

SCAN: A Structural Clustering Algorithm for Networks

Xiaowei Xu
Zhidan Feng

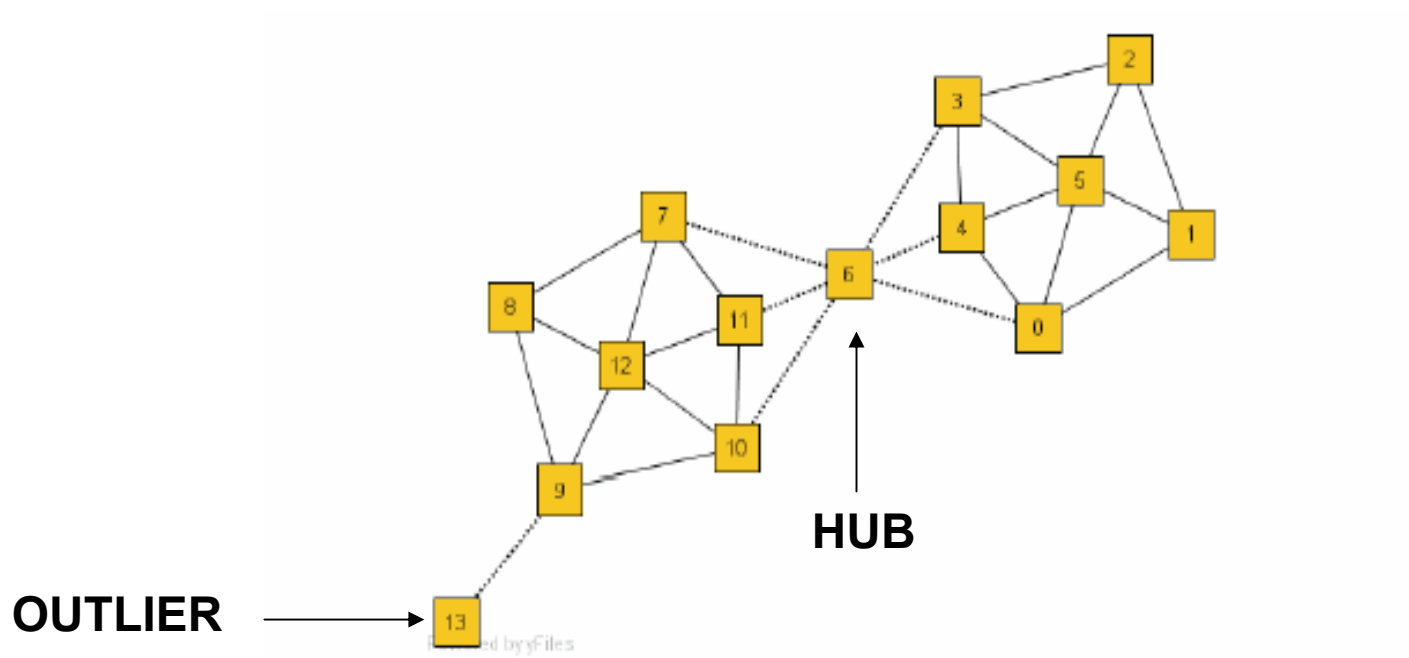
Nurcan Yuruk
Thomas Schweiger

Presented By
Maen Hammad

Features

- Identify Clusters
- Fast algorithm $O(m)$
- Identify hubs and outliers nodes

Example



Related Work

1. Min-max cut method

- Partition the graph into 2 clusters
- Minimize edges between the clusters
- Maximize edges within a cluster
- Optimization → Find the best cut (difficult ?)
- Some constraints to achieve optimality

Related work

2. Modularity clustering

$$Q = \sum_{s=1}^k \left[\frac{l_s}{L} - \left(\frac{d_s}{2L} \right)^2 \right]$$

L : number of edges

L_s : number of edges in cluster S

d_s : sum of degrees in cluster s

- Optimality \rightarrow Maximize $Q \rightarrow$ NP Complete
- Greedy method is used
- Running time is $O(md \log n)$

Structure-Connected Clusters

- Considering the neighborhood around two connected vertices
- Two vertices are assigned to a cluster according to how they share neighbors
- Some neighbors are outliers or hubs

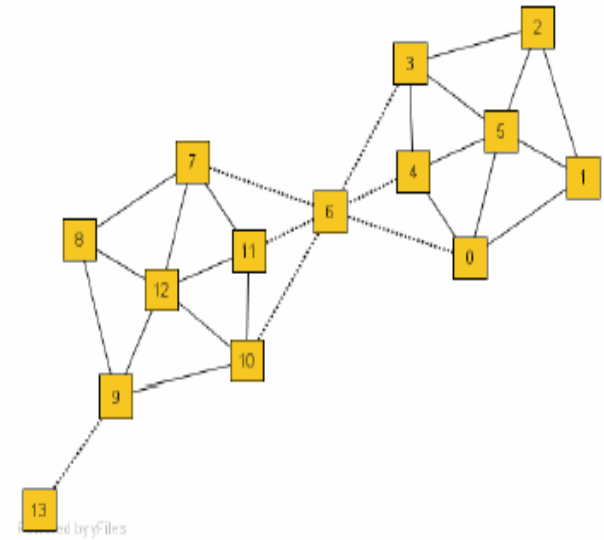
Definitions

- VERTEX STRUCTURE

$$\Gamma(v) = \{w \in V \mid (v,w) \in E\} \cup \{v\}$$

- STRUCTURAL SIMILARITY

$$\sigma(v,w) = \frac{|\Gamma(v) \cap \Gamma(w)|}{\sqrt{|\Gamma(v)| |\Gamma(w)|}}$$



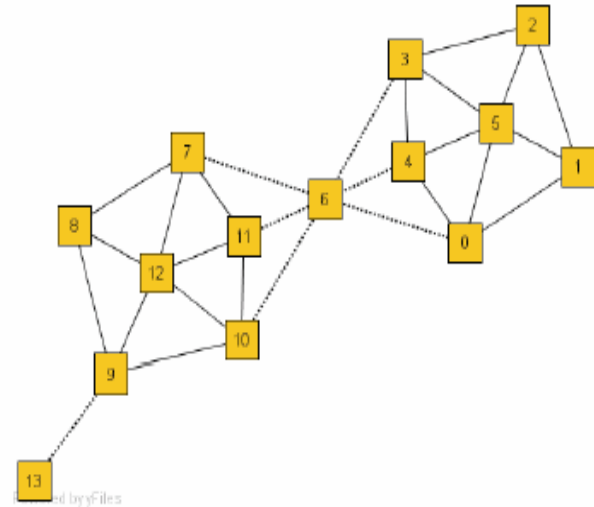
Definitions

- ε -NEIGHBORHOOD

$$N_{\varepsilon}(v) = \{w \in \Gamma(v) \mid \sigma(v, w) \geq \varepsilon\}$$

- CORE

$$CORE_{\varepsilon, \mu}(v) \Leftrightarrow |N_{\varepsilon}(v)| \geq \mu$$



- DIRECT STRUCTURE REACHABILITY

$$DirREACH_{\varepsilon, \mu}(v, w) \Leftrightarrow CORE_{\varepsilon, \mu}(v) \wedge w \in N_{\varepsilon}(v)$$

Definitions

- **HUB**

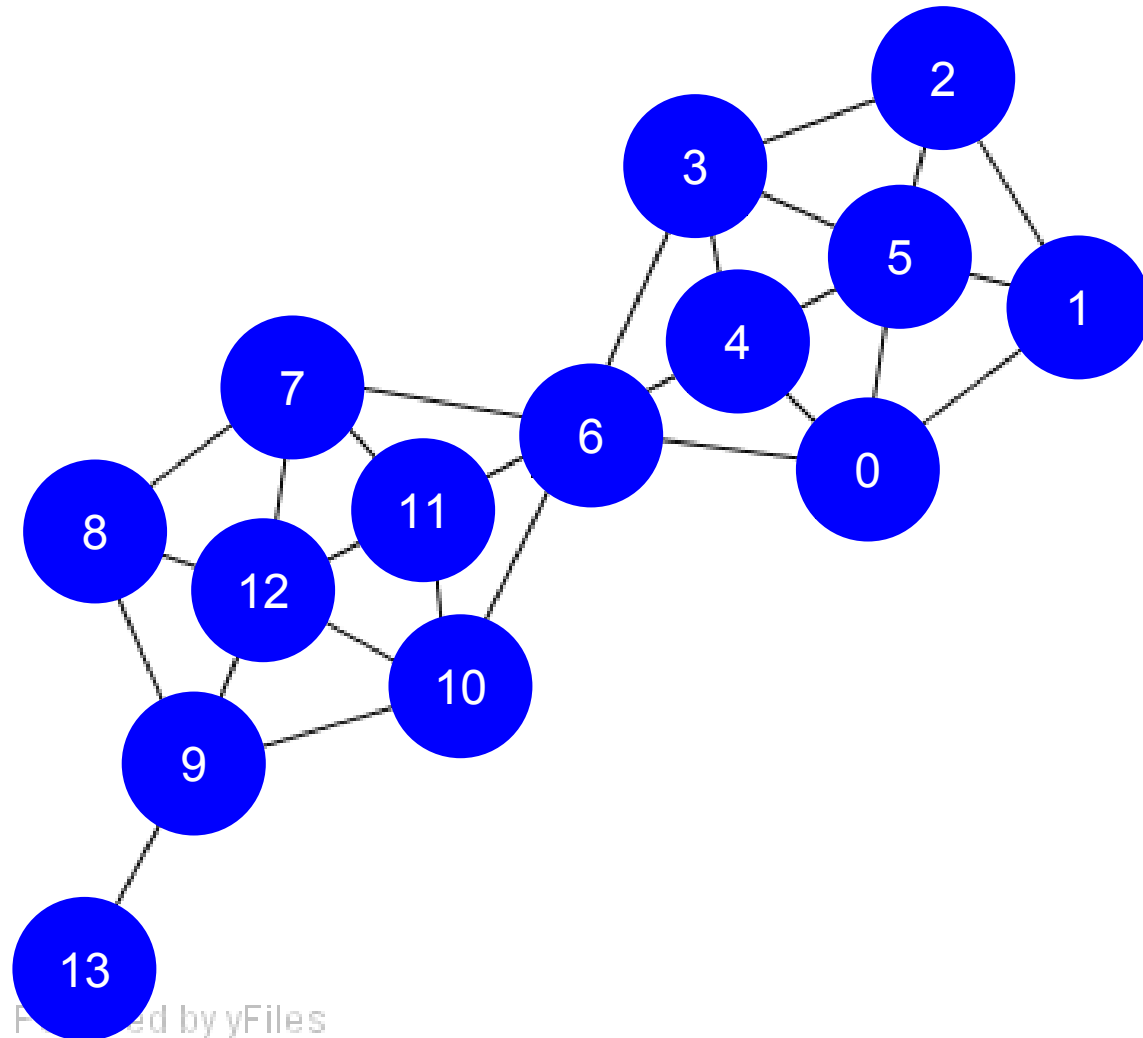
1. If v is not a member of any cluster
2. If v bridges different clusters

