

# Community Detection in Social Networks

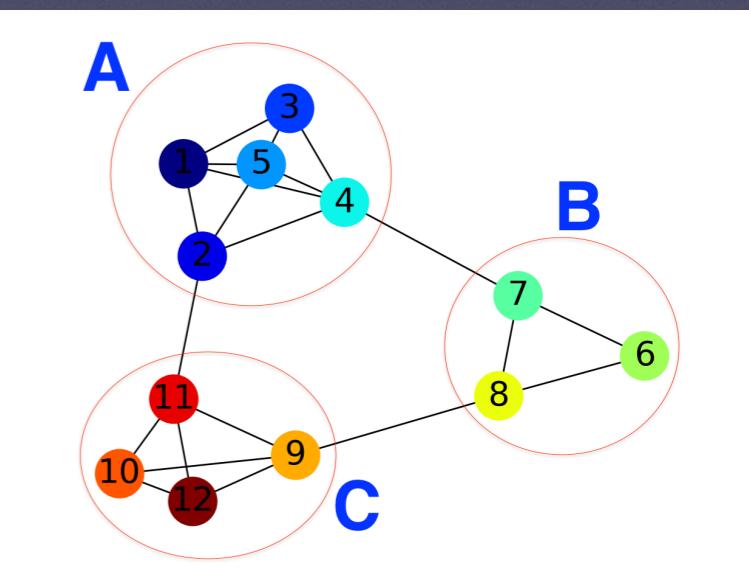
Yawen Sun  
1836786

# Outline

- Background
- Motivation and Objectives
- 4 types of Algorithms
- Results and Evaluation

# Background

- Communities in Social Networks
- Features of communities
  - Nodes inside are communicated frequently
  - Nodes outside are connected loosely
- A simple network with 3 communities



# Motivations and Objectives

Community detection in distinct fields

- Sociology
  - Helps in users behaviours analysis
- Biology
  - Helps in identification of a new population
- Computer Science
  - Can be used in recommender systems

# Community Detection Algorithms

- GN Algorithm
  - A divisive algorithm based on graph theory
- Louvain
  - An agglomerative algorithm using a greedy techniques
- Spectral Clustering
  - An algorithm based on clustering techniques
- GA-Net
  - An algorithm based on genetic algorithm inspired from nature

# GN Algorithm

Input: the network  $G=(V,E)$

while there is no edge:

for  $e \in E$  do

calculating its edge betweenness centrality

remove an edge with the highest value of edge betweenness centrality

if there is a tie, remove an edge in the tie set

calculating the modularity

Output: the partition with the highest modularity

**Betweenness centrality:**

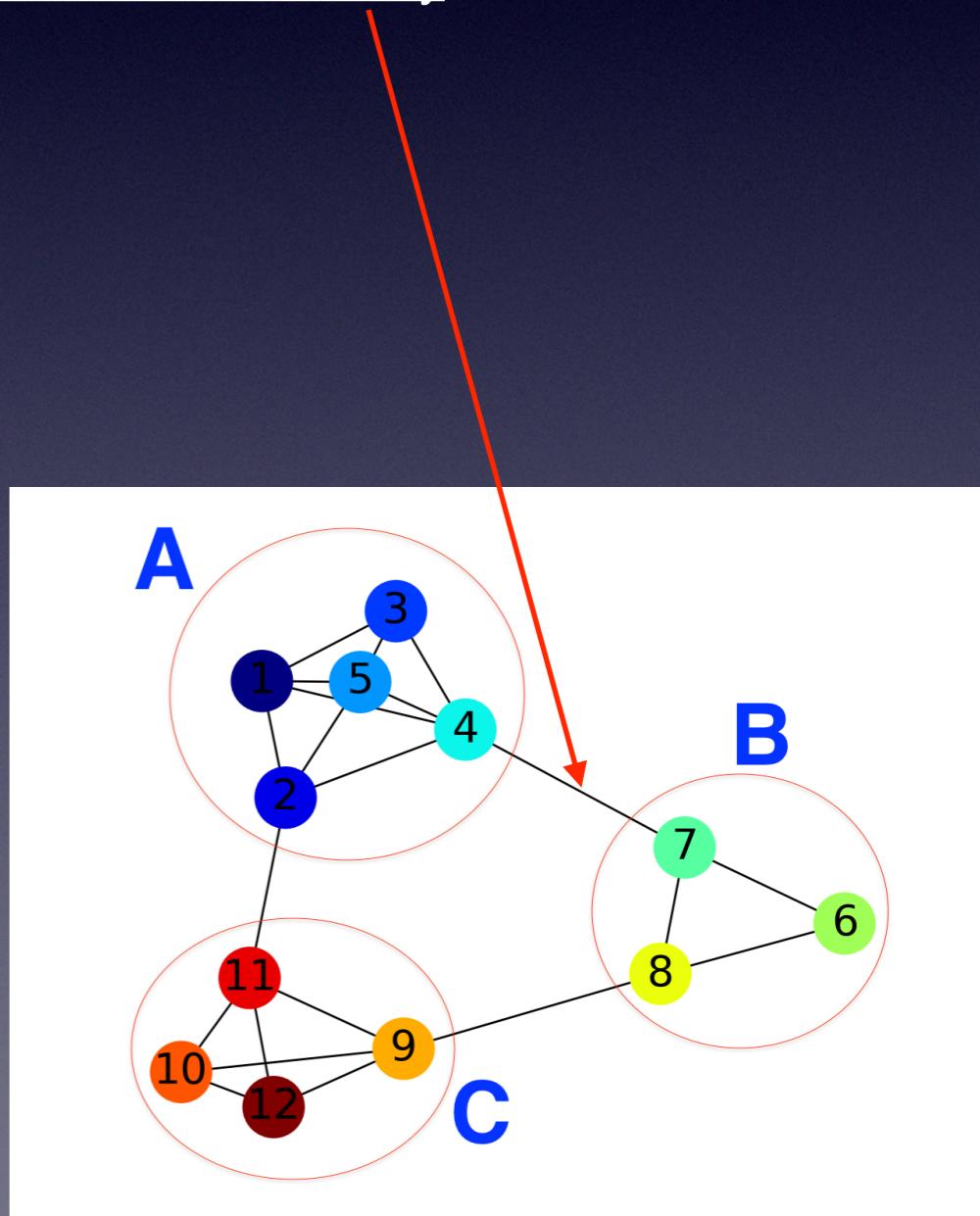
$$C_B(e) = \sum_{s \neq v \neq t \in V} \frac{\sigma_{st}(e)}{\sigma_{st}}$$

A fraction of the shortest paths between two nodes.

**Modularity:**

$$Q = \sum_i^n (e_{ii} - a_i^2), \quad a_i = \sum_j^n e_{ij}$$

A measure of how dense in communities.



# Louvain

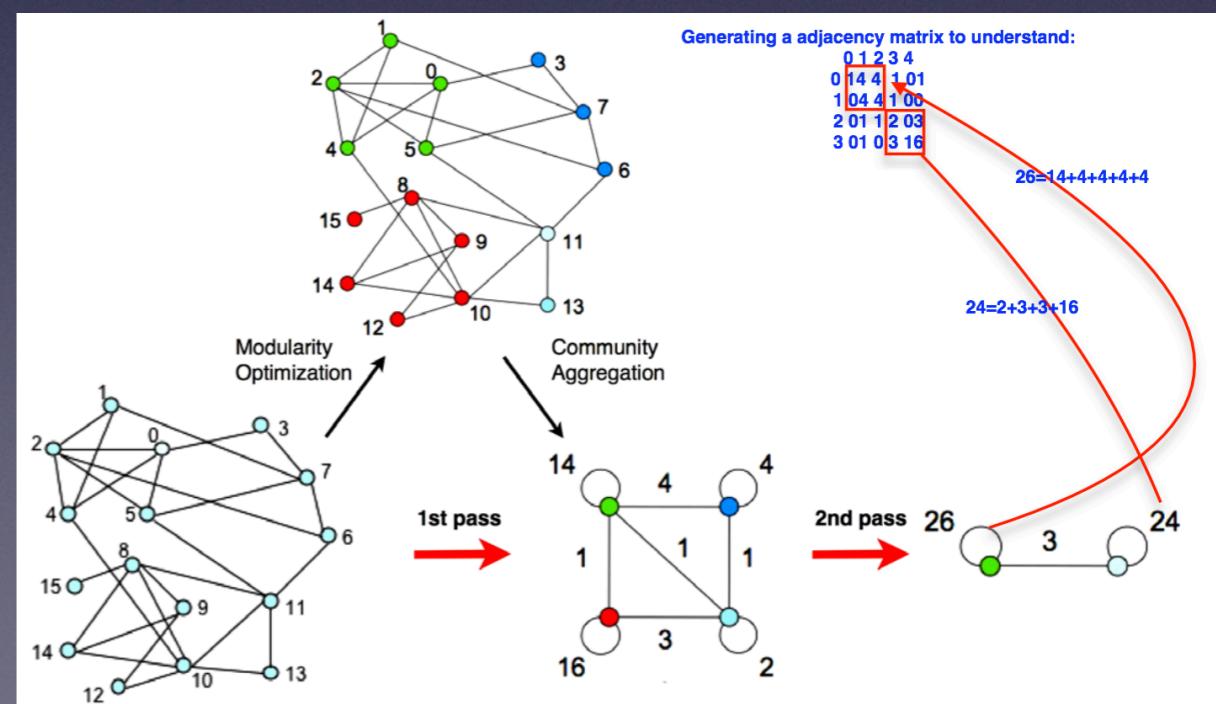
1. Assign every node as a community.
2. Combine:
  - a. combine every node with its neighbour community and calculate  $\Delta Q$
  - b. put this node into its neighbour whose value of  $\Delta Q$  is the highest
3. Set new nodes:
  - a. set one community as a new “node”
  - b. change the weights of links
4. Repeat 1. until Modularity is not changed.

