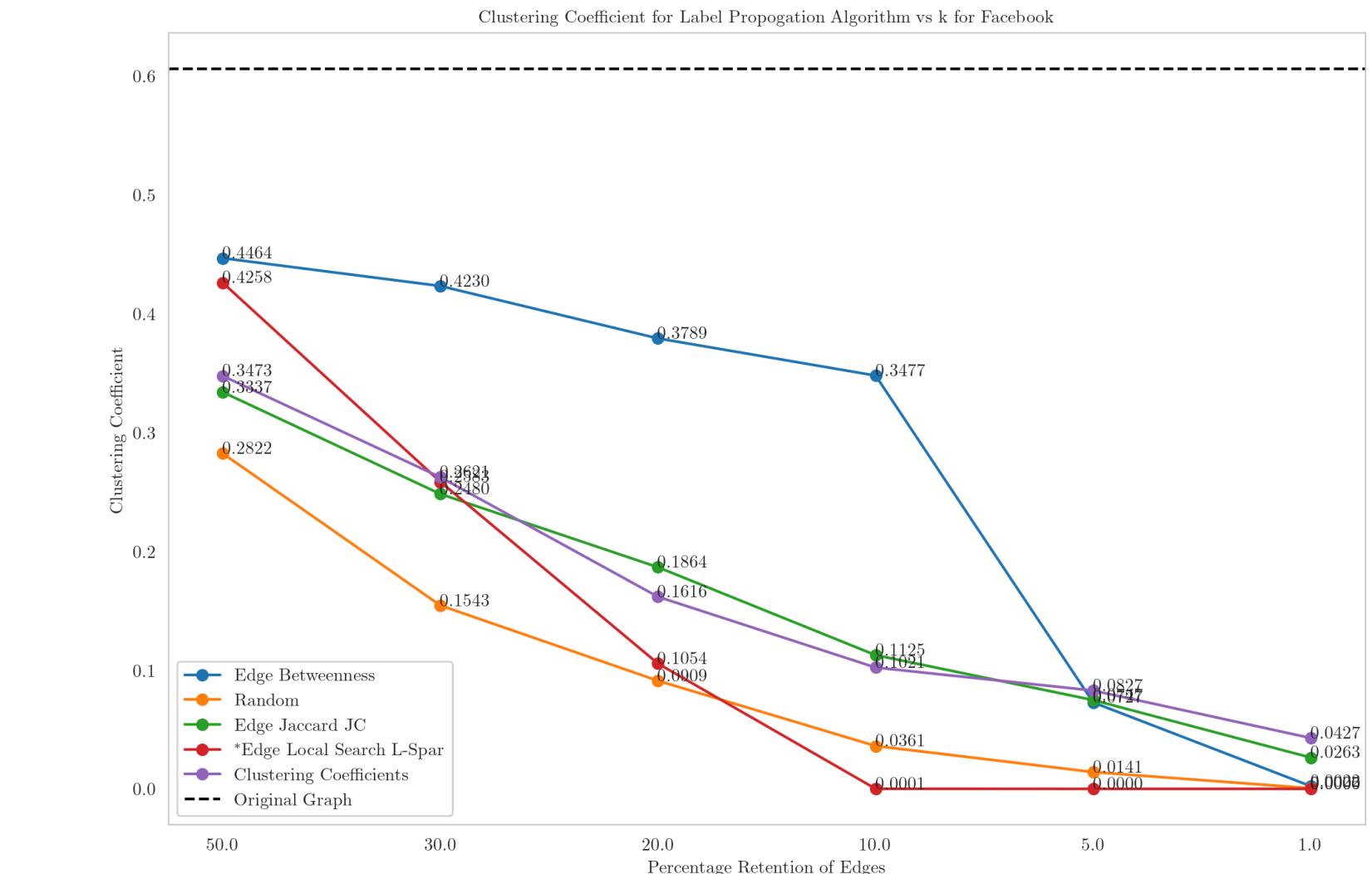
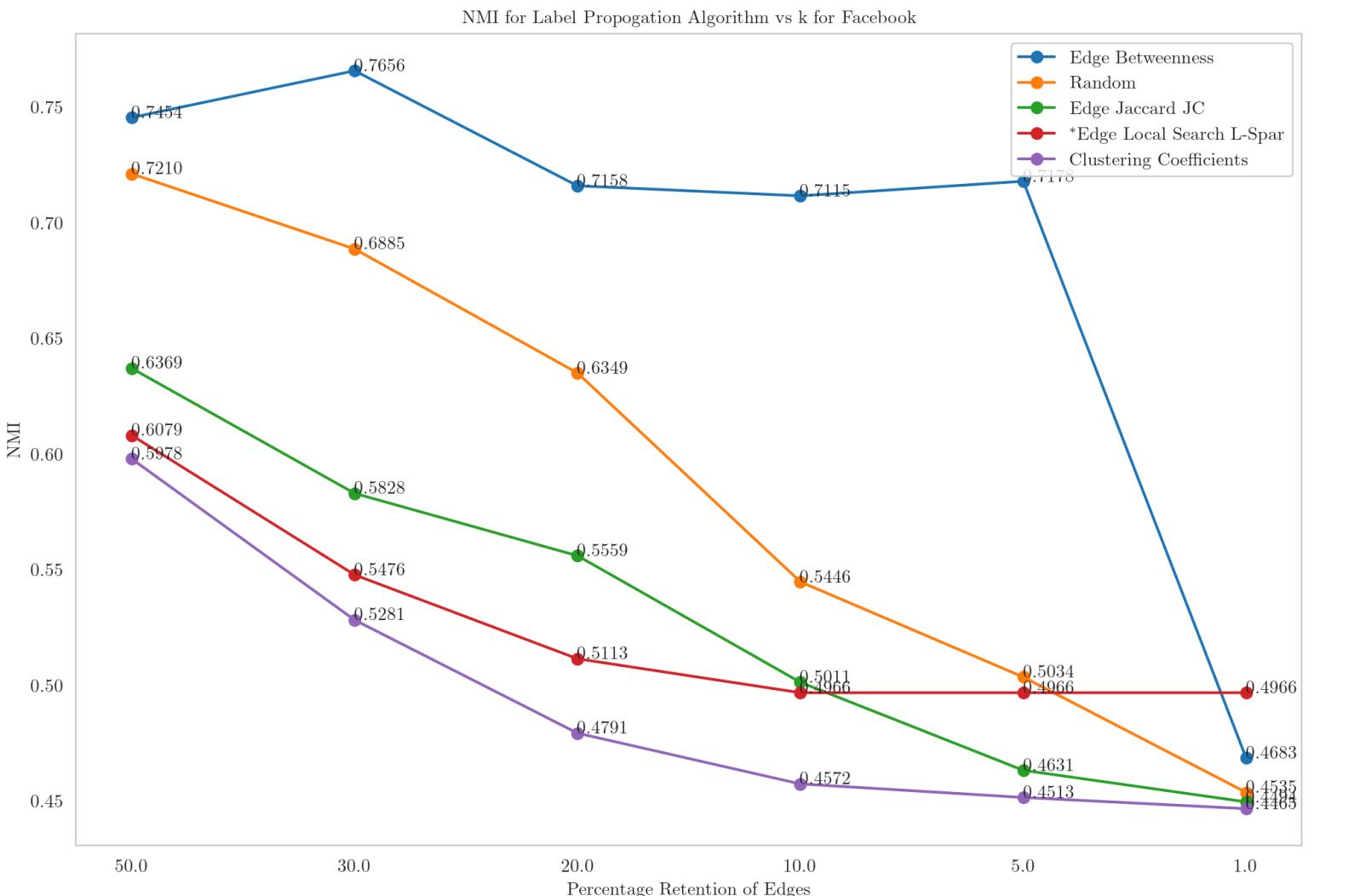