- **OUTLIER**

1. If v is not a member of any cluster
2. If v does not bridge different clusters

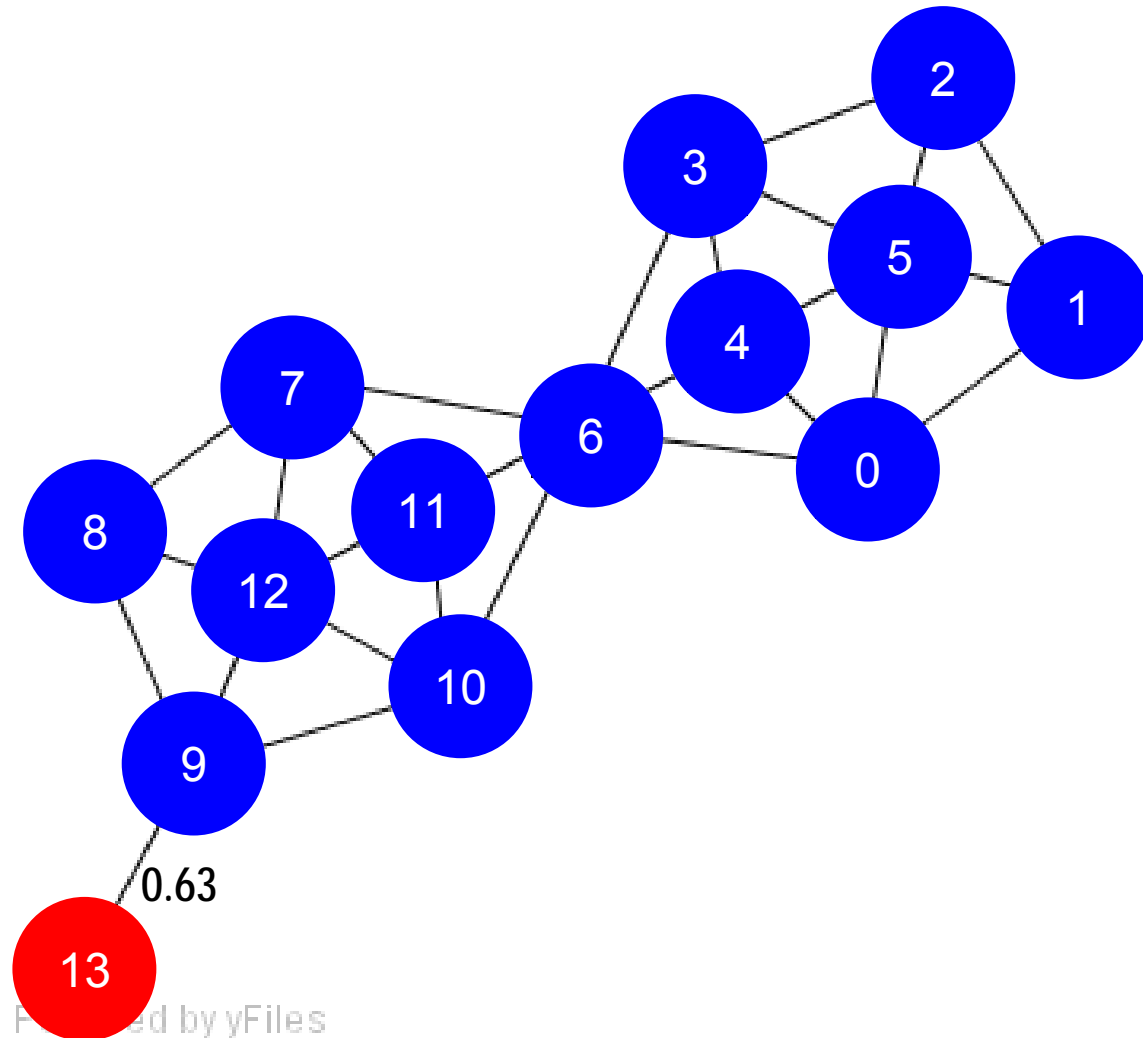
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



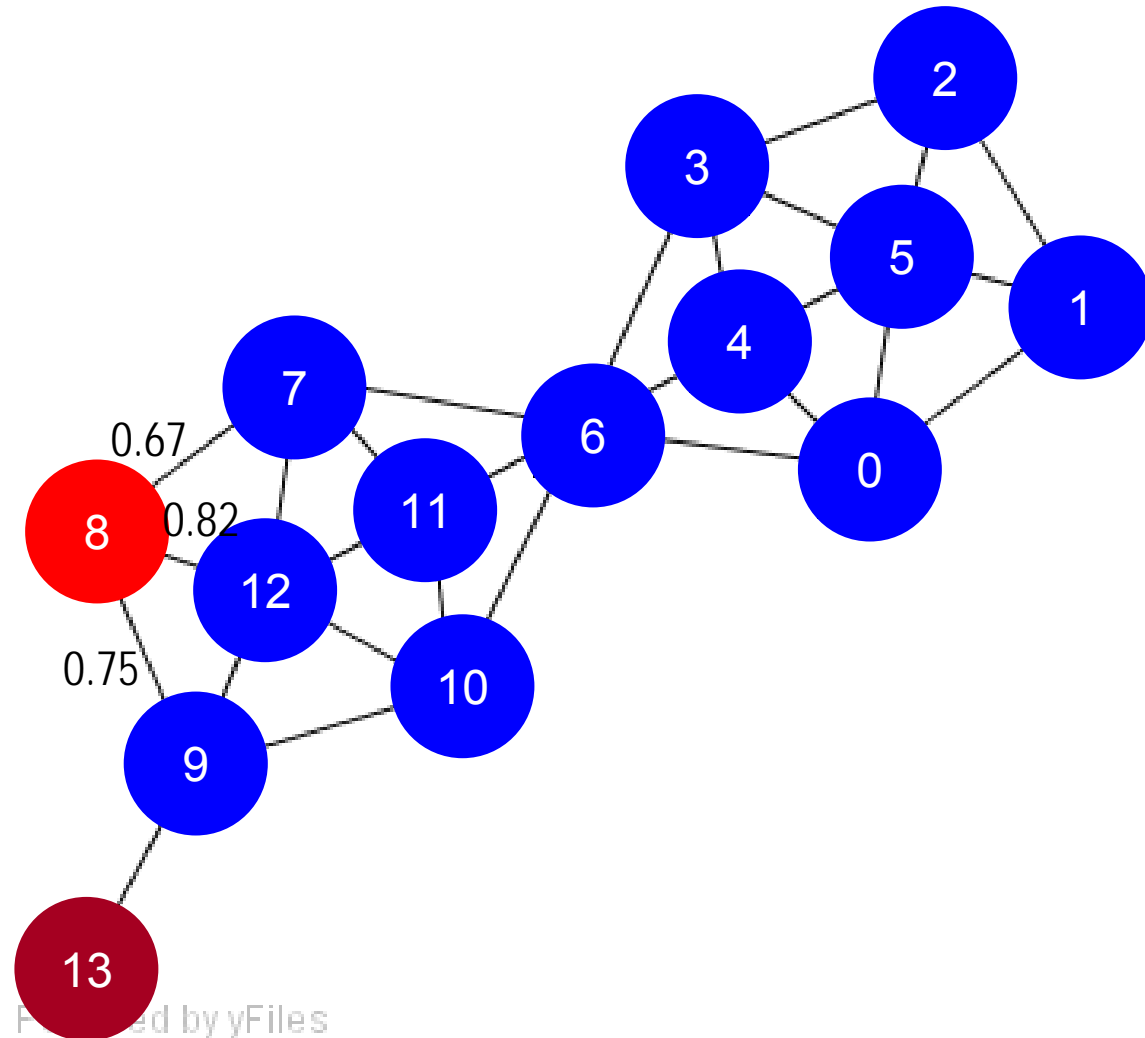
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



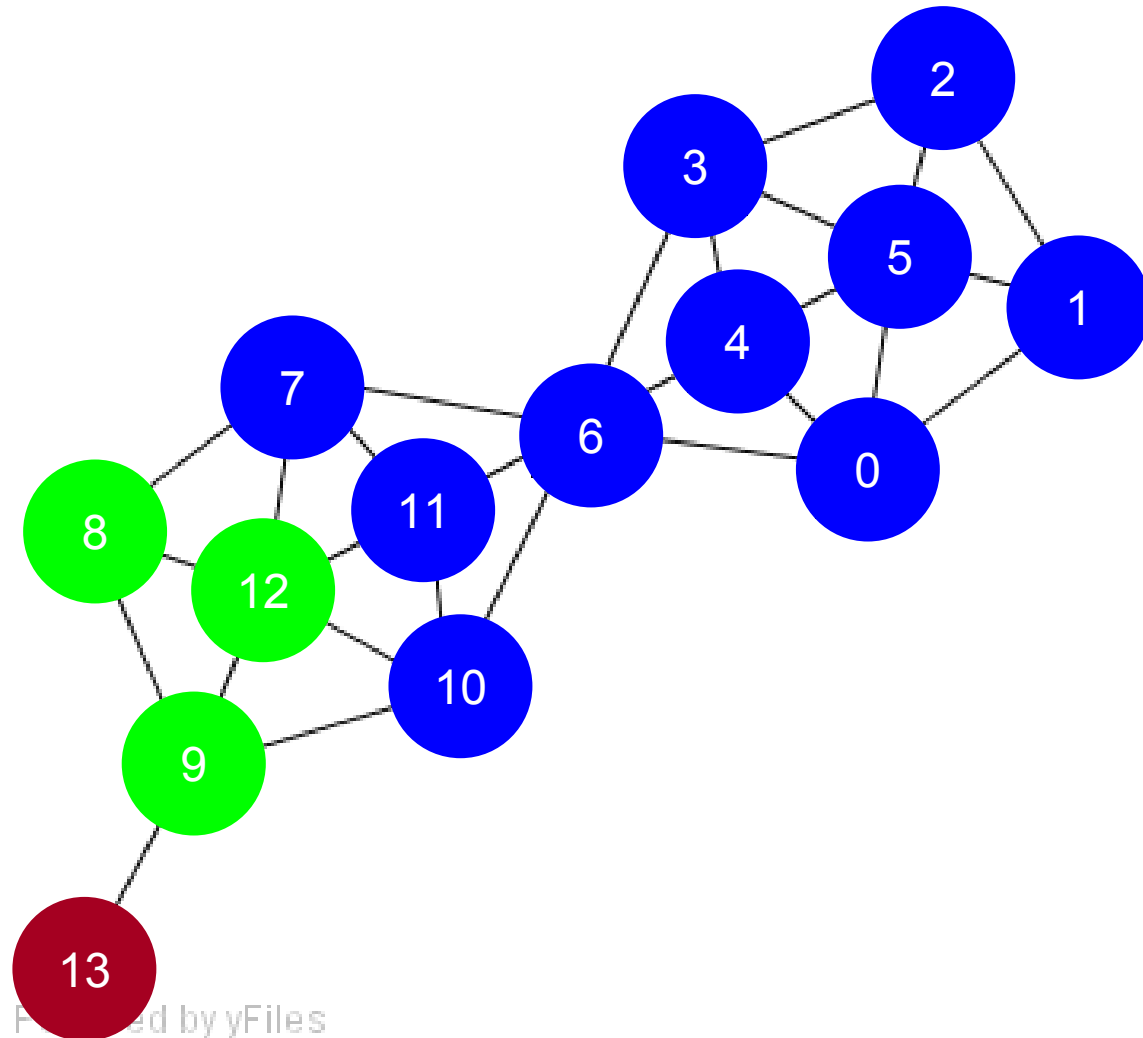
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



