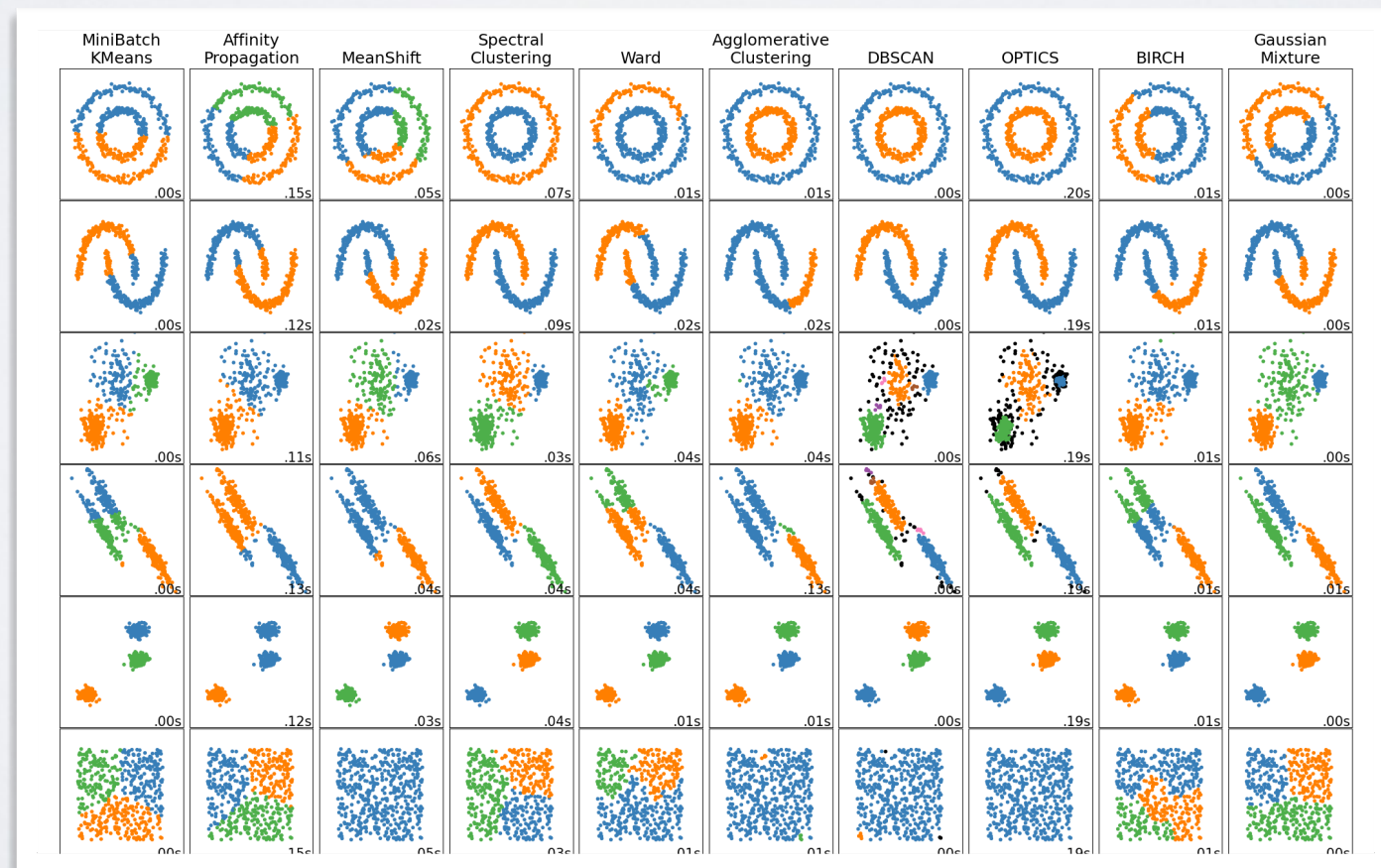


# COMMUNITY DETECTION (GRAPH CLUSTERING)

# COMMUNITY DETECTION

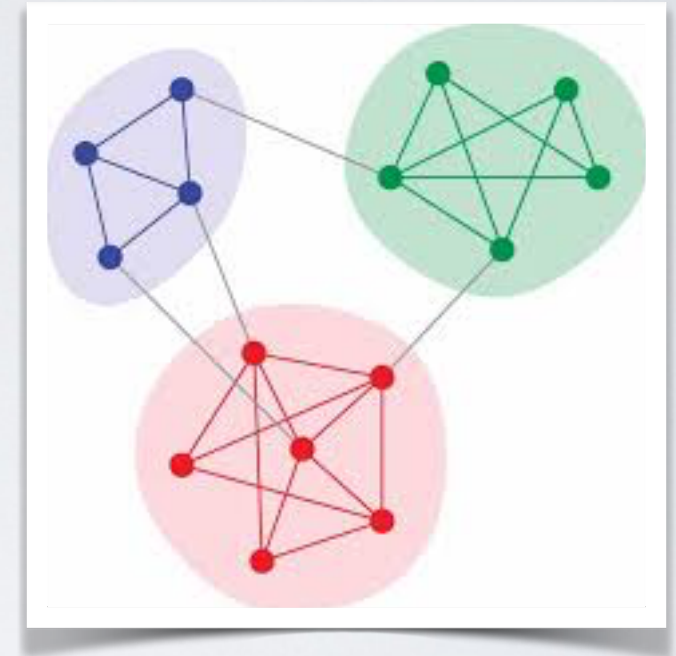
- Community detection is equivalent