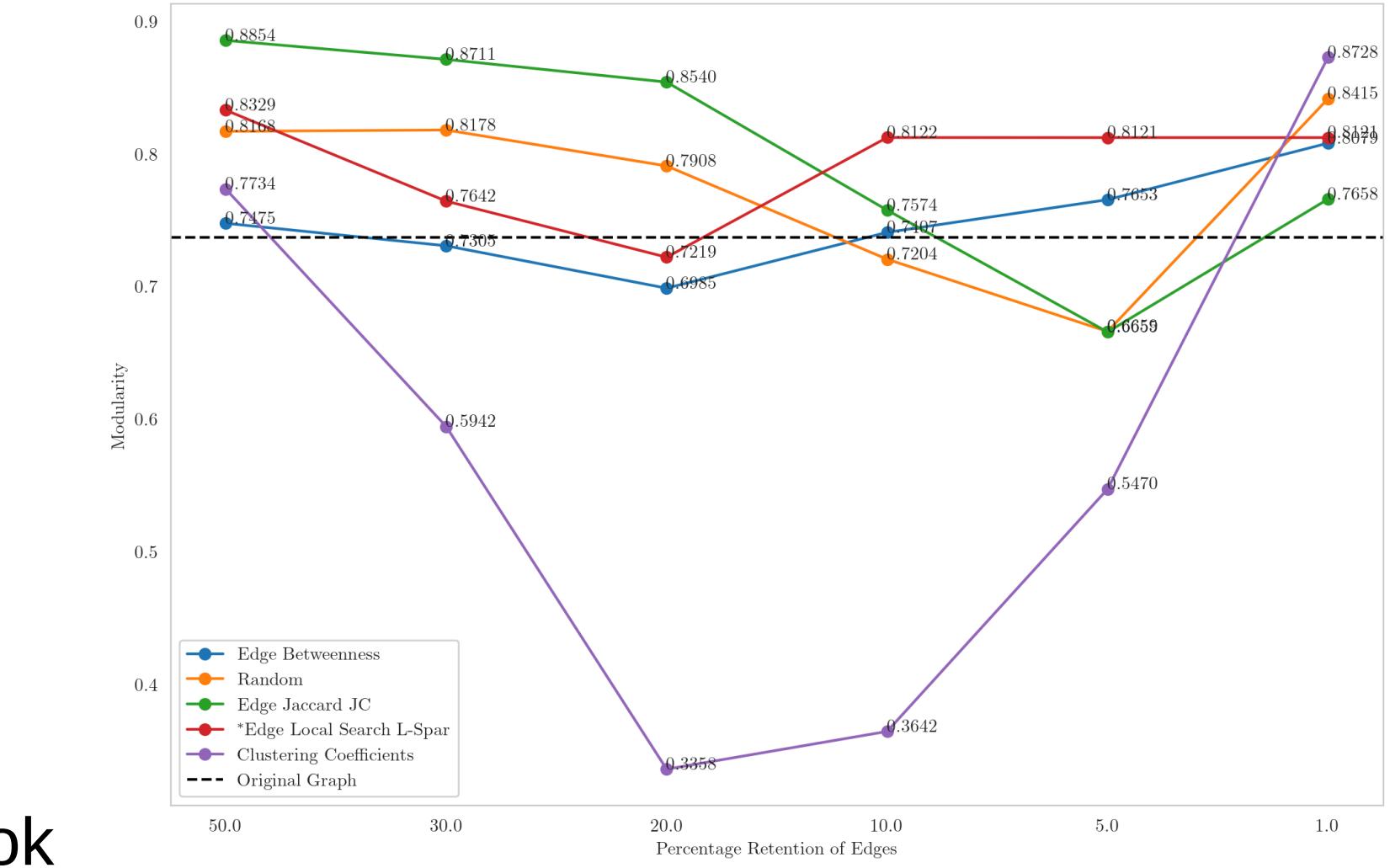
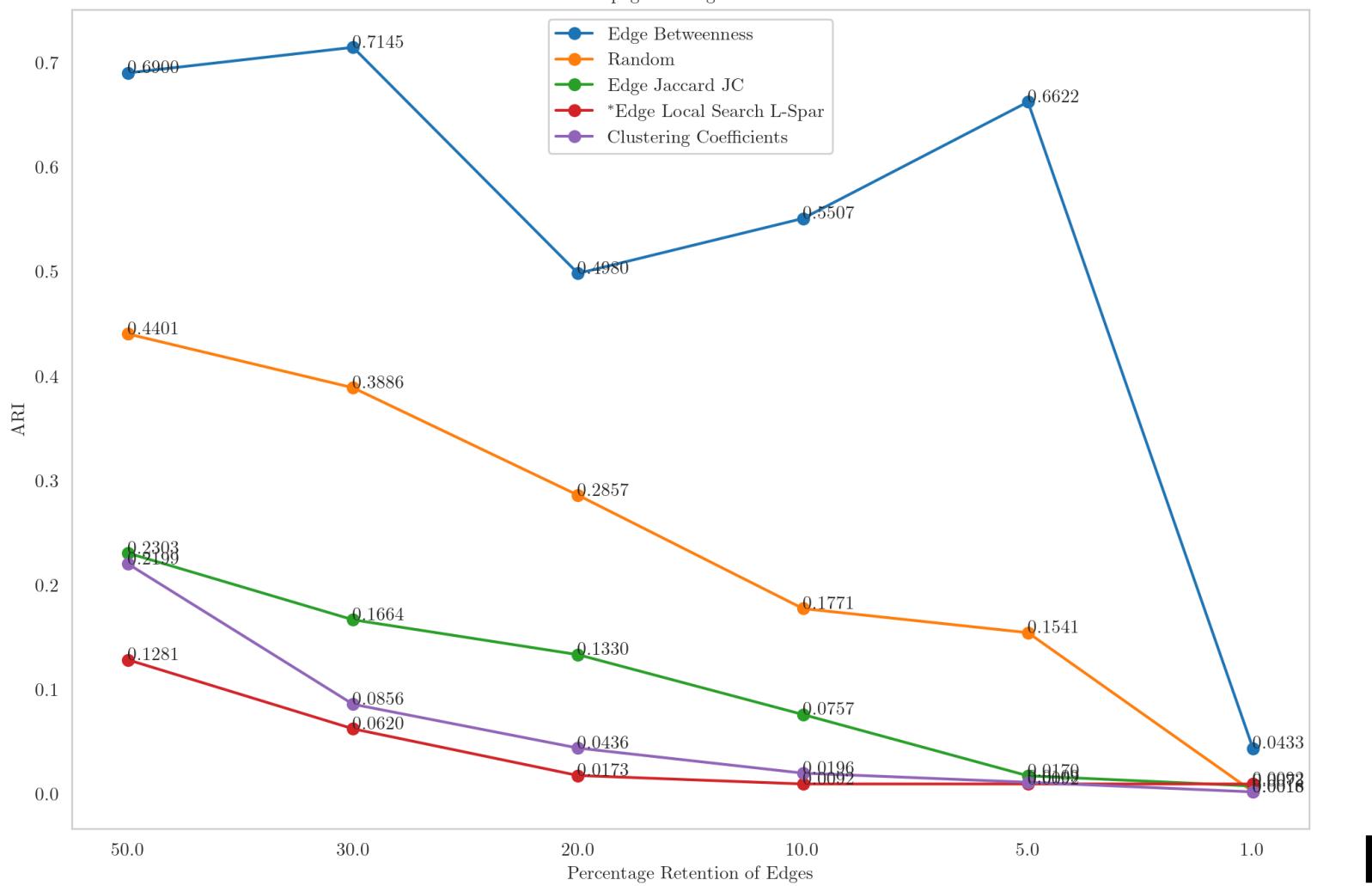