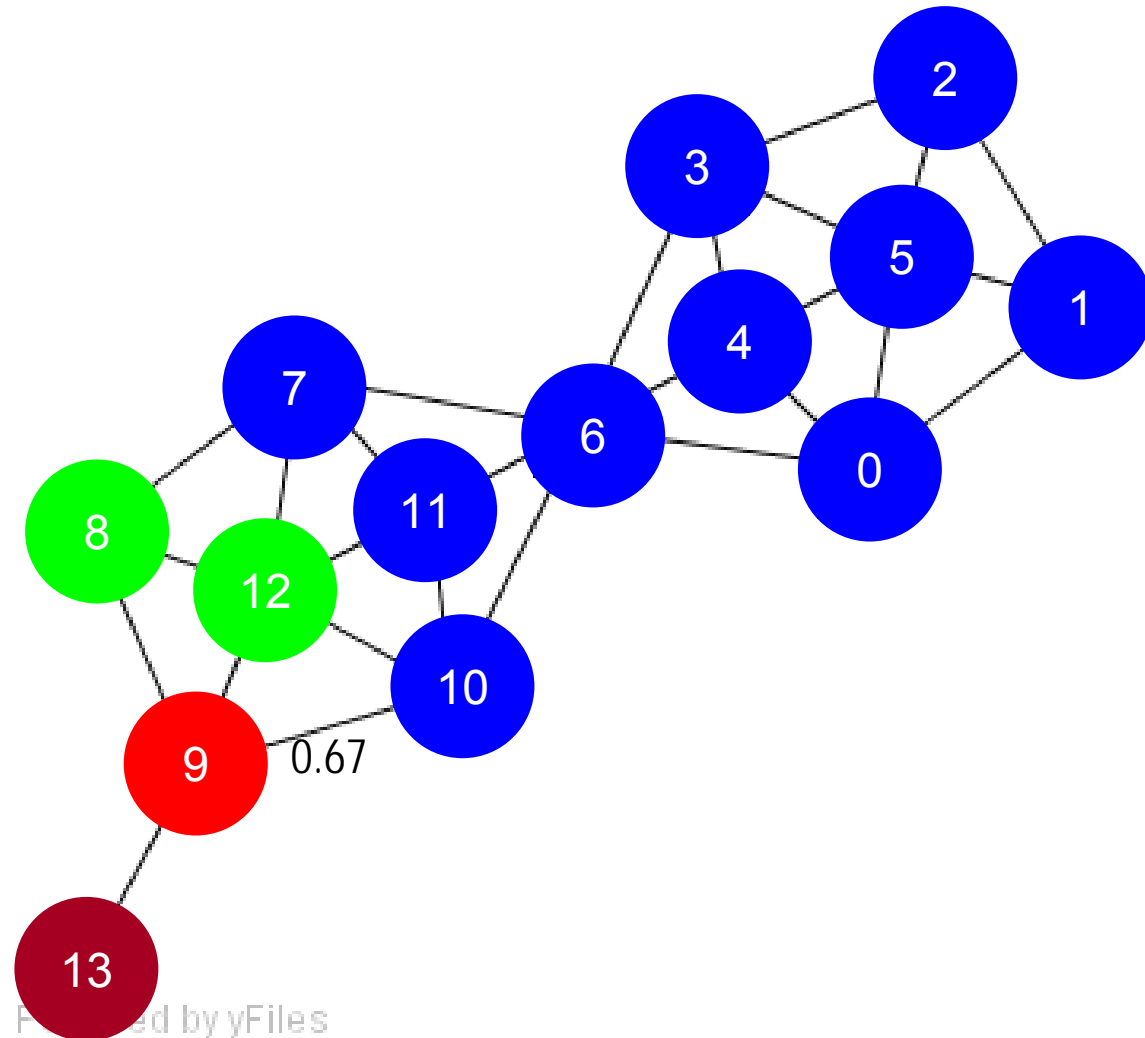
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



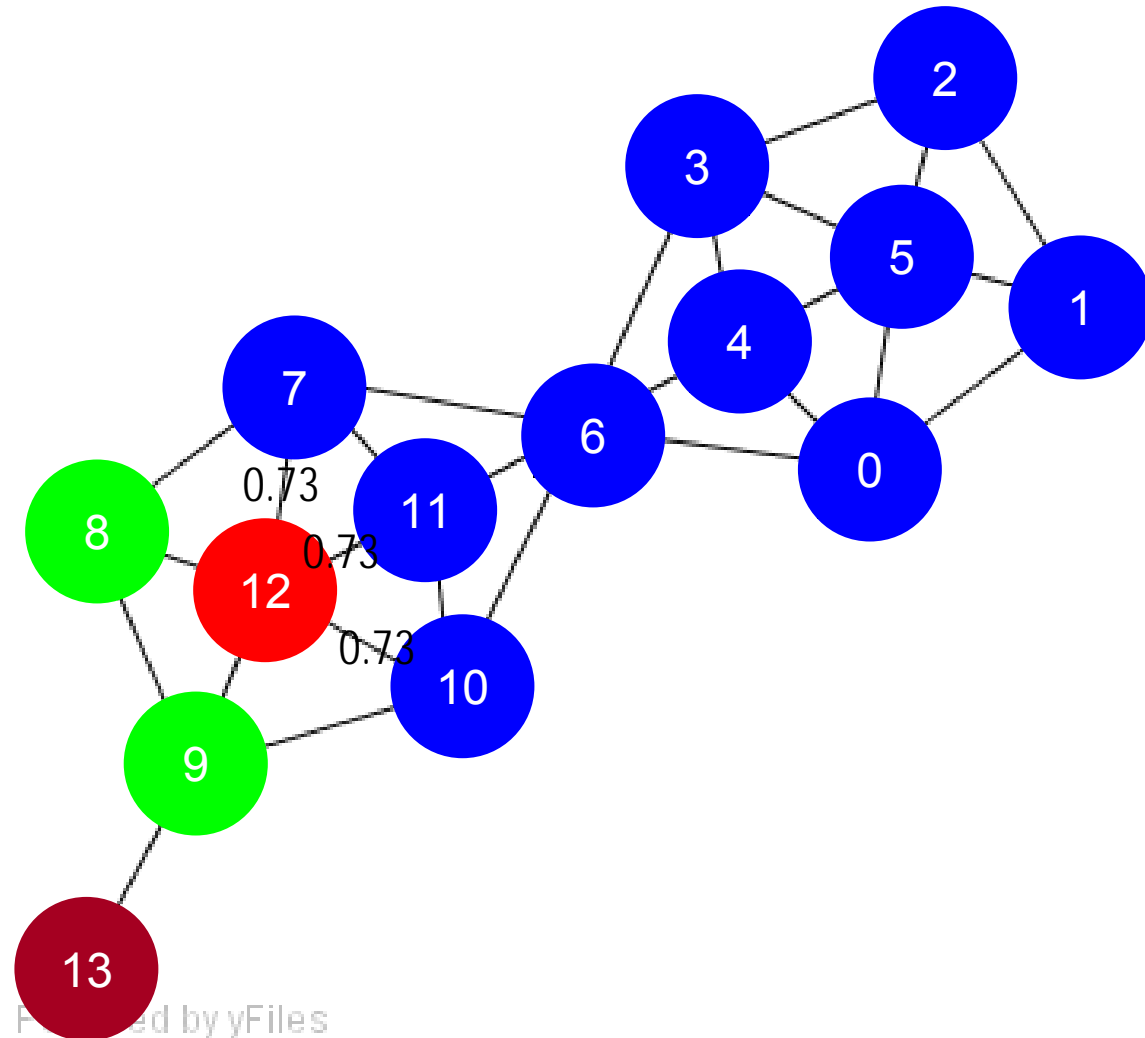
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



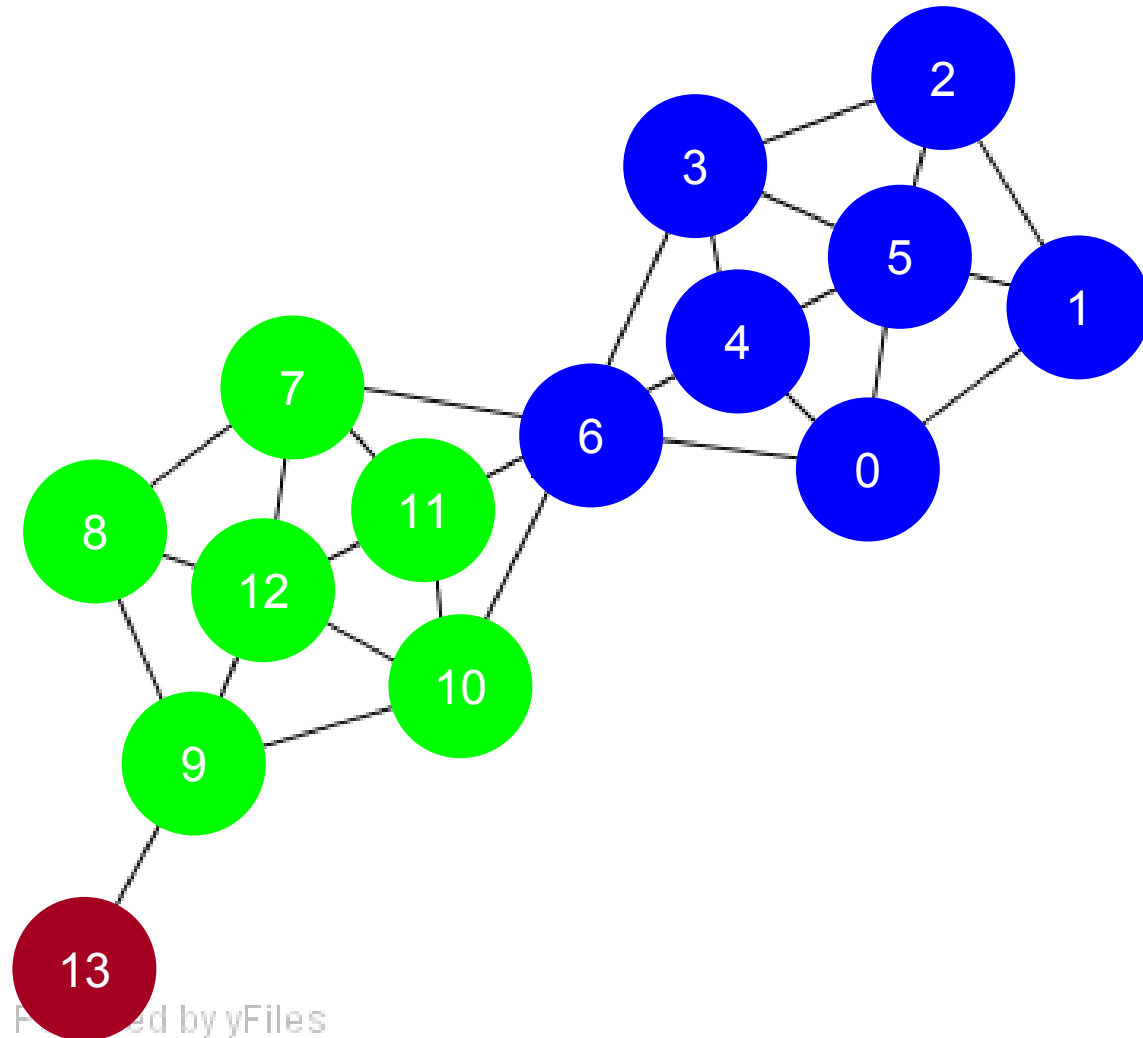
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