## $\Delta Q$ :

The difference between Modularity of before-combination community and after-combination community

$$\Delta Q = \left[ \frac{\Sigma_{in} + 2k_{i,in}}{2m} - \left( \frac{\Sigma_{tot} + 2k_i}{2m} \right)^2 \right] - \left[ \frac{\Sigma_{in}}{2m} - \left( \frac{\Sigma_{tot}}{2m} \right)^2 - \left( \frac{k_i}{2m} \right)^2 \right] = \frac{1}{2m} \left( k_{i,in} - \frac{\Sigma_{tot} k_i}{m} \right)$$

$$\text{maximising } \Delta Q = \text{maximising } \left( k_{i,in} - \frac{\Sigma_{tot} k_i}{m} \right)$$



# Spectral Clustering

## Unnormalised Spectral Clustering

Input: **the value of k**, the network

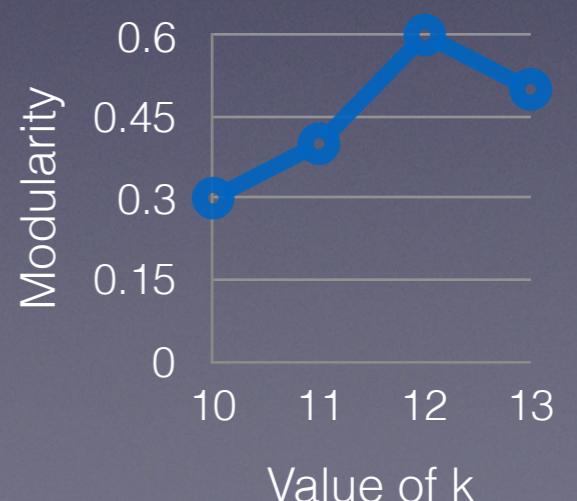
1. construct the unnormalised Laplacian matrix, L from the network
2. calculate the first k eigenvectors  $u_1, \dots, u_k$  from Laplacian matrix
3. generate matrix U (a  $N \times k$  matrix) whose columns are k eigenvectors
4. apply **k-means** for N points in U with k dimensions to get k clusters

Output: k clusters

## Laplacian Matrix L:

$$L = D(\text{degree matrix}) - A(\text{adjacency matrix})$$

## To find the optimal k



# Spectral Clustering

## Normalised Spectral Clustering

Input: **the value of k**, the network

1. calculate the normalised Laplacian matrix, L<sub>sym</sub>
2. computer first k eigenvectors u<sub>1</sub>,...,u<sub>k</sub> from normalised Laplacian matrix, L<sub>sym</sub>
3. construct matrix U (a N\*k matrix) whose columns are k eigenvectors from last step
4. generate matrix T from U through normalising, to be specific, tip
5. apply **k-means** for N points in T with k dimensions to get k clusters