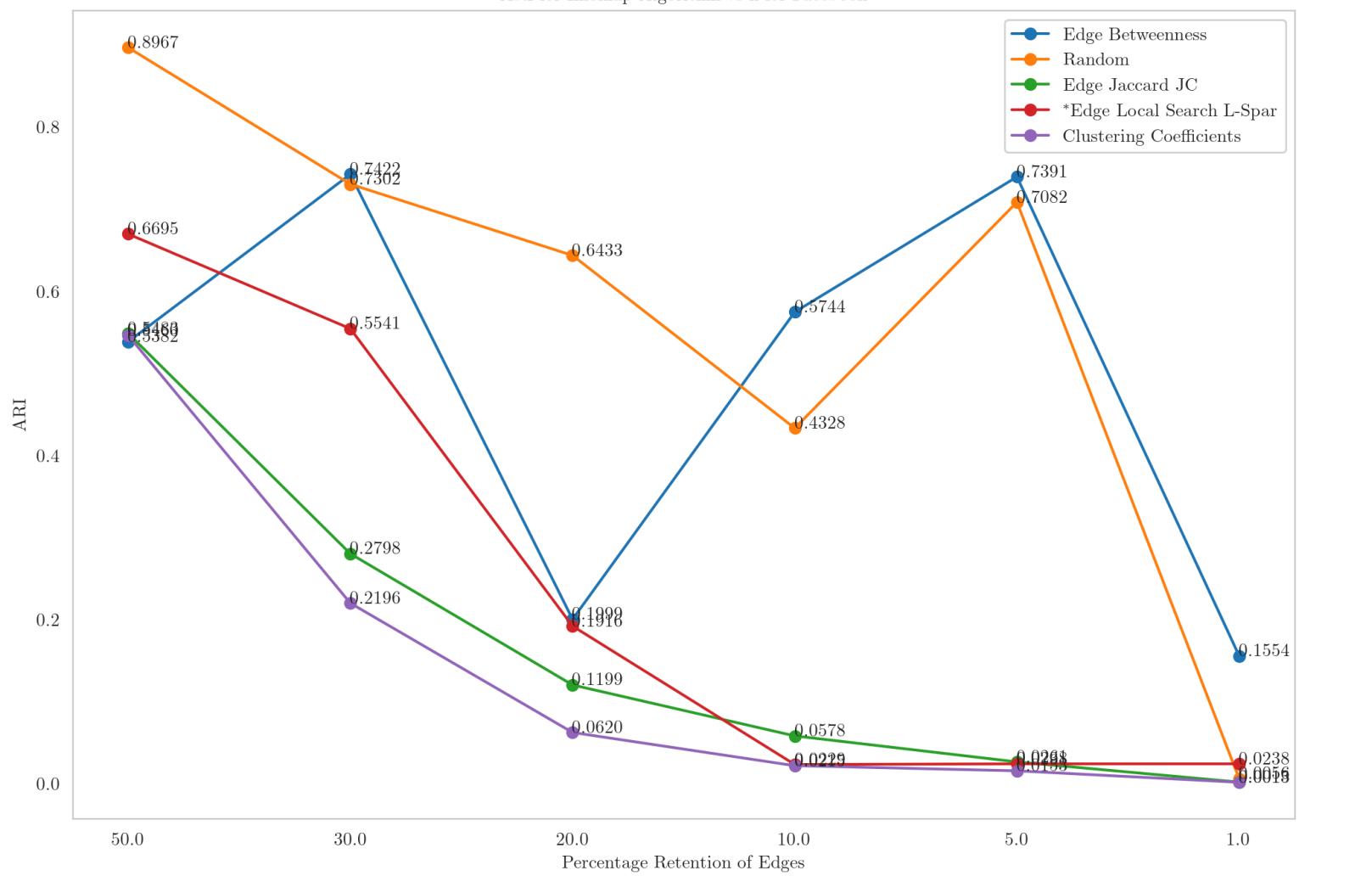
Facebook Louvain



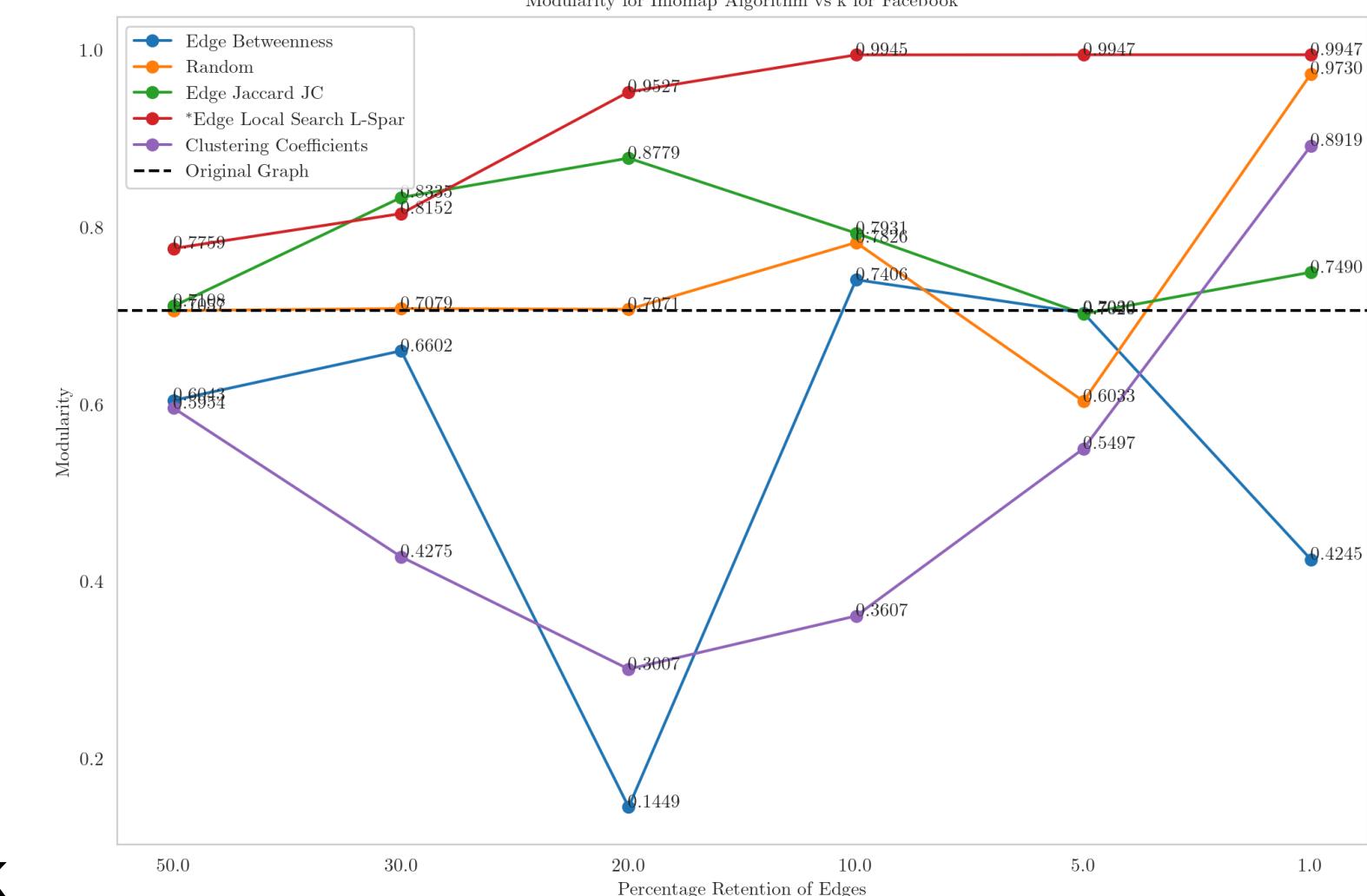
Facebook LPA



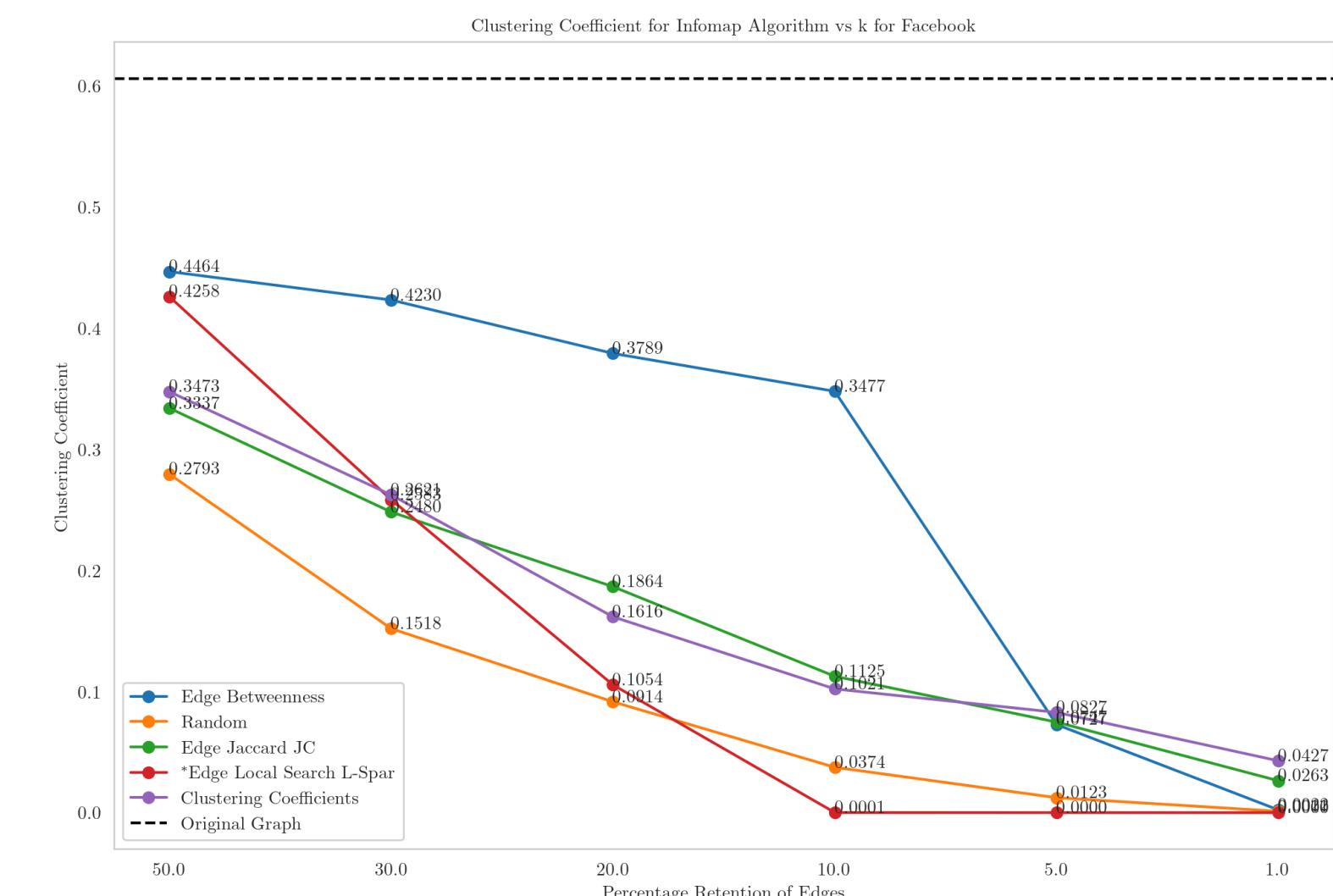