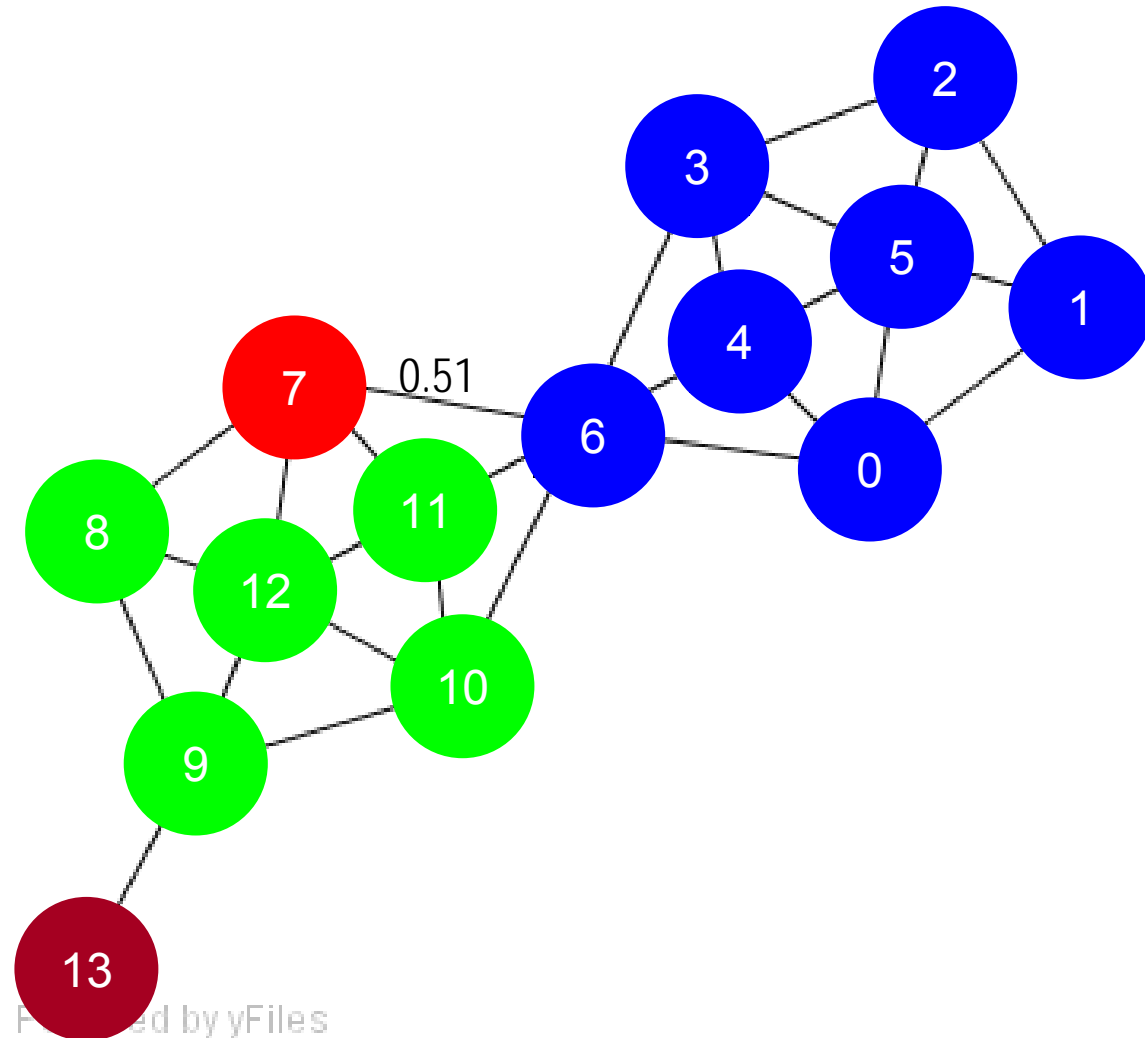
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



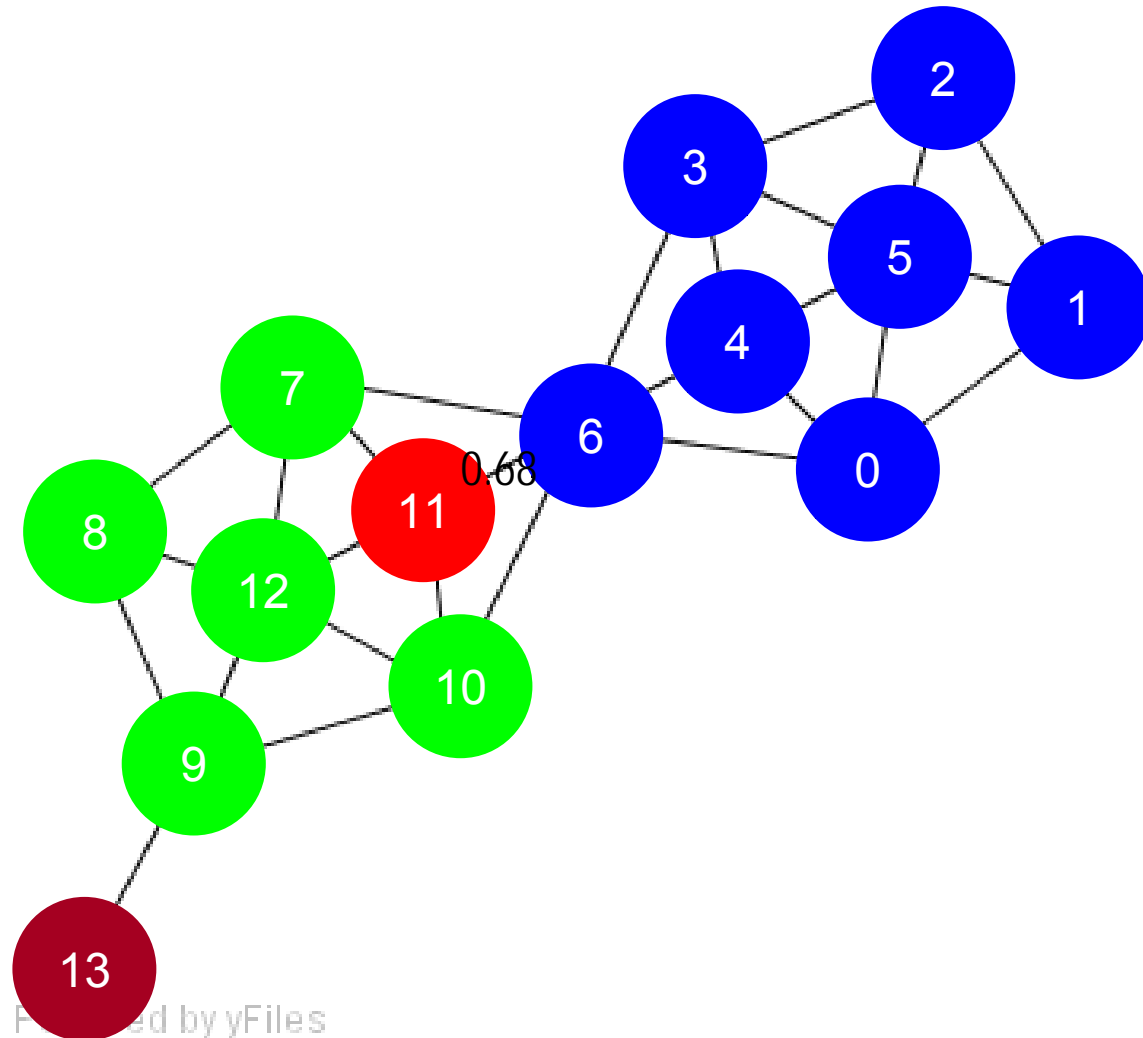
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