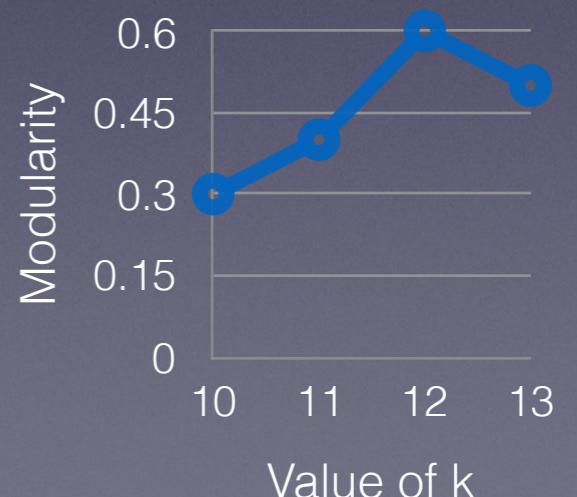
Output: k clusters

## Normalised Laplacian Matrix L:

$$L_{sym} = D^{-\frac{1}{2}}LD^{-\frac{1}{2}} = D^{-\frac{1}{2}}(D - A)D^{-\frac{1}{2}}$$

Matrix T:  $t_{ij} = \frac{u_{ij}}{(\sum_k u_{ik}^2)^{\frac{1}{2}}}$

To find the optimal k



# GA-Net

## 1. Initialisation:

initialise population in the first generation “*the locus-based adjacency representation*” (*left part*)

## 2. compute fitness value for individuals

## 3. Selection:

select two individuals through **roulette selection method**

## 4. Uniform Crossover:

according to a generated binary string to crossover (*right part*)

## 5. Mutation:

change value of the gene randomly if the changed one is appropriate

## 6. replace the old generation

## 7. repeat from step 2 to step 6 until the terminal condition is satisfied

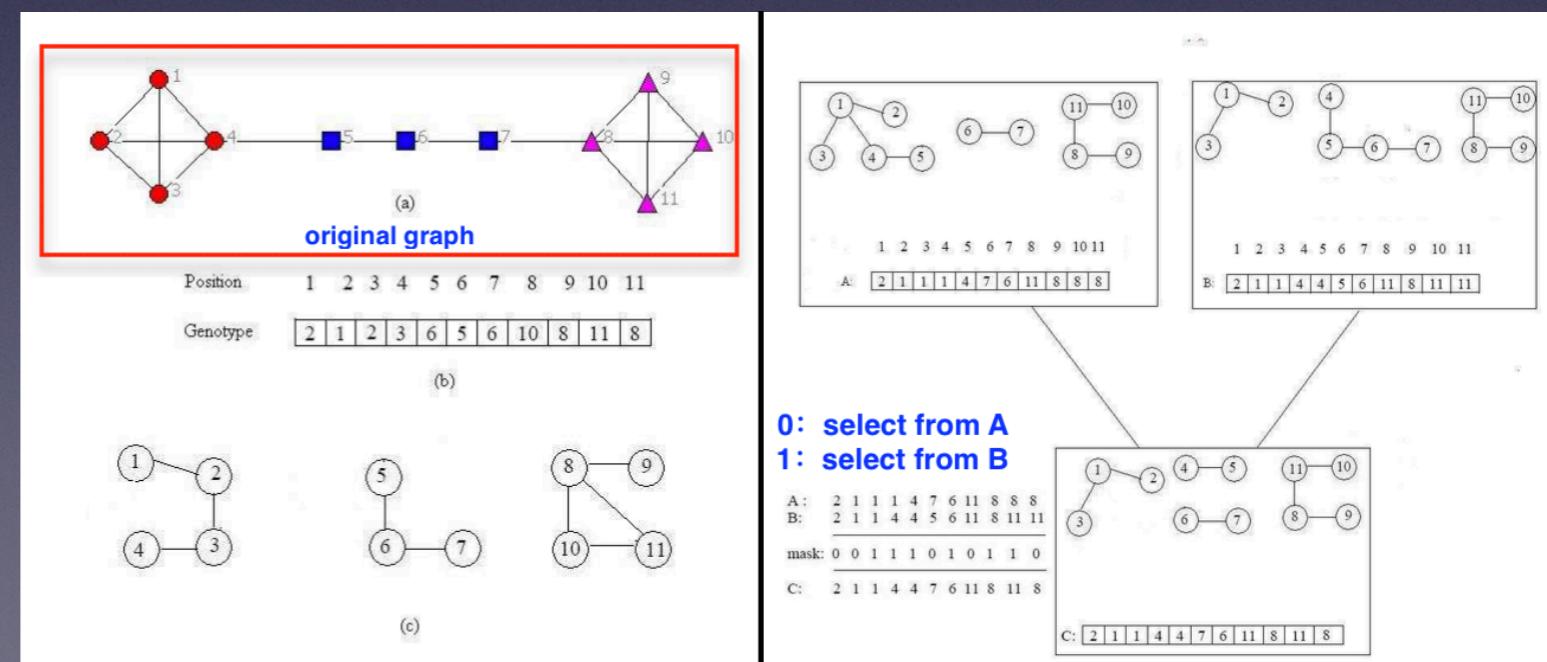
## Objective Function:

Community Score (CS)

$$CS(r) = \sum_i^k Q(S_i, r) = \sum_i^k M(S_i, r) * V_s \quad (10)$$

$$M(S, r) = \frac{\Sigma(a_{ij})^r}{|I|} \quad (10.1), \quad V_s = \sum_{i \in I, j \in J} A_{ij} \quad (10.2)$$

$$a_{ij} = \frac{1}{|J|} \sum_{j \in J} A_{ij} \quad (10.3), \quad a_{Ij} = \frac{1}{|I|} \sum_{i \in I} A_{ij} \quad (10.4)$$



# Results and Evaluation

## Results of the football team dataset

	Original network	GN Algorithm	Louvain
Number of clusters	10	10	10
Average modularity	0.35	0.35	0.35
Execution time (s)	~10	~10	~10

A network graph illustrating connections between numbered nodes. The nodes are clustered and colored based on their groupings. The clusters are as follows:

- Green Cluster (Top):** Nodes 76, 95, 27, 96, 65, 56, 113, 87, 66, 112, 20, 17, 44, 63, 36, 103, 1, 33, 45, 105.
- Yellow Cluster (Top Right):** Nodes 37, 105, 109, 409.
- Blue Cluster (Center):** Nodes 48, 86, 91, 92, 58, 106, 10, 38, 13, 60, 64, 15, 6, 10, 39, 32, 47.
- Orange Cluster (Left):** Nodes 11, 98, 107, 52, 28, 0, 3, 102, 93, 6, 52, 88, 4, 81, 50, 74, 110, 140, 90, 84, 41, 40, 23, 9, 67, 2, 73, 46, 8, 78, 77, 51, 108, 1, 10, 22.
- Purple Cluster (Bottom Left):** Nodes 114, 110, 111, 112, 113.
- Pink Cluster (Bottom Center):** Nodes 19, 12, 79, 35, 26, 30, 29, 10, 94, 61, 155, 18, 34, 96, 31, 54.
- Red Cluster (Right):** Nodes 42, 19, 12, 79, 35, 26, 30, 29, 10, 94, 61, 155, 18, 34, 96, 31, 54.

The graph shows a complex web of connections between these numbered nodes across the different clusters.

The figure is a network graph with 100 nodes. Nodes are colored according to their cluster membership:

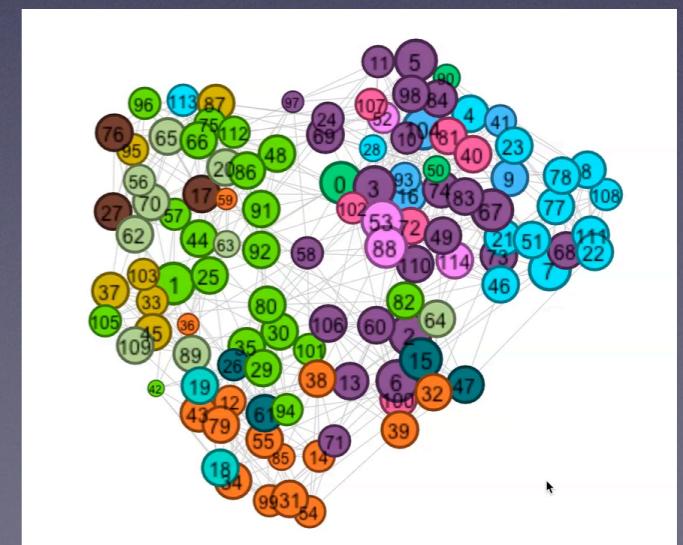
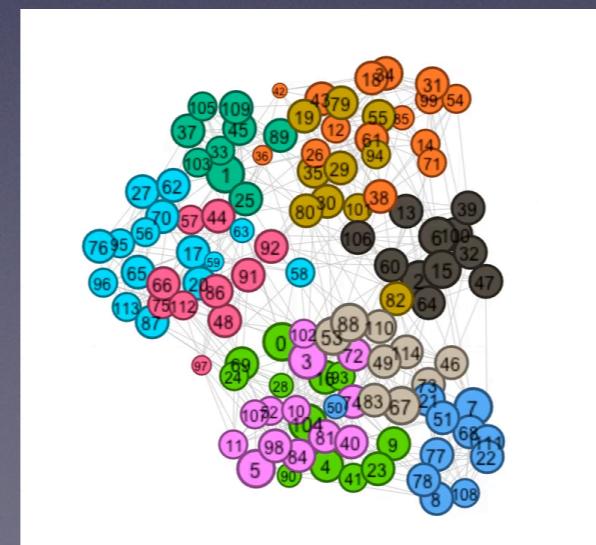
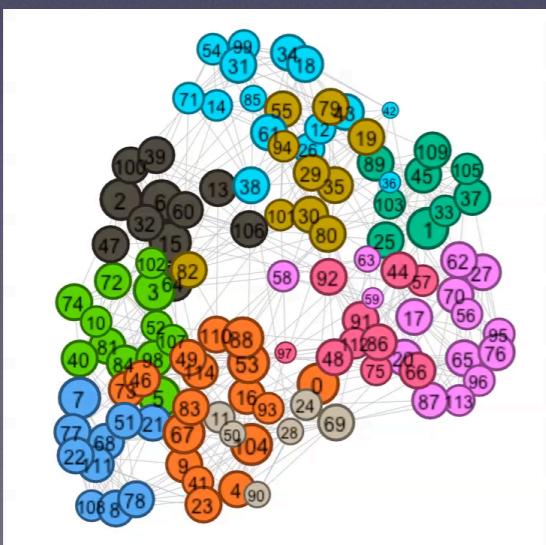
- Blue cluster: 27, 62, 105, 109, 45, 37, 66, 56, 70, 113, 96, 65, 17, 59, 29, 44, 103, 33, 1, 25, 63, 58, 48, 91, 92, 86, 112, 87, 79, 11, 55, 24, 97, 0, 102, 53, 88, 3, 28, 10, 104, 93, 6, 49, 110, 14, 82, 64, 2, 15, 32, 60, 61, 106, 13, 6, 100, 39, 47.
- Pink cluster: 43, 18, 34, 99, 31, 64, 26, 35, 29, 10, 38, 14, 71.
- Yellow cluster: 19, 79, 61, 55, 65, 94.
- Green cluster: 51, 7, 68, 112, 81, 40, 41, 23, 9, 67, 83, 77, 114, 8, 108.
- Orange cluster: 5, 98, 84, 4, 101, 50, 74, 11, 102, 103, 104, 90.
- Grey cluster: 12, 16, 18, 20, 21, 22, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113.

A network graph illustrating connections between 110 numbered nodes. The nodes are clustered by color: orange (top left), teal (bottom left), pink (bottom center), blue (center), dark grey/black (right side), and light green (far right). The nodes are interconnected by a web of lines representing their relationships.