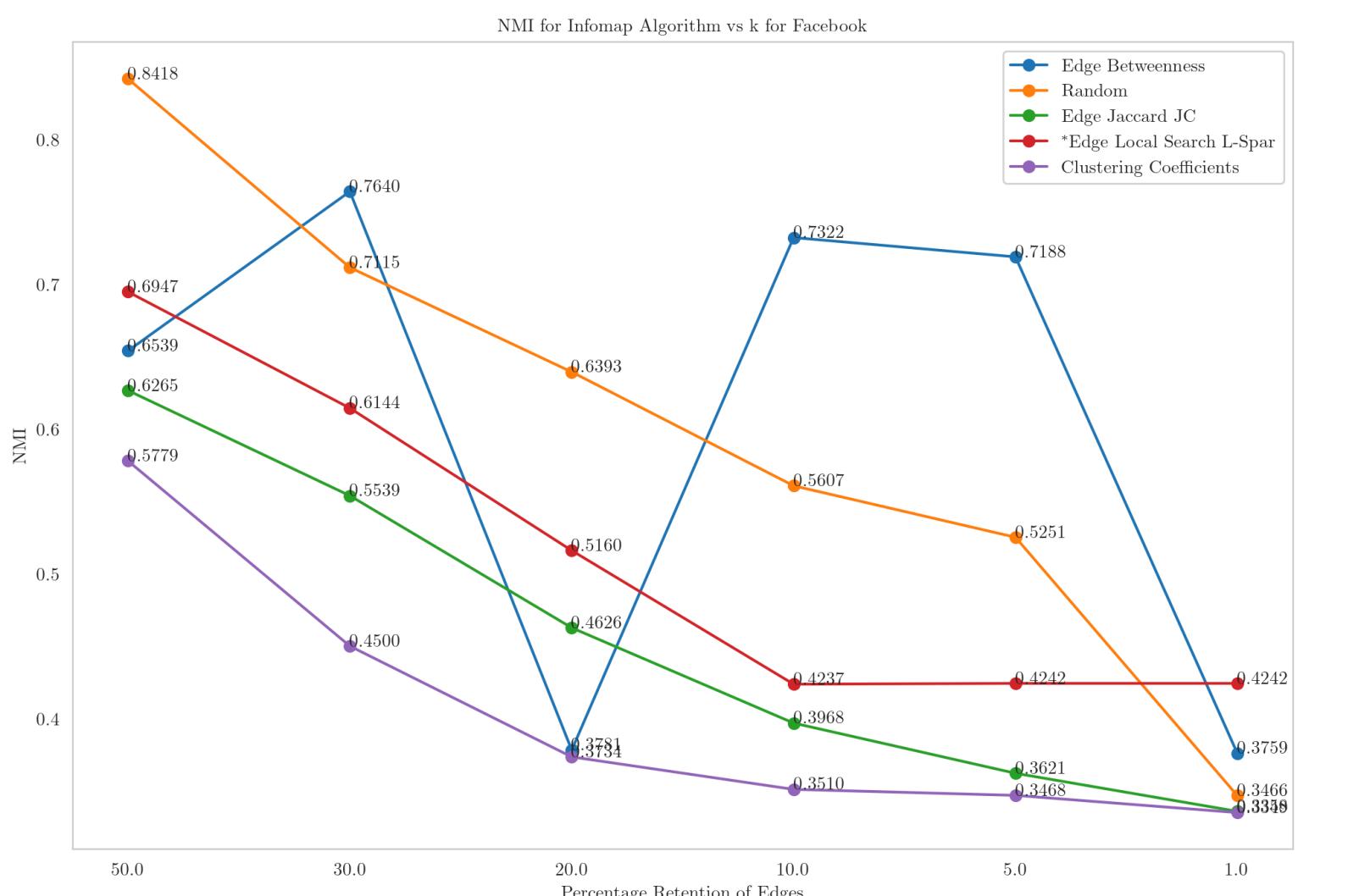
ARI for Infomap Algorithm vs k for Facebook



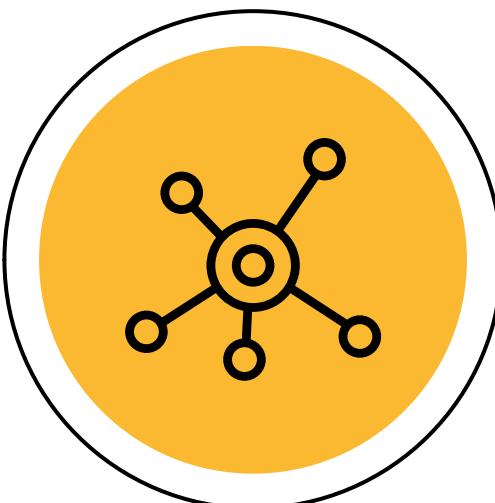
Modularity for Infomap Algorithm vs k for Facebook