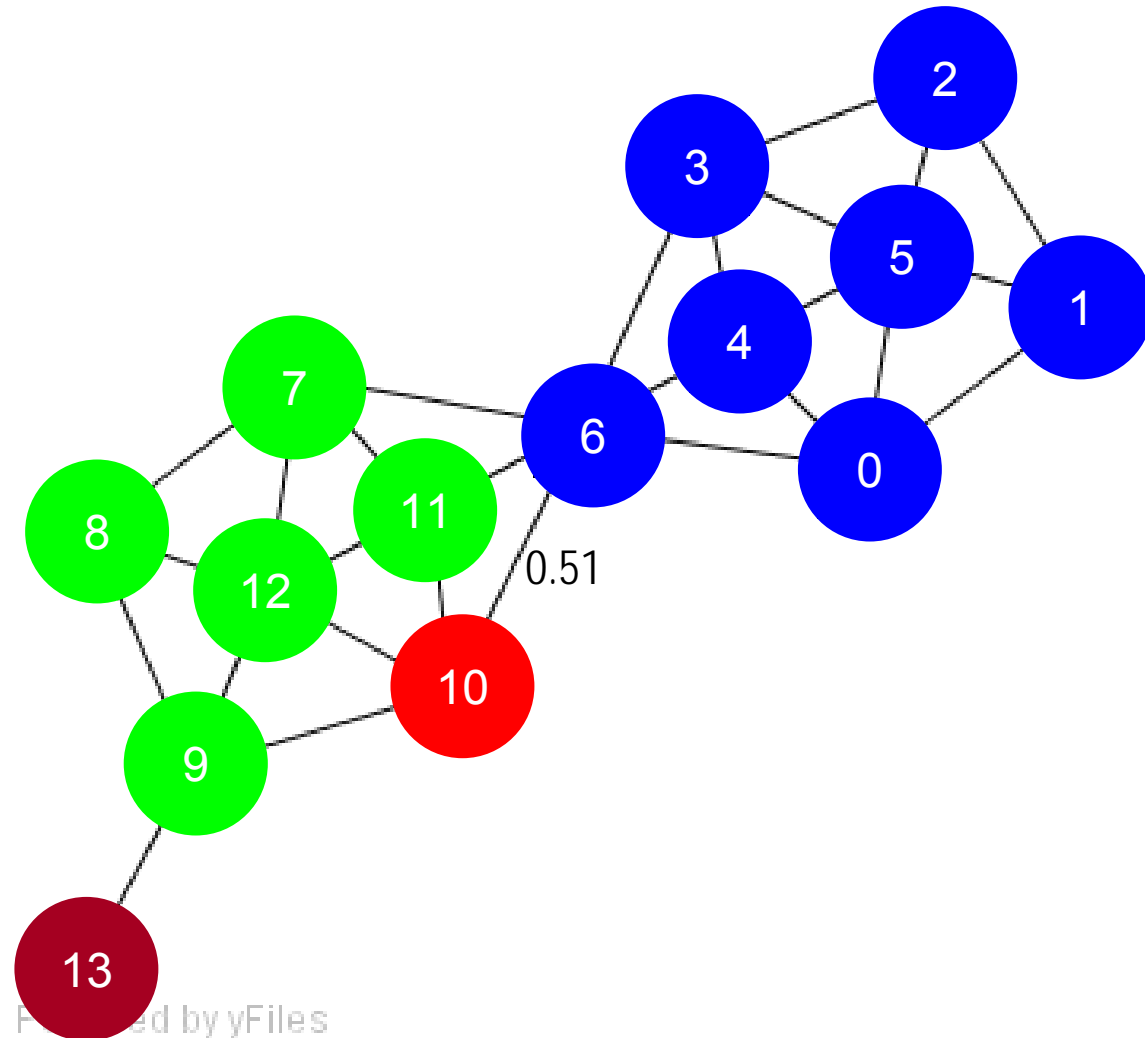
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



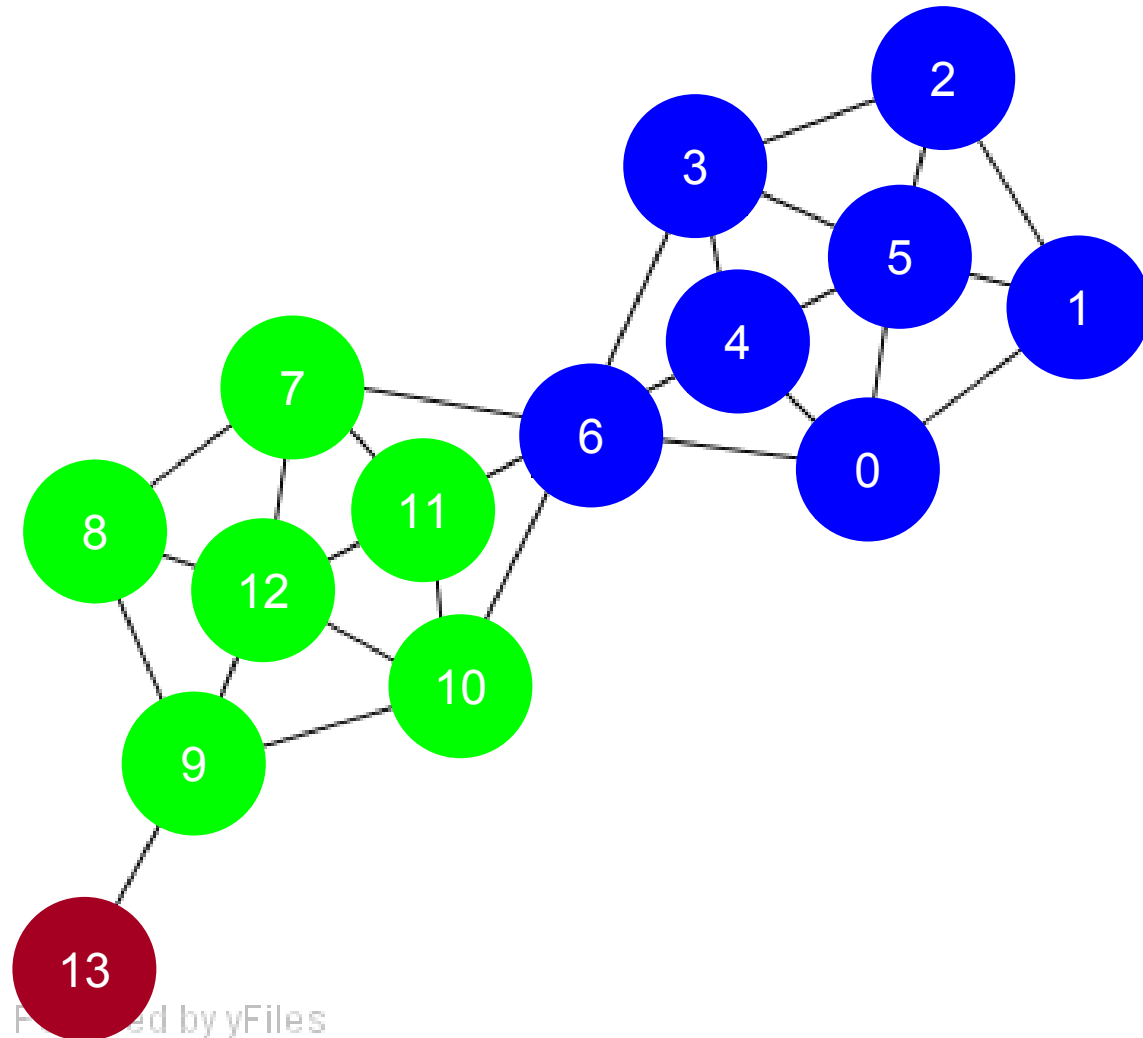
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



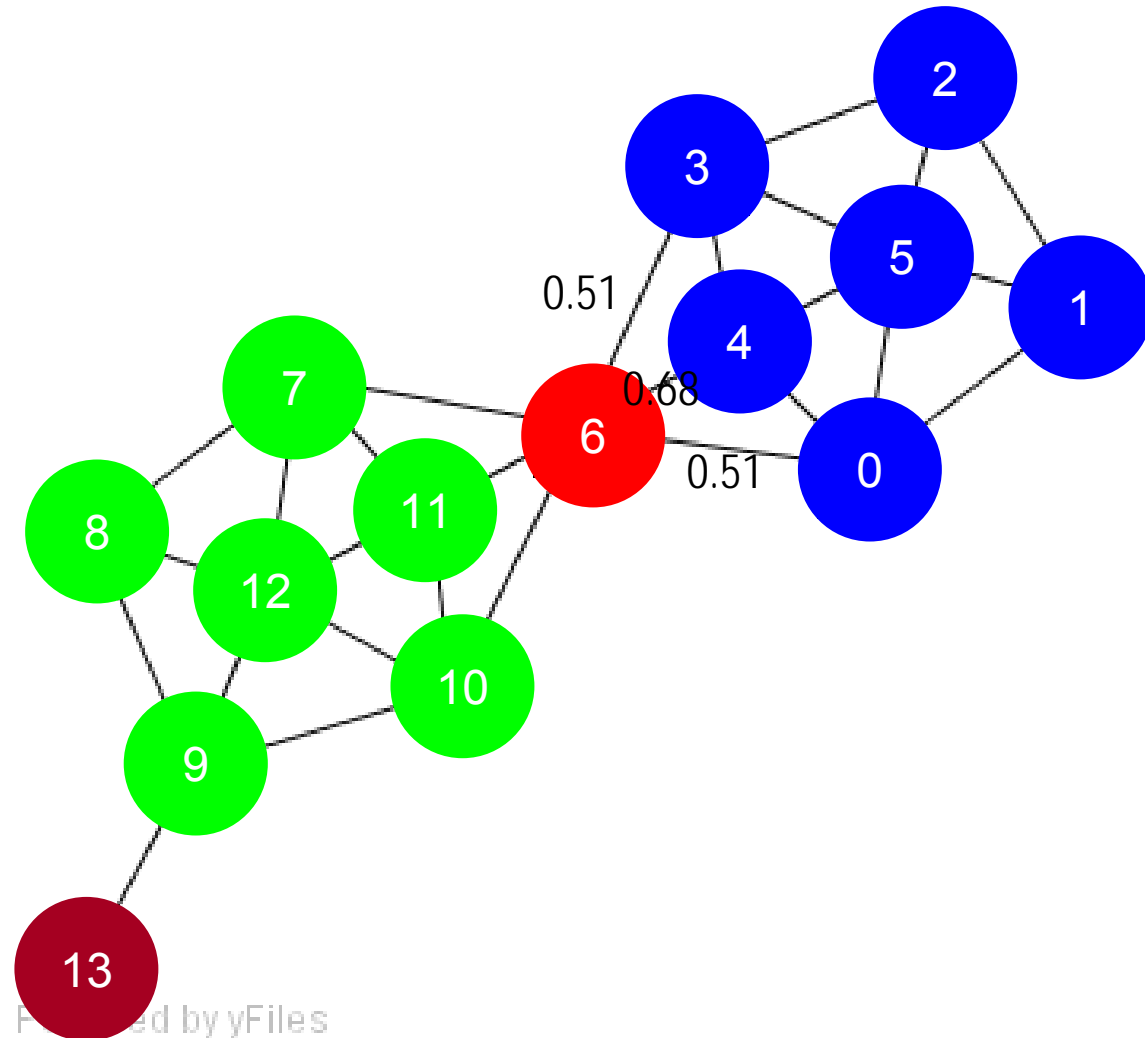
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