# Unnormalised Spectral Clustering

# Normalised Spectral Clustering

# GA-Net



# Results and Evaluation

## Evaluation using Metrics in the football team dataset

- Internal Metrics
- External Metrics

- Modularity

$$Q = \sum_i^n (e_{ii} - a_i^2) , \quad a_i = \sum_j^n e_{ij}$$

- Internal Edge Density

$$IED = \frac{m_s}{n_s(n_s-1)/2}$$

- Inverse Conductance

$$\text{Inverse Conductance} = 1 - \frac{c_s}{2m_s + c_s}$$

- Normalised Mutual Information

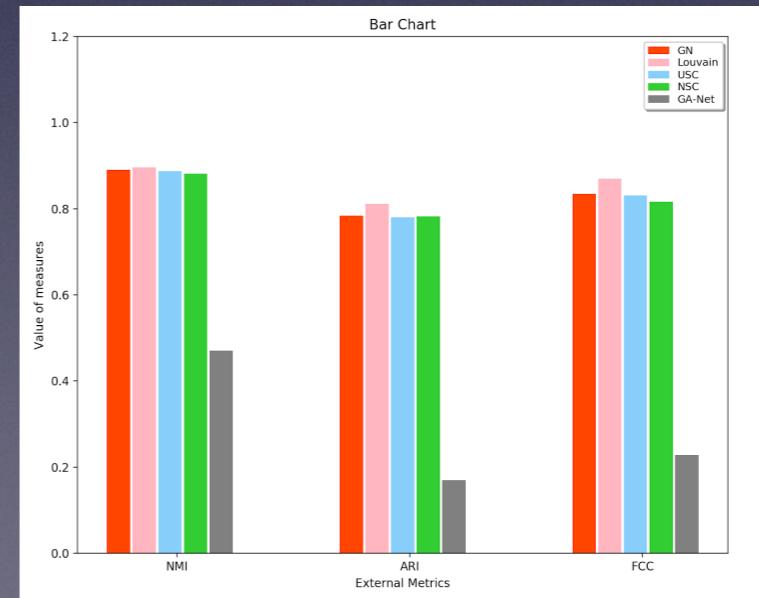
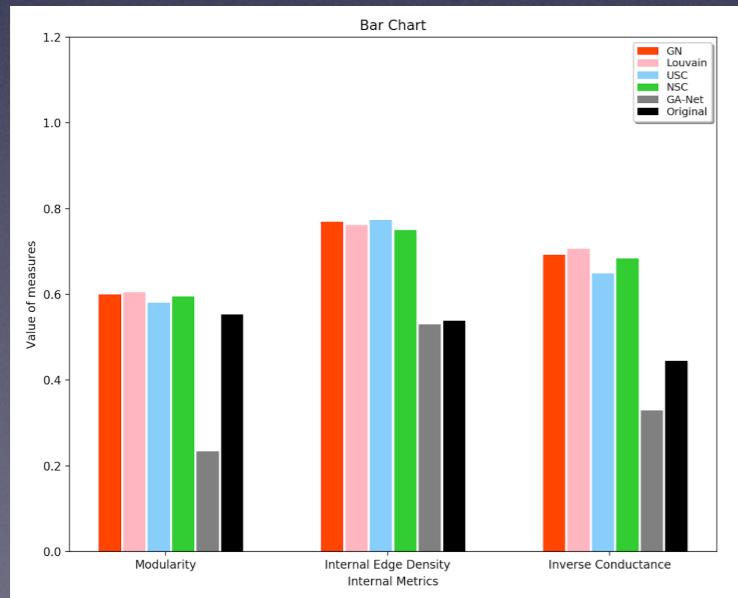
$$I(A, B) = \frac{-2 \sum_{i=1}^{c_A} \sum_{j=1}^{c_B} N_{ij} \log(N_{ij}N/N_i.N_j)}{\sum_{i=1}^{c_A} N_i \log(N_i/N) + \sum_{j=1}^{c_B} N_j \log(N_j/N)}$$

- Adjusted Rand Index

$$ARI = \frac{\binom{n}{2} \sum_{r=1}^R \sum_{c=1}^C \binom{t_{rc}}{2} - \left[ \sum_{r=1}^R \binom{t_r}{2} \sum_{c=1}^C \binom{t_c}{2} \right]}{\frac{1}{2} \binom{n}{2} \left[ \sum_{r=1}^R \binom{t_r}{2} + \sum_{c=1}^C \binom{t_c}{2} \right] - \left[ \sum_{r=1}^R \binom{t_r}{2} \sum_{c=1}^C \binom{t_c}{2} \right]}$$