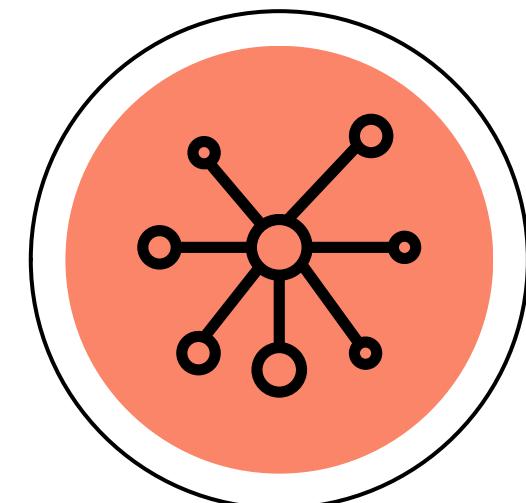
Facebook InfoMap



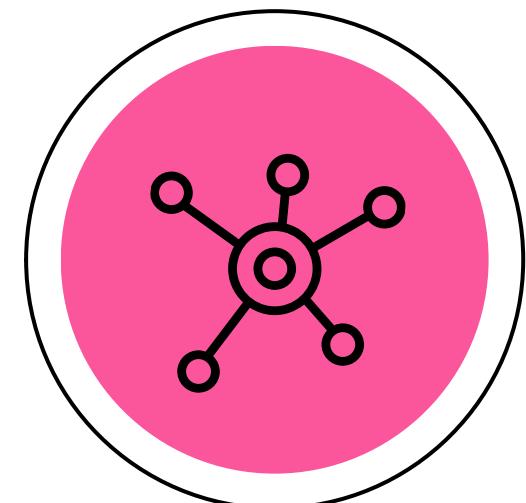
Some more sparsifying techniques



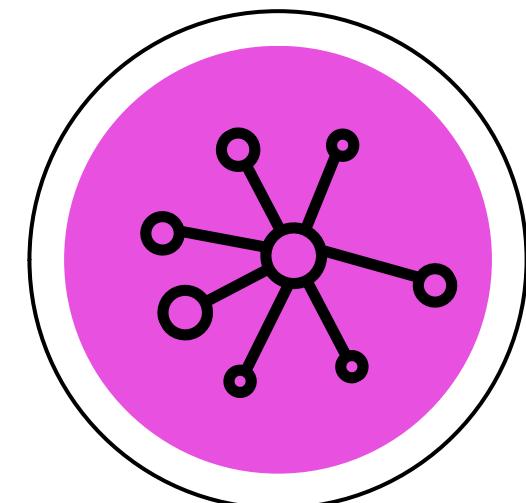
**Degree
Based Edge
Sampling**