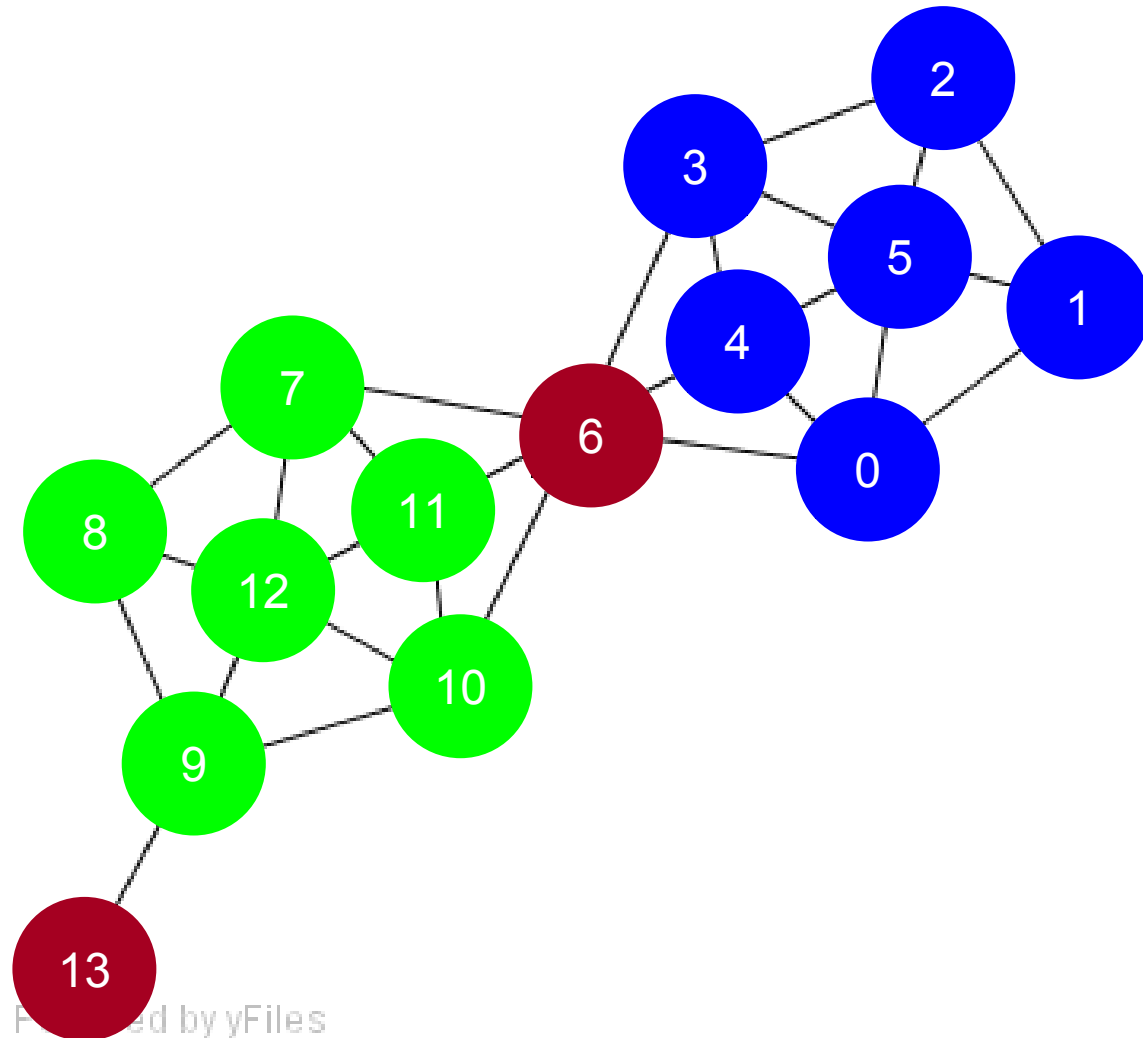
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



```

ALGORITHM SCAN( $G=\langle V, E \rangle, \varepsilon, \mu$ )
// all vertices in  $V$  are labeled as unclassified;
for each unclassified vertex  $v \in V$  do
// STEP 1. check whether  $v$  is a core;
    if  $\text{CORE}_{\varepsilon, \mu}(v)$  then
// STEP 2.1. if  $v$  is a core, a new cluster is expanded;
        generate new clusterID;
        insert all  $x \in N_{\varepsilon}(v)$  into queue  $Q$ ;
        while  $Q \neq \emptyset$  do
             $y = \text{first vertex in } Q$ ;
             $R = \{x \in V \mid \text{DirREACH}_{\varepsilon, \mu}(y, x)\}$ ;
            for each  $x \in R$  do
                if  $x$  is unclassified or non-member then
                    assign current clusterID to  $x$ ;
                if  $x$  is unclassified then
                    insert  $x$  into queue  $Q$ ;
            remove  $y$  from  $Q$ ;
        else
// STEP 2.2. if  $v$  is not a core, it is labeled as non-member
            label  $v$  as non-member;
    end for.
// STEP 3. further classifies non-members
for each non-member vertex  $v$  do
    if ( $\exists x, y \in \Gamma(v) \ (x.\text{clusterID} \neq y.\text{clusterID})$ ) then
        label  $v$  as hub
    else
        label  $v$  as outlier;
    end for.
end SCAN.

```

Complexity Analysis

- Given a graph with m edges and n vertices
- Check the core:
 - Compute structural similarity
 - Examine all the vertex's neighbors
 - $O(\deg(V_1) + \deg(V_2) + \dots + \deg(V_n)) = O(2m)$
 - $O(m)$