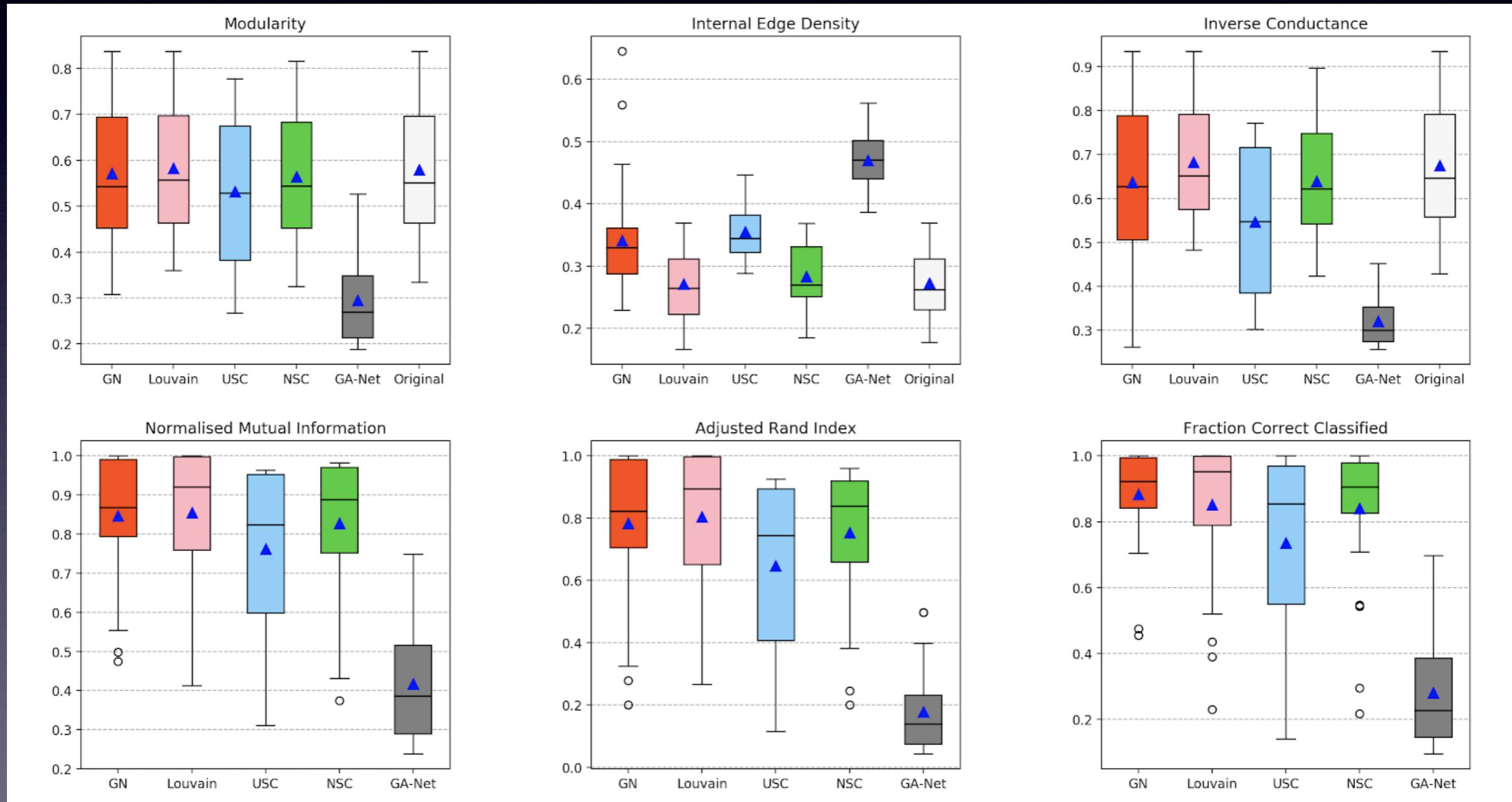
- Fraction of Correctly Classified nodes

an element in estimated community is assigned correctly if at least half neighbours assigned correctly



# Results and Evaluation

Evaluation using Metrics in several synthetic networks



# Results and Evaluation

Comparing Algorithms using plain language

Name	Advantages	Disadvantages
GN	Simple to realise; Do not need to set the number of communities	Time-consuming
Louvain	The fastest one among these; Do not need to set the number of communities	A complex algorithm
USC/NSC	Perform good, but the normalised one are better	Needs to set the number of communities
GA-Net	Do not need to set the number of communities	Has low accuracy

# Results and Evaluation

## Comparing Algorithms using statistical test

Oneway test (Anova, F-test)

### Internal metric -Modularity

Algorithms VS Original	T values	Degree fo freedom	p-values
GN	3.8723	29	0.0005652
Louvain	-2.2519	29	0.03207
USC	9.3677	29	2.843E-10
NSC	12.686	29	2.321E-13
GA-Net	26.926	29	2.2E-16

### External metric -Normalised Mutual Information

Comparison	T values	Degree fo freedom	p-values
USC vs NSC	-6.4289	29	4.946E-07
GN vs Louvain	-0.63826	29	0.5283
Louvain vs GA	19.445	29	2.2E-16

### Anova for all metrics

Results from GN and Louvain are similar.

Metrics	F values	Degree of freedom	p-values
Modularity	18.838	(5, 174)	6.027E-15
Internal Edge Density	49.872	(5, 174)	2.2E-16
Inverse Conductance	28.029	(5, 174)	2.2E-16
Normalised Mutual Information	31.271	(4, 145)	2.2E-16
Adjusted Rand Index	42.39	(4, 145)	2.2E-16
Fraction Correct Classified	44.53	(4, 145)	2.2E-16

Thanks for watching

Yawen Sun  
1836786