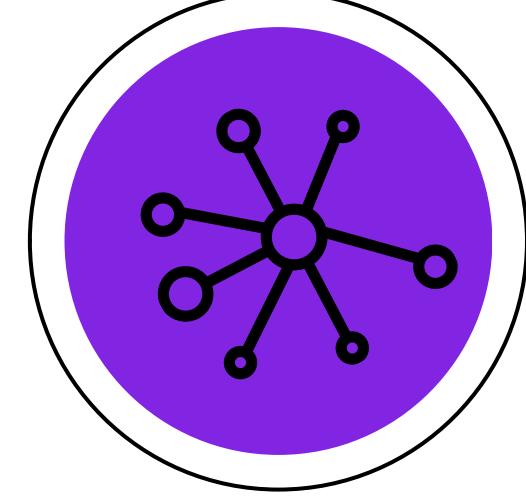
**Modified
Metropolis
Hastings**



**Forest Fire
Sampling**

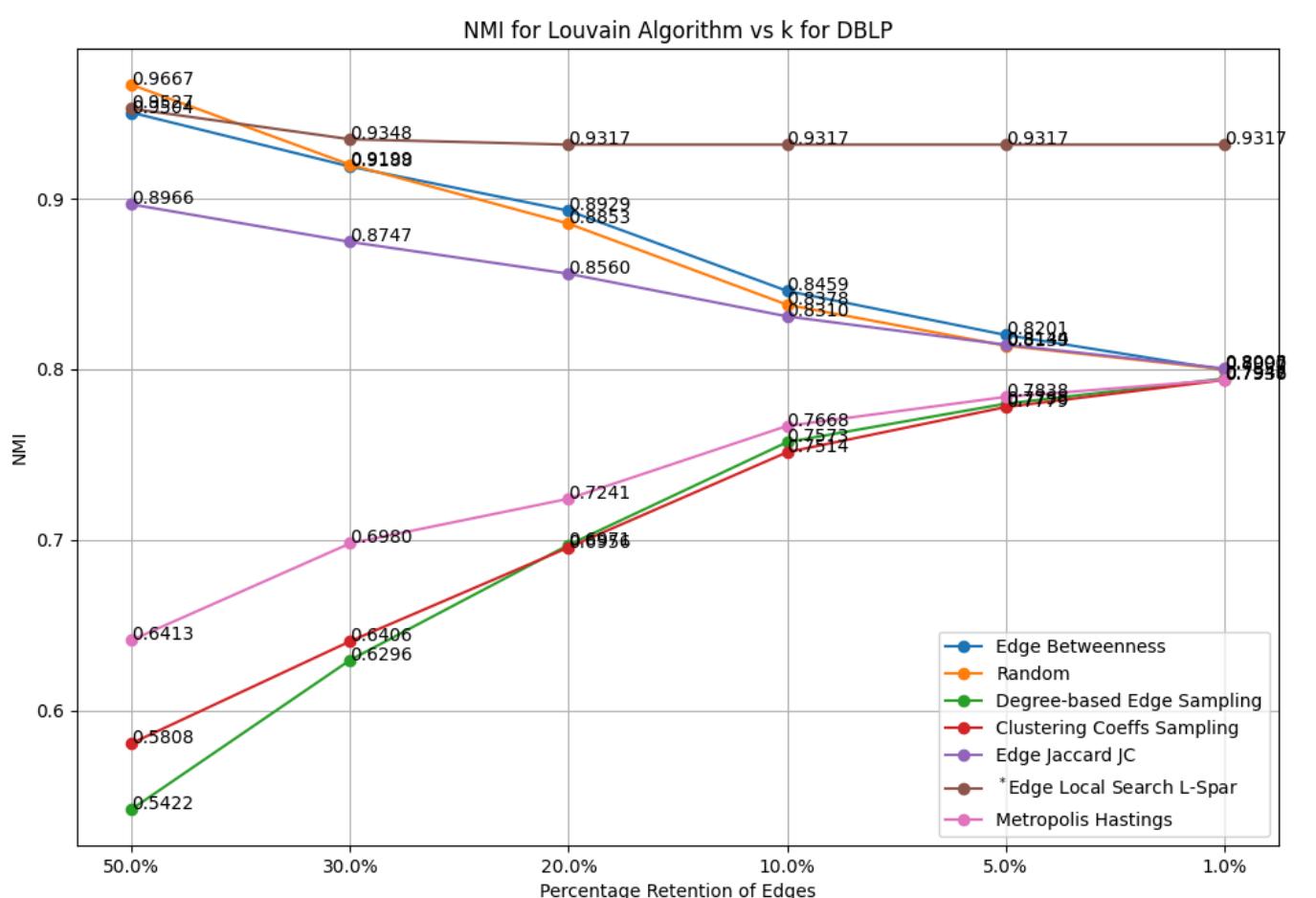
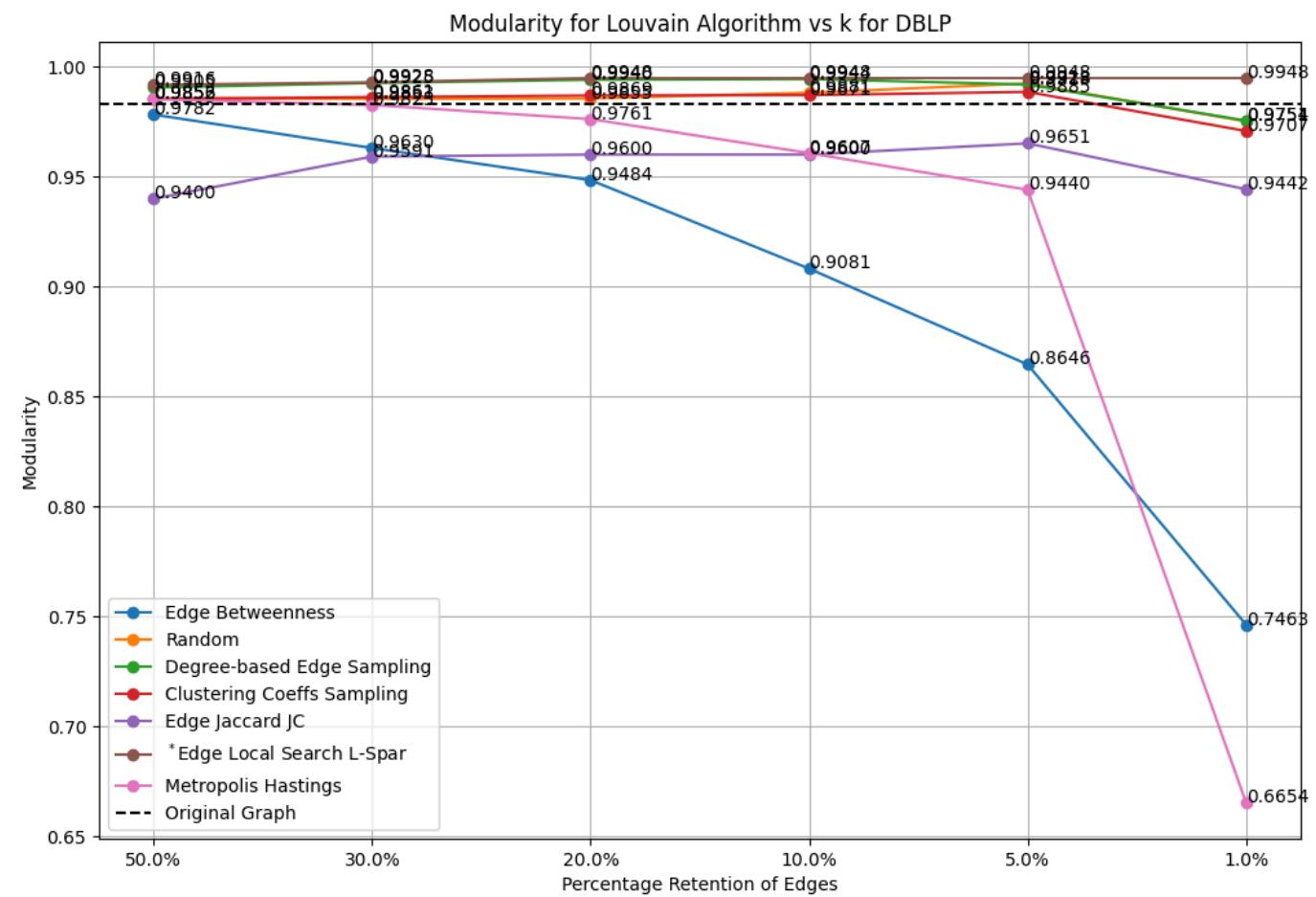
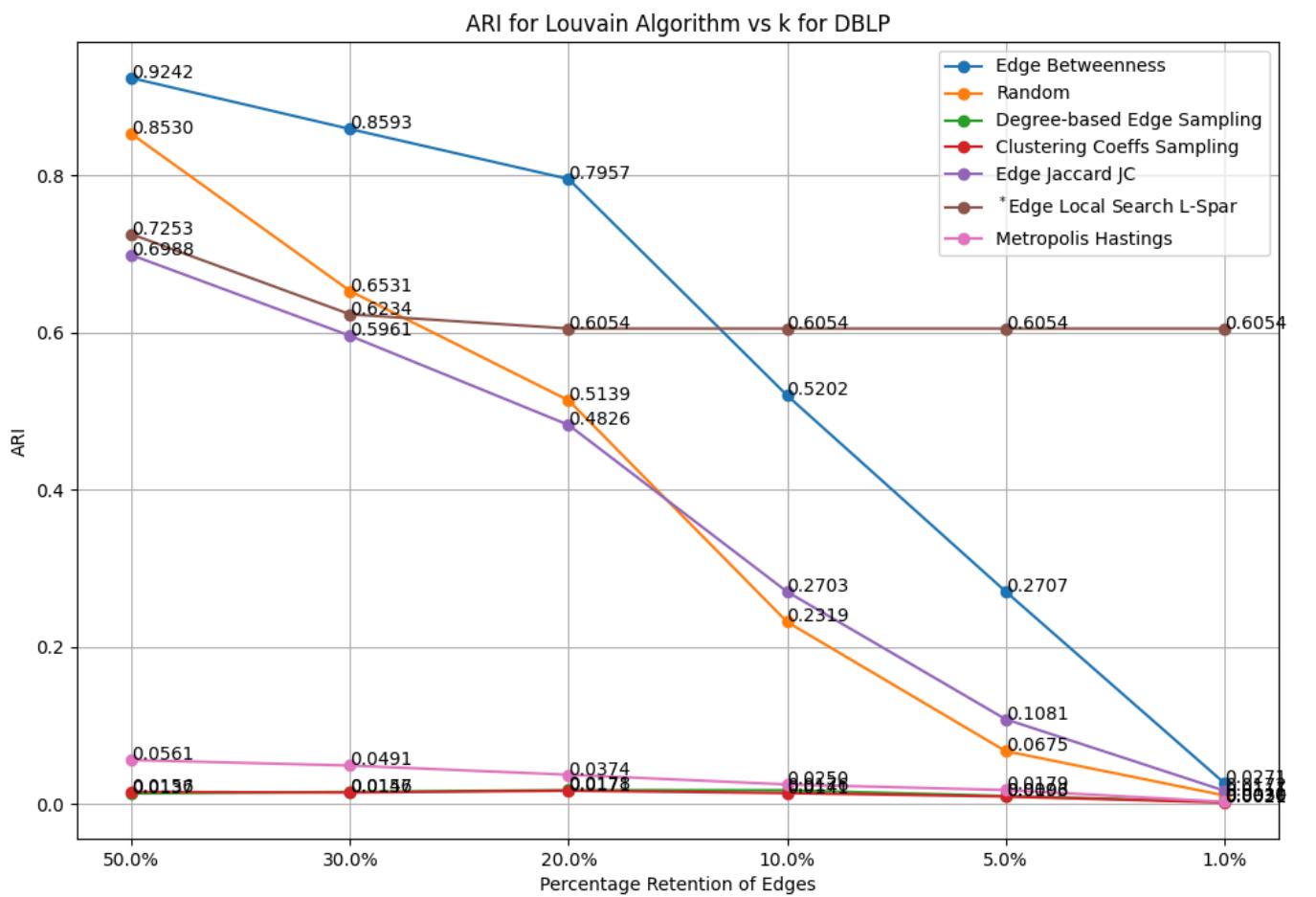