Efficiency

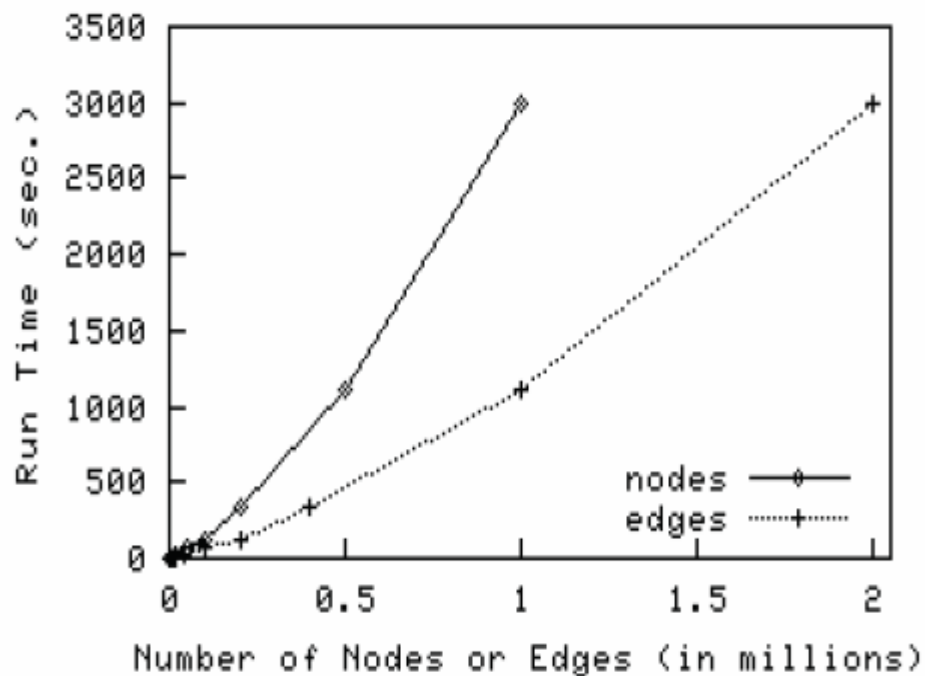


Figure 4. Running Time for FastModularity

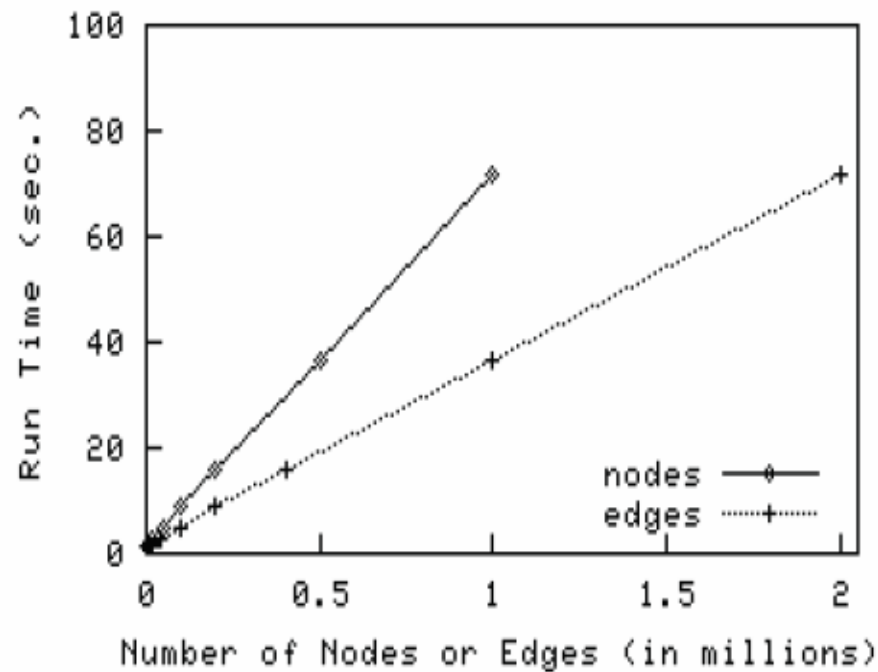


Figure 5. Running Time for SCAN

Effectiveness

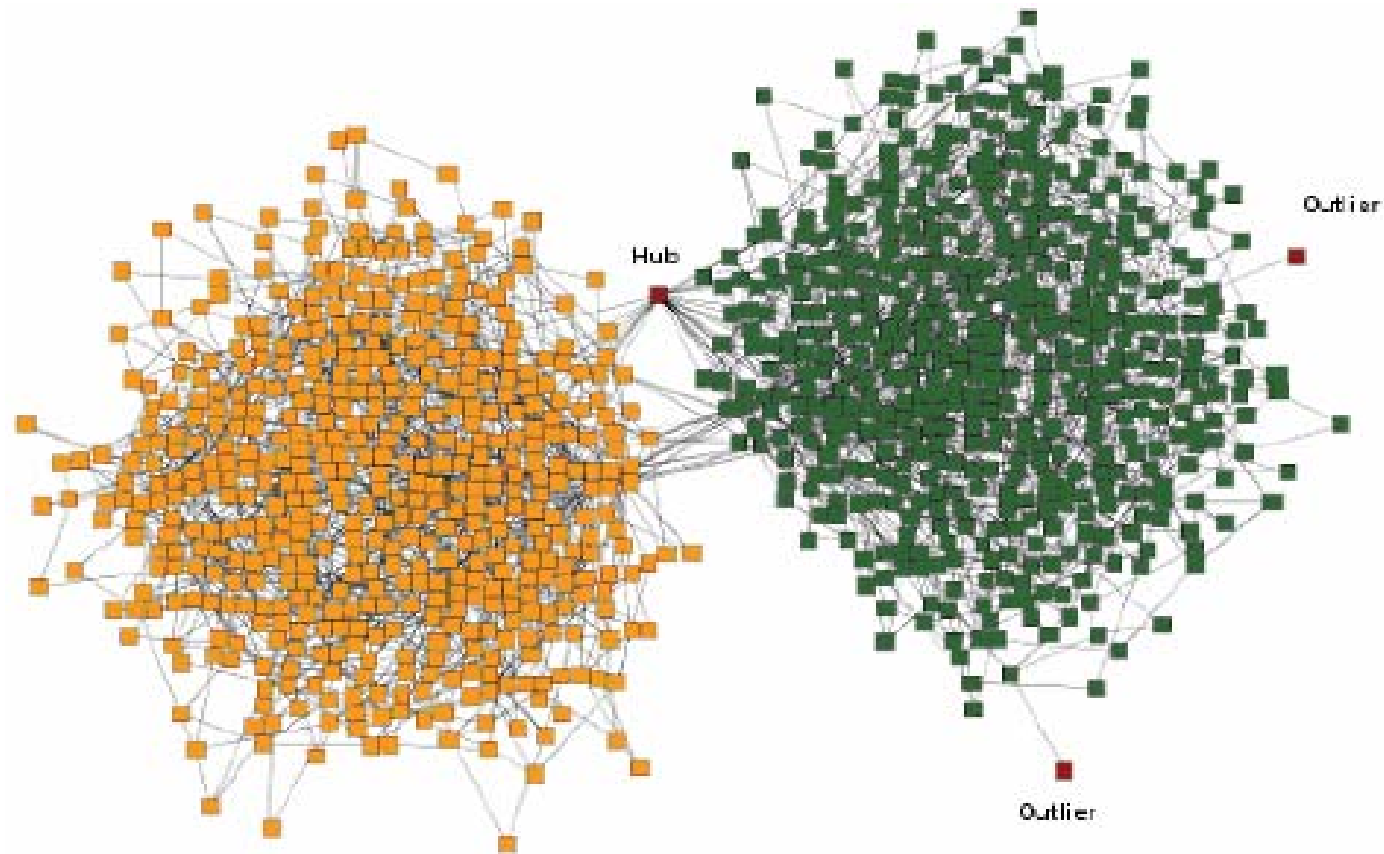
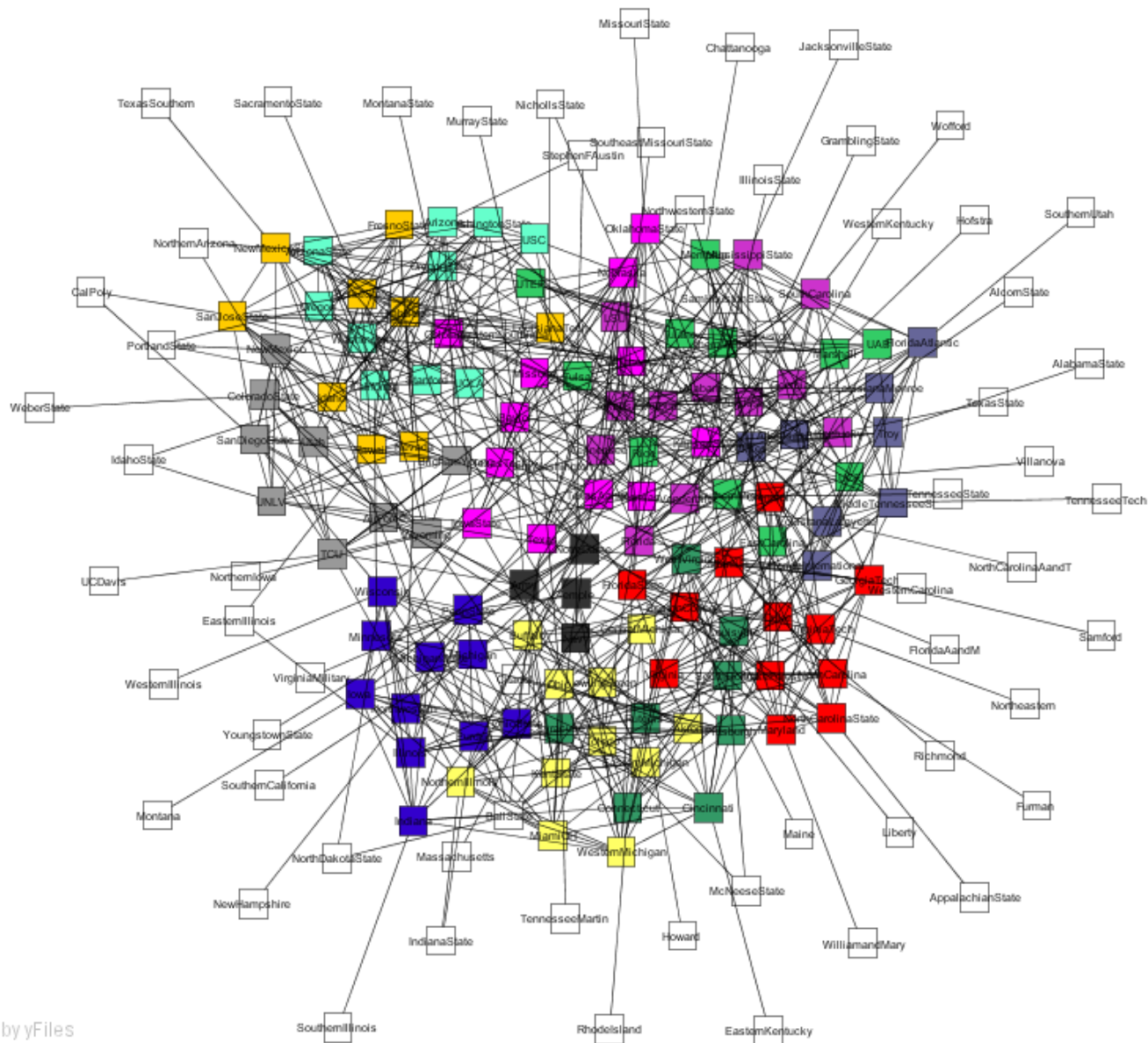
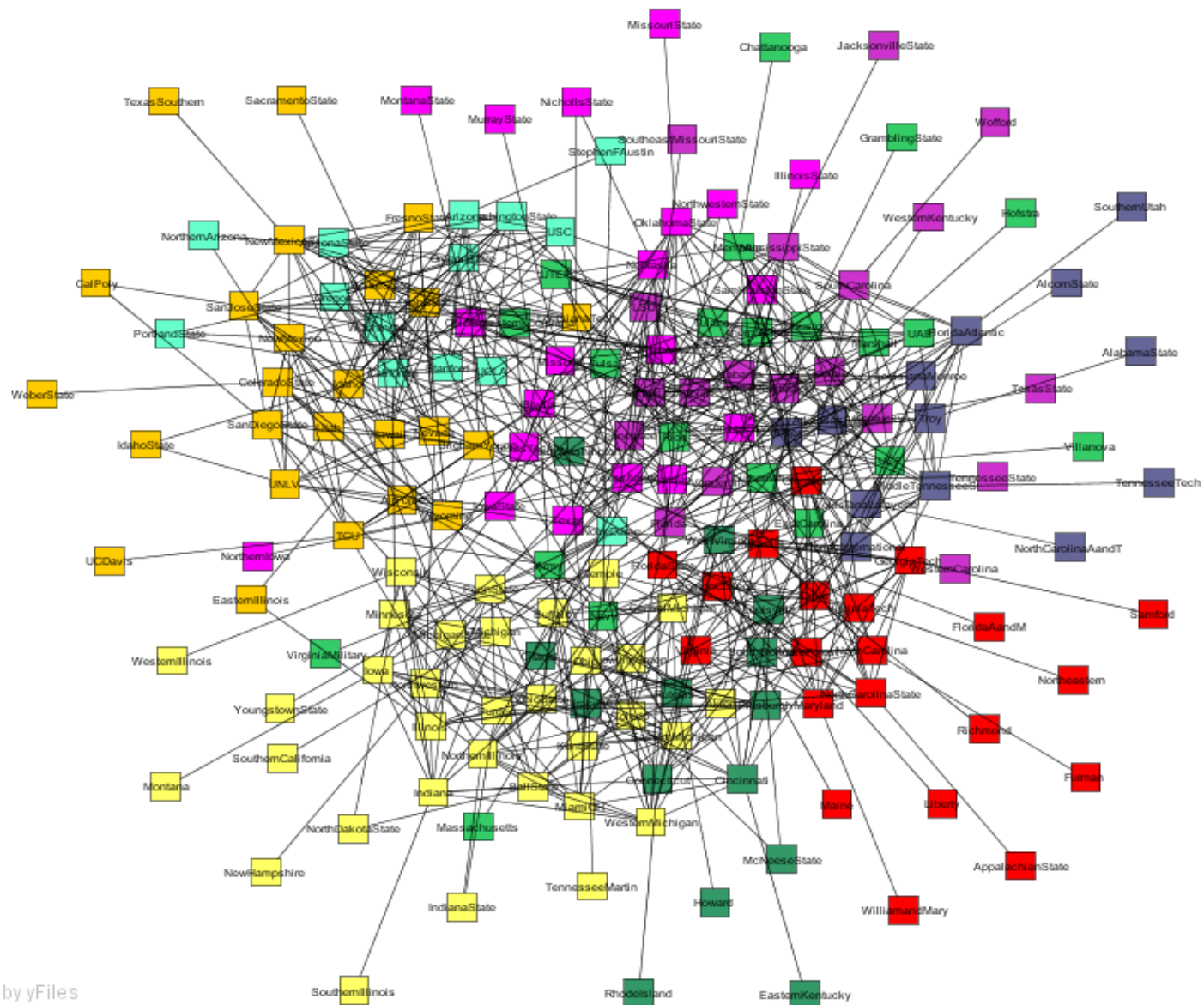


Figure 3. A Synthetic Graph with 1,000 Vertices



Powered by yFiles

Figure 6. NCAA Football Bowl Subdivision schedule as a network, showing the 12 conferences in color, independent schools in black, and lower division schools in white.



Powered by yFiles

Figure 7. NCAA Football Bowl Subdivision schedule as clustered by FastModularity Algorithm.