**Triangle Count
Based Sampling**

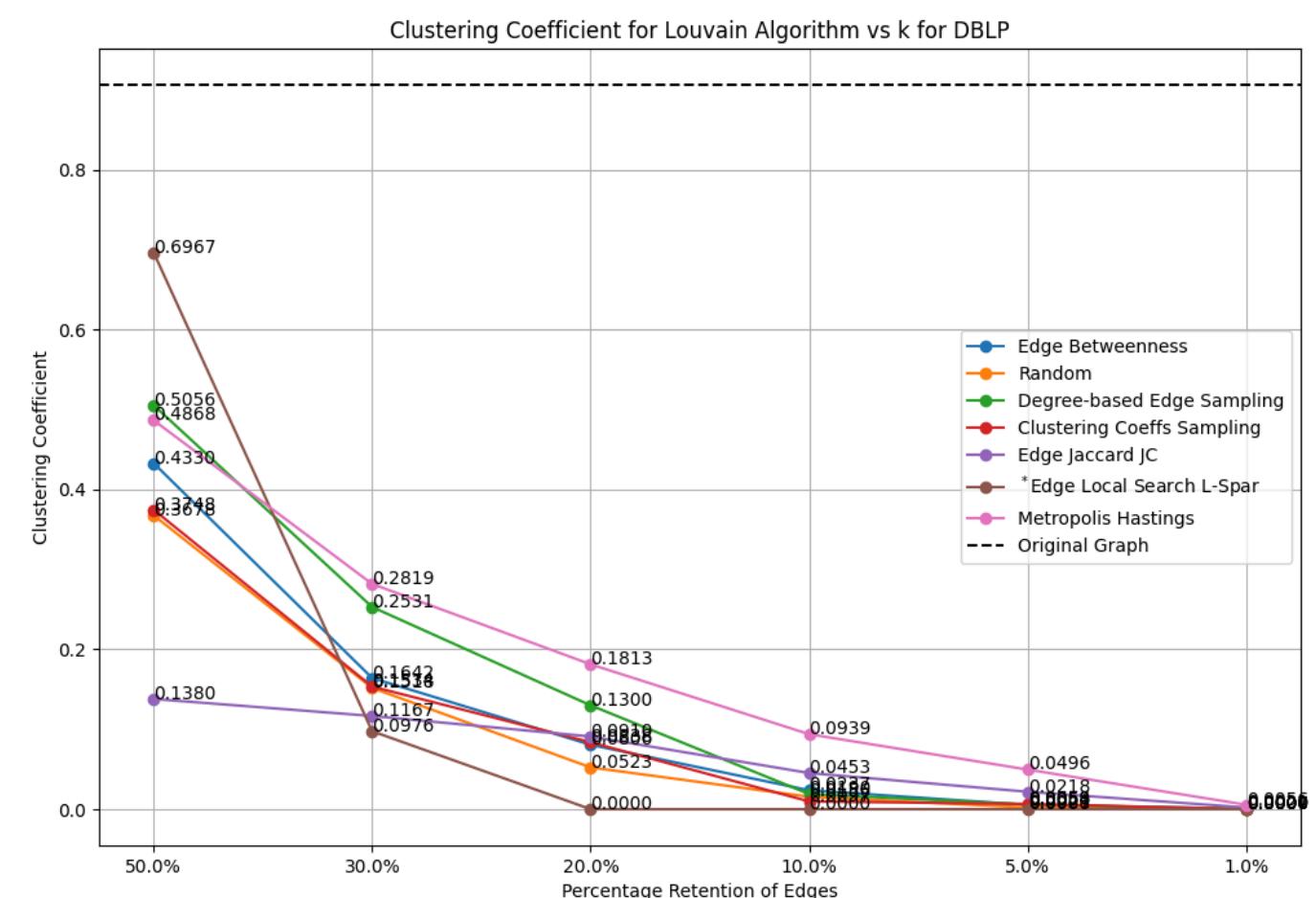


**Effective
Resistance**

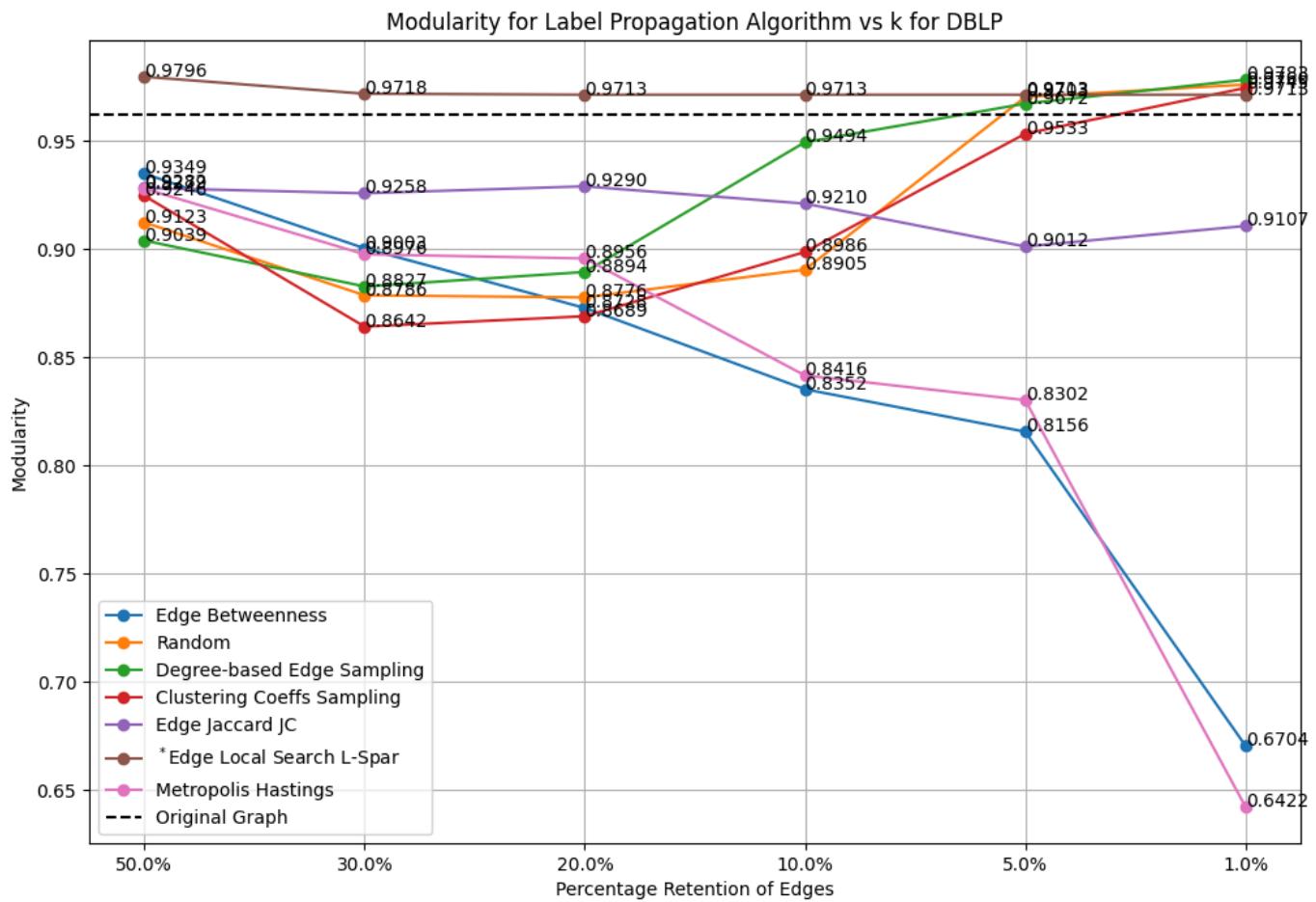
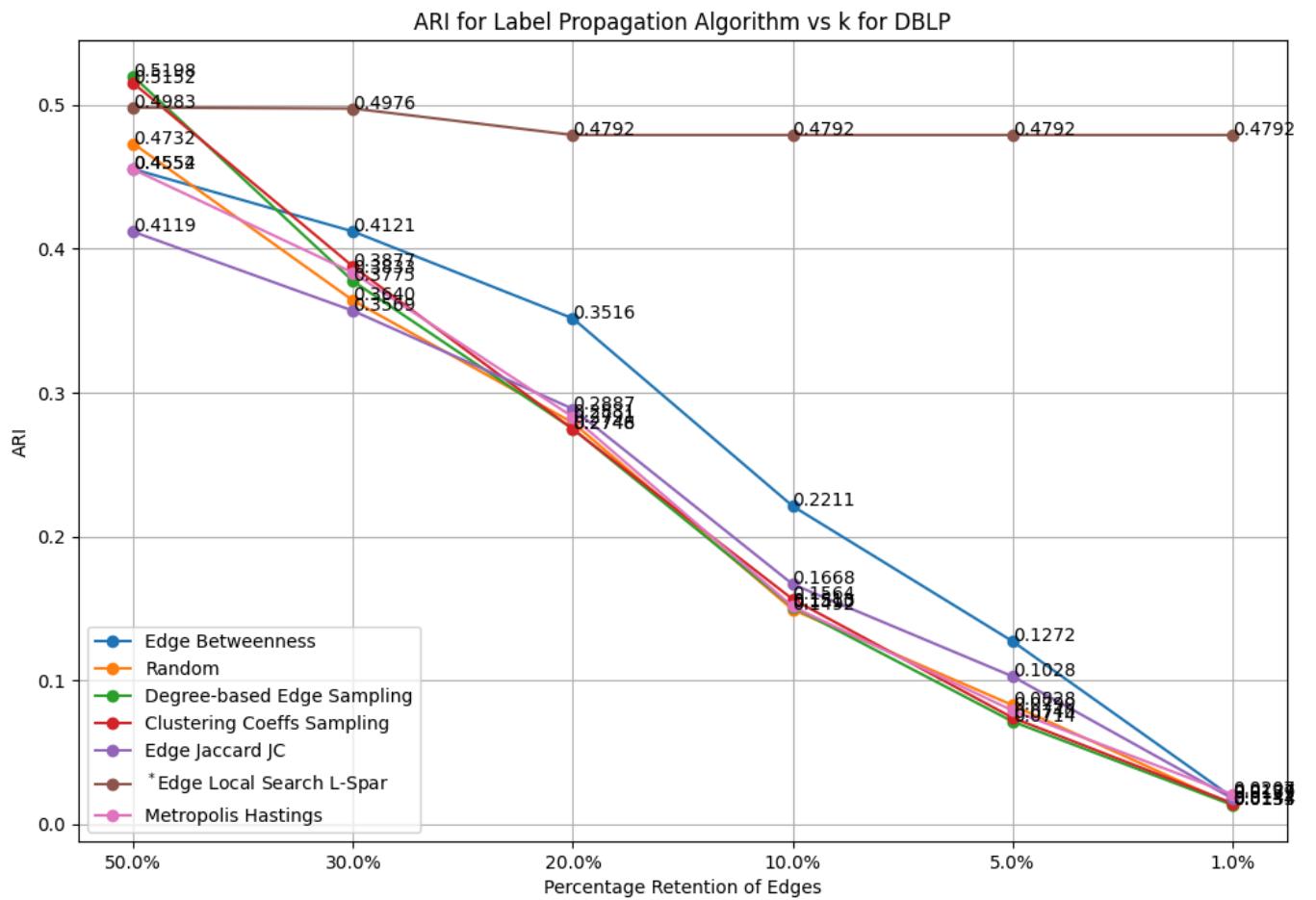
Glimpses...



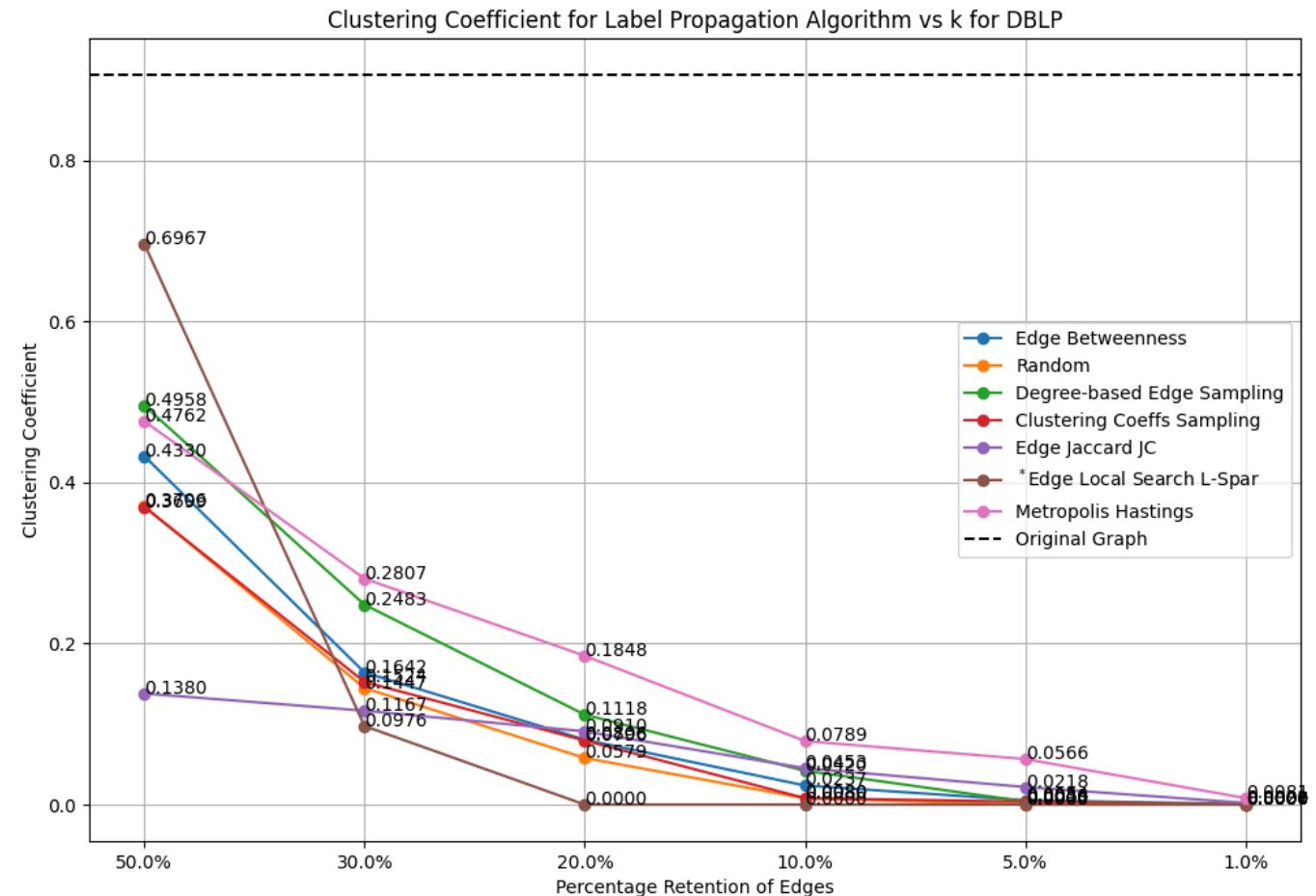
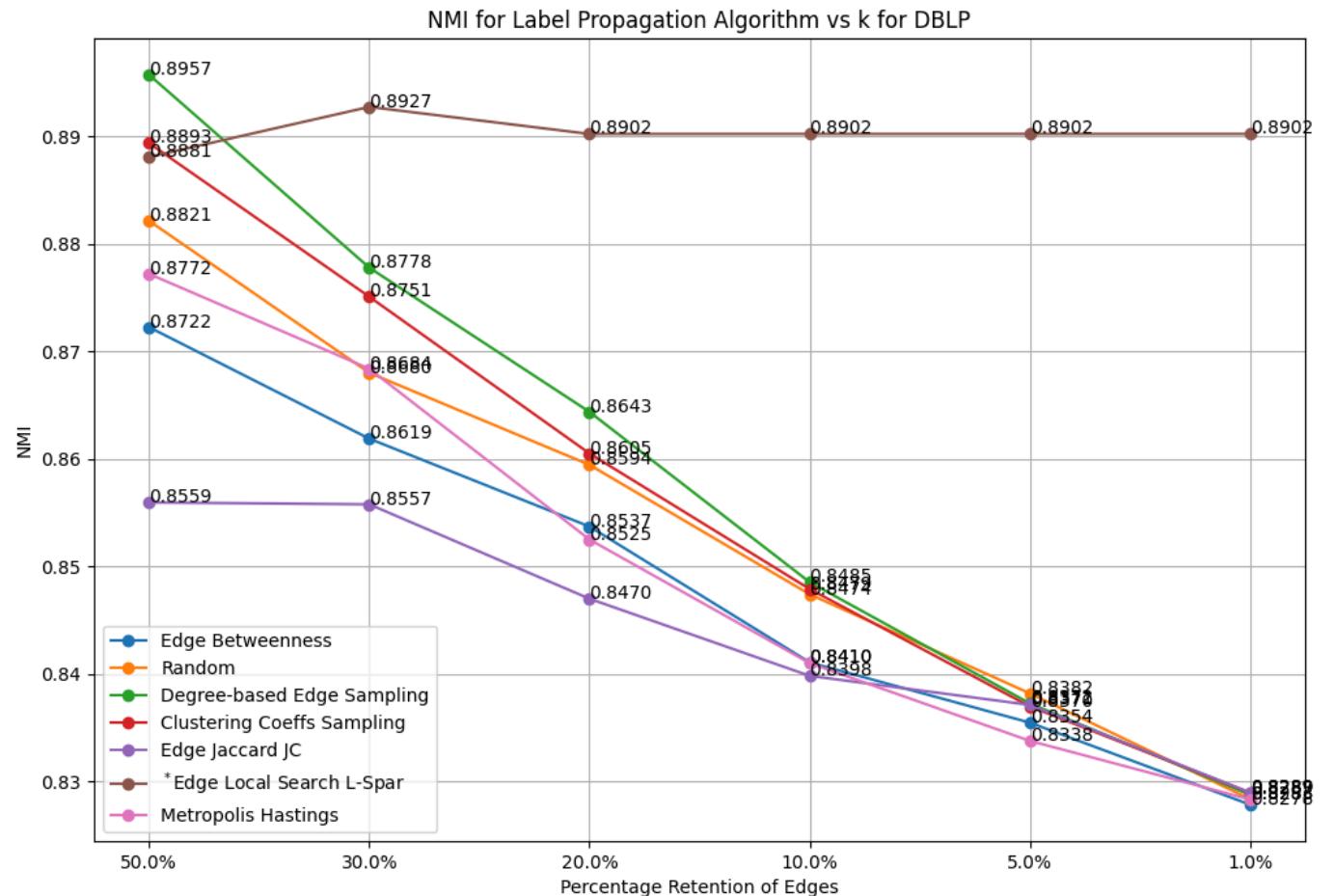
DBLP Louvain



Glimpses...



DBLP Label Propagation



References

- [1.] Hu, P., & Lau, W. C. (2013). A survey and taxonomy of graph sampling. arXiv preprint arXiv:1308.5865.
- [2.] Satuluri, V., Parthasarathy, S., & Ruan, Y. (2011, June). Local graph sparsification for scalable clustering. In Proceedings of the 2011 ACM SIGMOD International Conference on Management of data (pp. 721-732).
- [3.] Wu, H. Y., & Chen, Y. L. (2020, November). Graph sparsification with generative adversarial network. In 2020 IEEE International Conference on Data Mining (ICDM) (pp. 1328-1333). IEEE.
- [4.] Spielman, D. A., & Srivastava, N. (2008, May). Graph sparsification by effective resistances. In Proceedings of the fortieth annual ACM symposium on Theory of computing (pp. 563-568).
- [5.] Wickman, R., Zhang, X., & Li, W. (2021). A Generic Graph Sparsification Framework using Deep Reinforcement Learning. arXiv preprint arXiv:2112.01565.
- [6.] Hamann, M., Lindner, G., Meyerhenke, H., Staudt, C. L., & Wagner, D. (2016). Structure-preserving sparsification methods for social networks. Social Network Analysis and Mining, 6, 1-22.
- [7.] Yang, J., & Leskovec, J. (2012, August). Defining and evaluating network communities based on ground-truth. In Proceedings of the ACM SIGKDD Workshop on Mining Data Semantics (pp. 1-8).
- [8.] Lindner, G., Staudt, C. L., Hamann, M., Meyerhenke, H., & Wagner, D. (2015, August). Structure-preserving sparsification of social networks. In Proceedings of the 2015 IEEE/ACM International conference on advances in social networks analysis and mining 2015 (pp. 448-454).

Graph Libraries Used

[1.] NetworKit

[2.] Networkx

[3.] CDlib

[4.] Pyvis

Community Structure in Large Networks: Natural Cluster Sizes and the Absence of Large Well-Defined Clusters

Jure Leskovec, Kevin J. Lang, Anirban Dasgupta, and Michael W. Mahoney

Abstract. A large body of work has been devoted to defining and identifying clusters or communities in social and information networks, i.e., in graphs in which the nodes represent underlying social entities and the edges represent some sort of interaction between pairs of nodes. Most such research begins with the premise that a community or a cluster should be thought of as a set of nodes that has more and/or better connections between its members than to the remainder of the network. In this paper, we explore from a novel perspective several questions related to identifying meaningful communities in large social and information networks, and we come to several striking conclusions.

Rather than defining a procedure to extract sets of nodes from a graph and then attempting to interpret these sets as “real” communities, we employ approximation algorithms for the graph-partitioning problem to characterize as a function of size the statistical and structural properties of partitions of graphs that could plausibly be interpreted as communities. In particular, we define the *network community profile plot*, which characterizes the “best” possible community—according to the conductance measure—over a wide range of size scales. We study over one hundred large real-world networks, ranging from traditional and online social networks, to technological and information networks and web graphs, and ranging in size from thousands up to tens of millions of nodes.

Our results suggest a significantly more refined picture of community structure in large networks than has been appreciated previously. Our observations agree with previous work on small networks, but we show that large networks have a very different structure. In particular, we observe tight communities that are barely connected to the rest of the network at very small size scales (up to ≈ 100 nodes); and communities of size scale beyond ≈ 100 nodes gradually “blend into” the expander-like core of the network and thus become less “community-like,” with a roughly inverse relationship between community size and optimal community quality. This observation agrees well