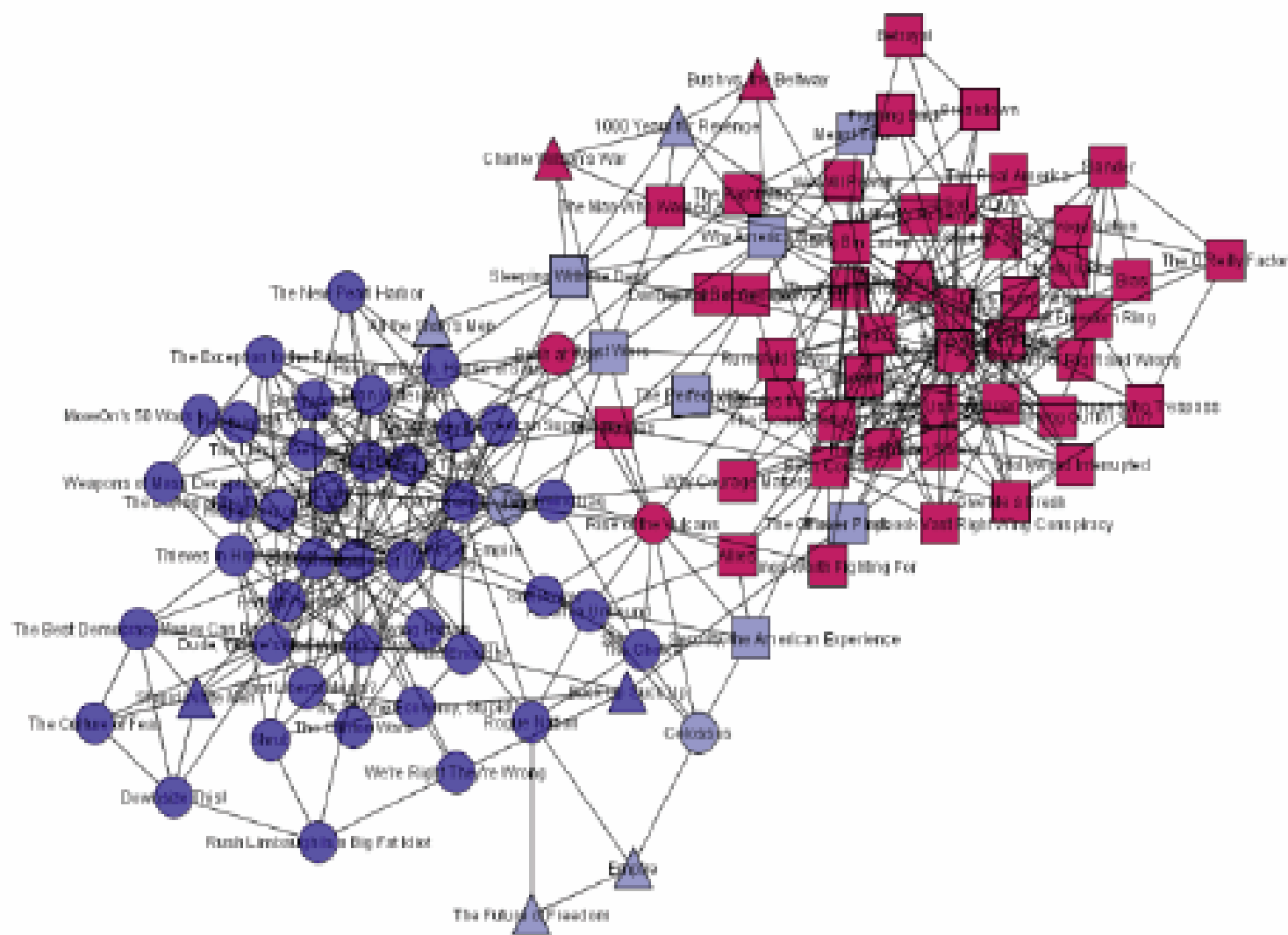


Figure 9. The Result of SCAN on Political Book Graph.

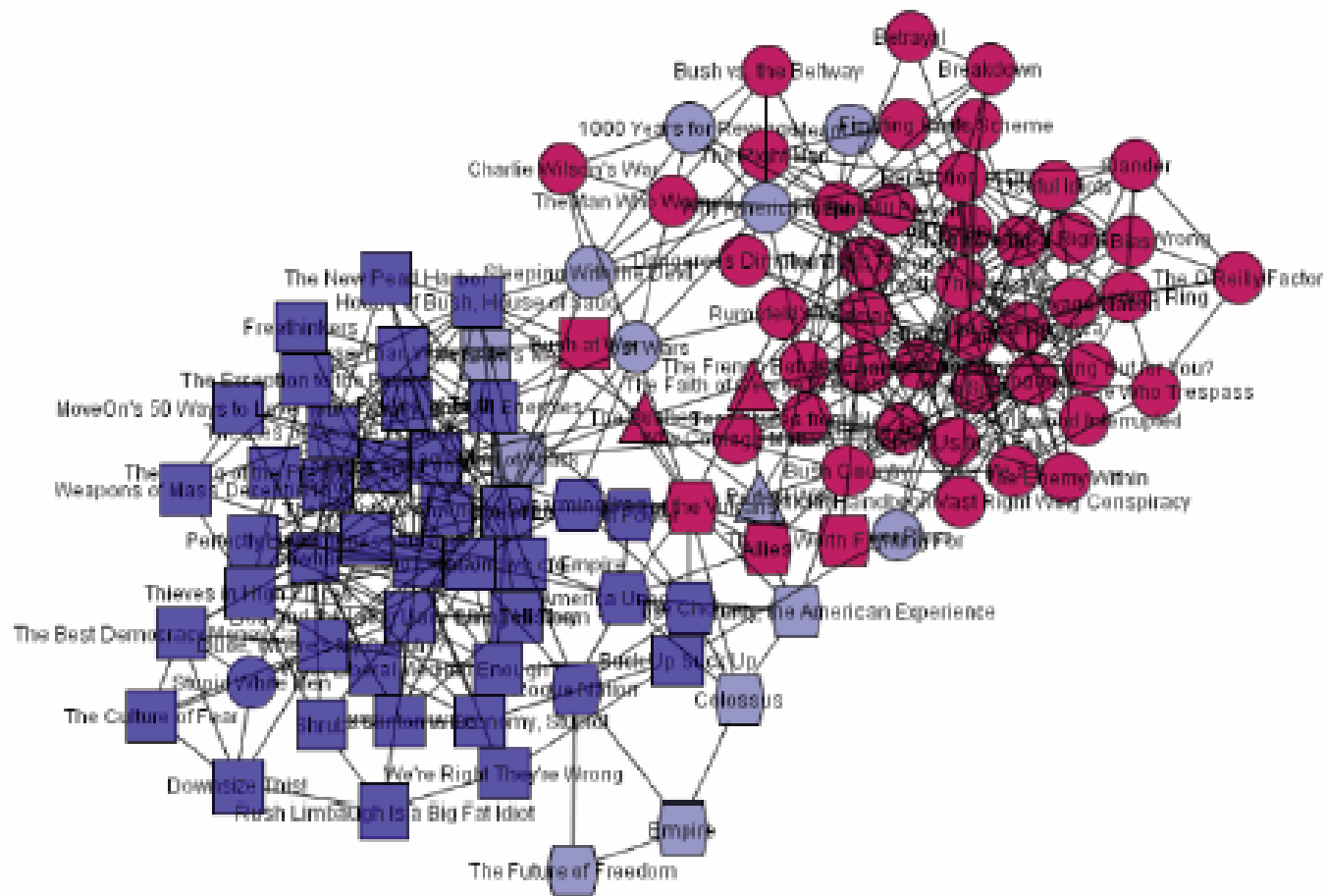


Figure 10. The Result of FastModularity on Political Book Graph.

Choosing ϵ

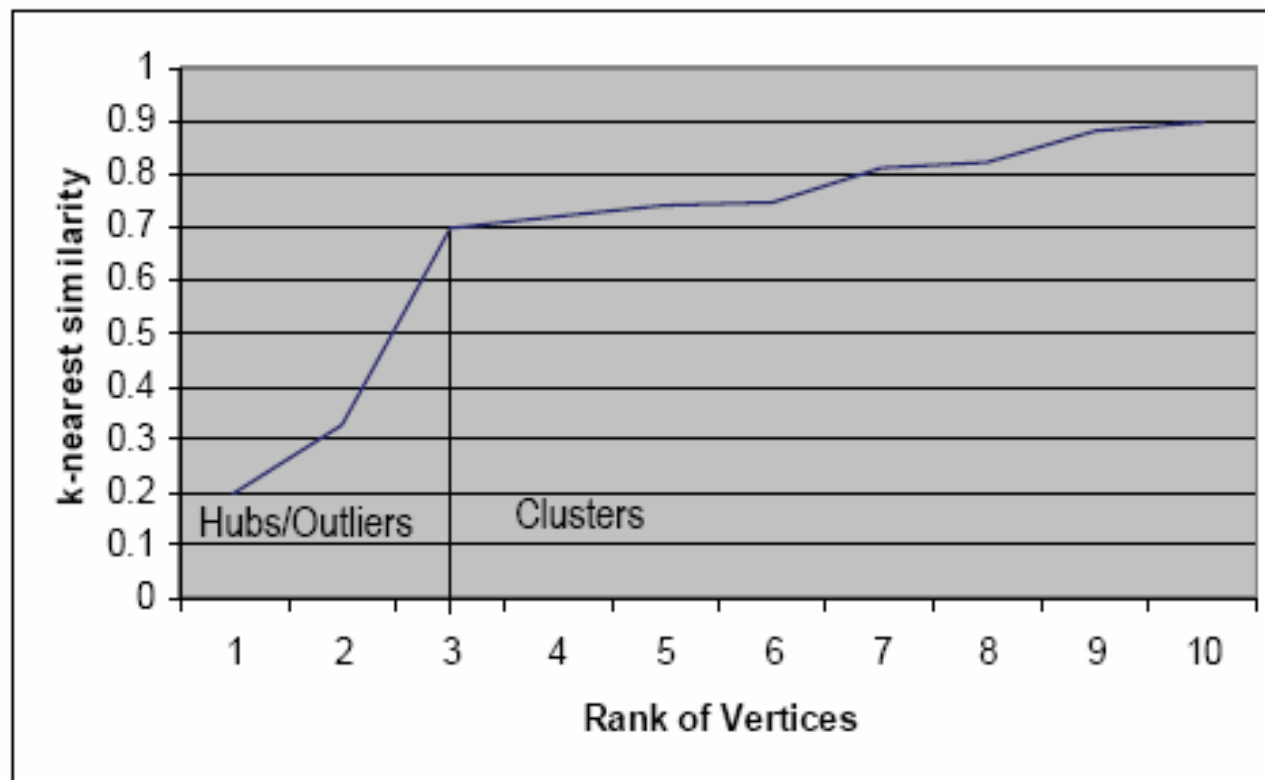


Figure 14. Sorted k-Nearest Structural Similarity.

Adjusted Rand Index (ARI)

$$\frac{\sum_{i,j} \binom{n_{ij}}{2} - \left[\sum_i \binom{n_{i.}}{2} \sum_j \binom{n_{.j}}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[\sum_i \binom{n_{i.}}{2} + \sum_j \binom{n_{.j}}{2} \right] - \left[\sum_i \binom{n_{i.}}{2} \sum_j \binom{n_{.j}}{2} \right] / \binom{n}{2}}$$

	SCAN	FastModularity
College football	1	0.24
Political books	0.71	0.64