to “clustering” in unstructured data
- Similar problems: what is a good community ?



# COMMUNITY DETECTION

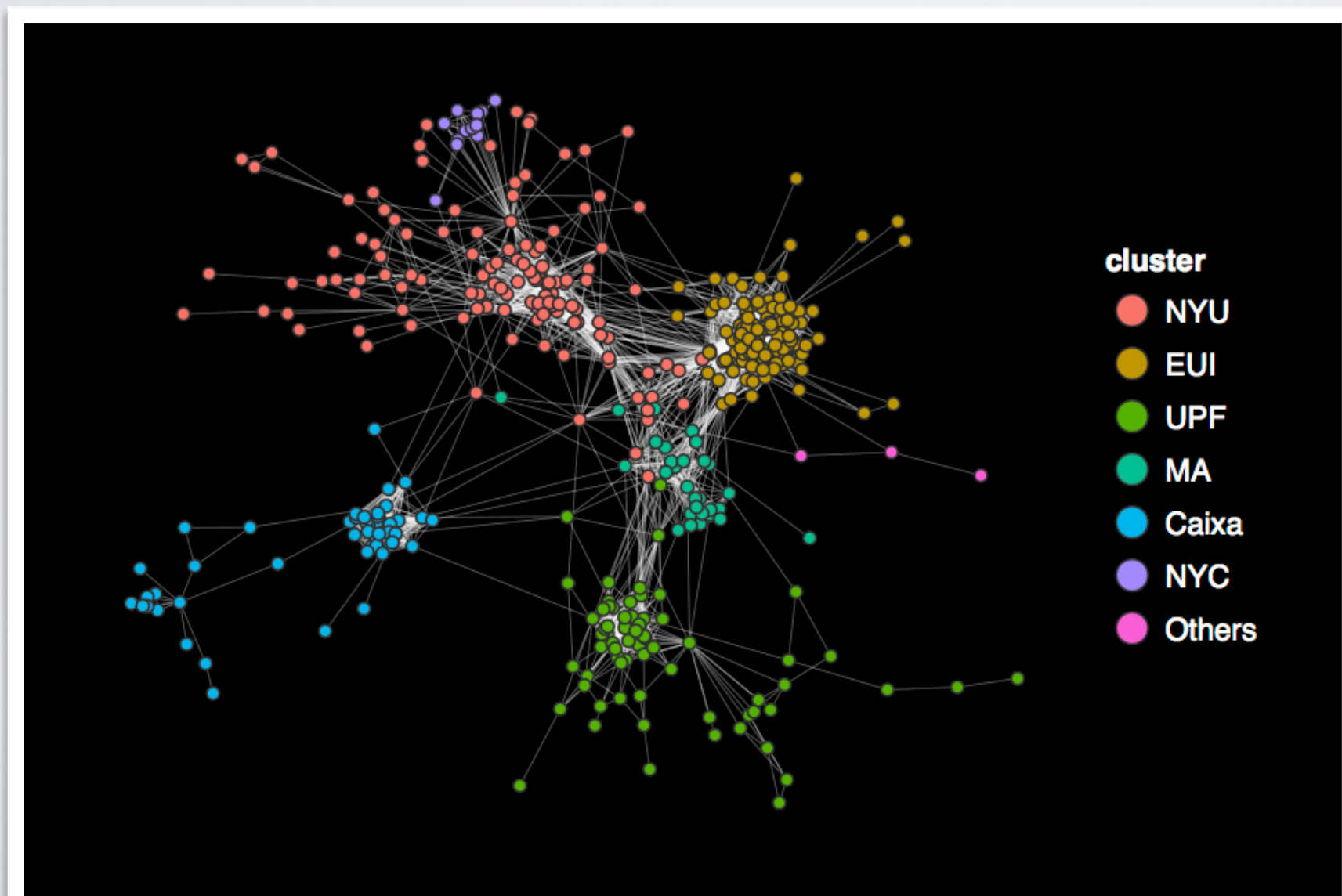
- Community detection:

- ▶ Find groups of nodes that are:
  - Strongly connected to each other
  - Weakly connected to the rest of the network
  - Ideal form: each community is 1) A clique, 2) A separate connected component
- ▶ No formal definition
- ▶ Hundreds of methods published since 2003



# COMMUNITY STRUCTURE IN REAL GRAPHS

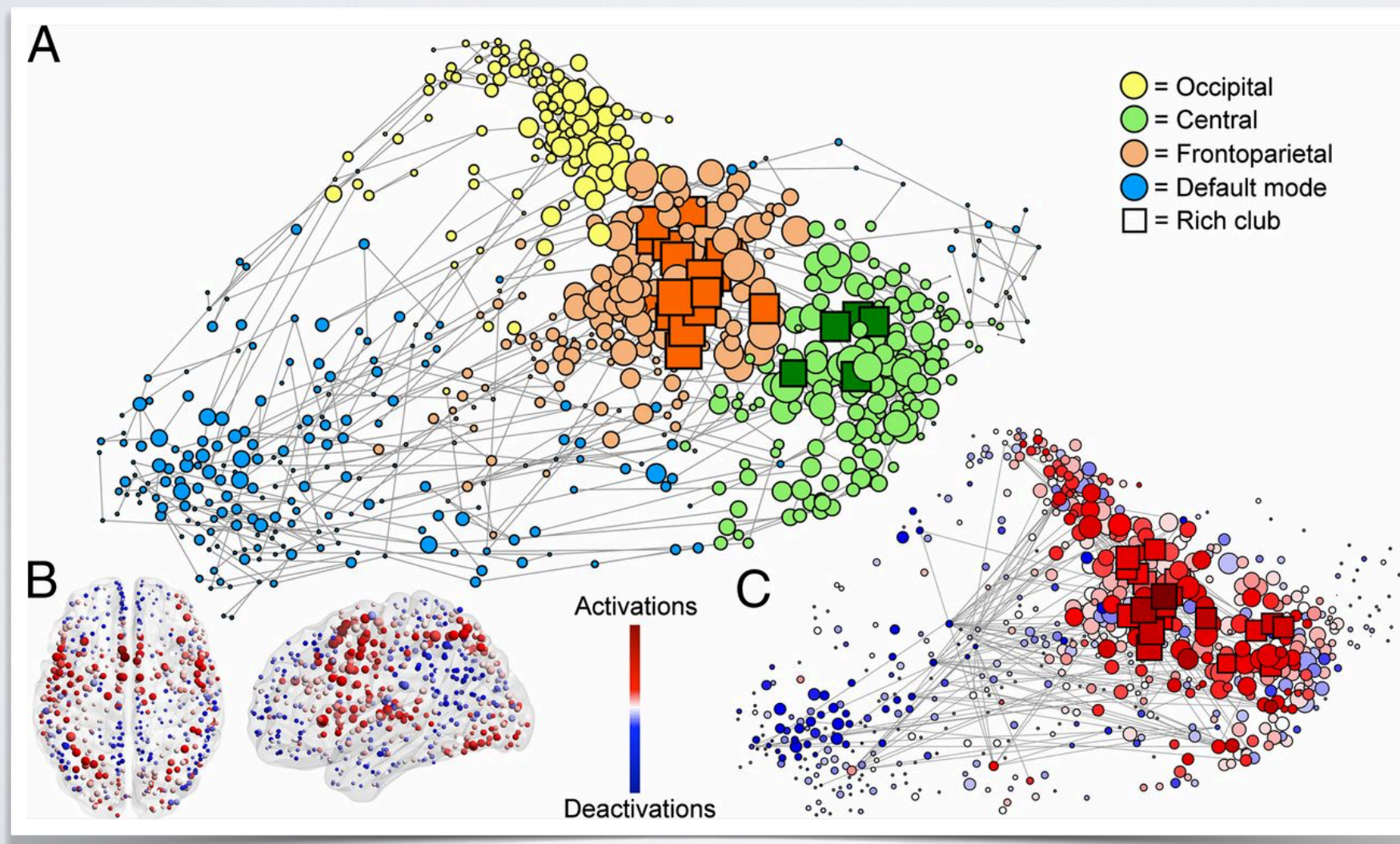
- If you plot the graph of your facebook friends, it looks like this





# COMMUNITY STRUCTURE IN REAL GRAPHS

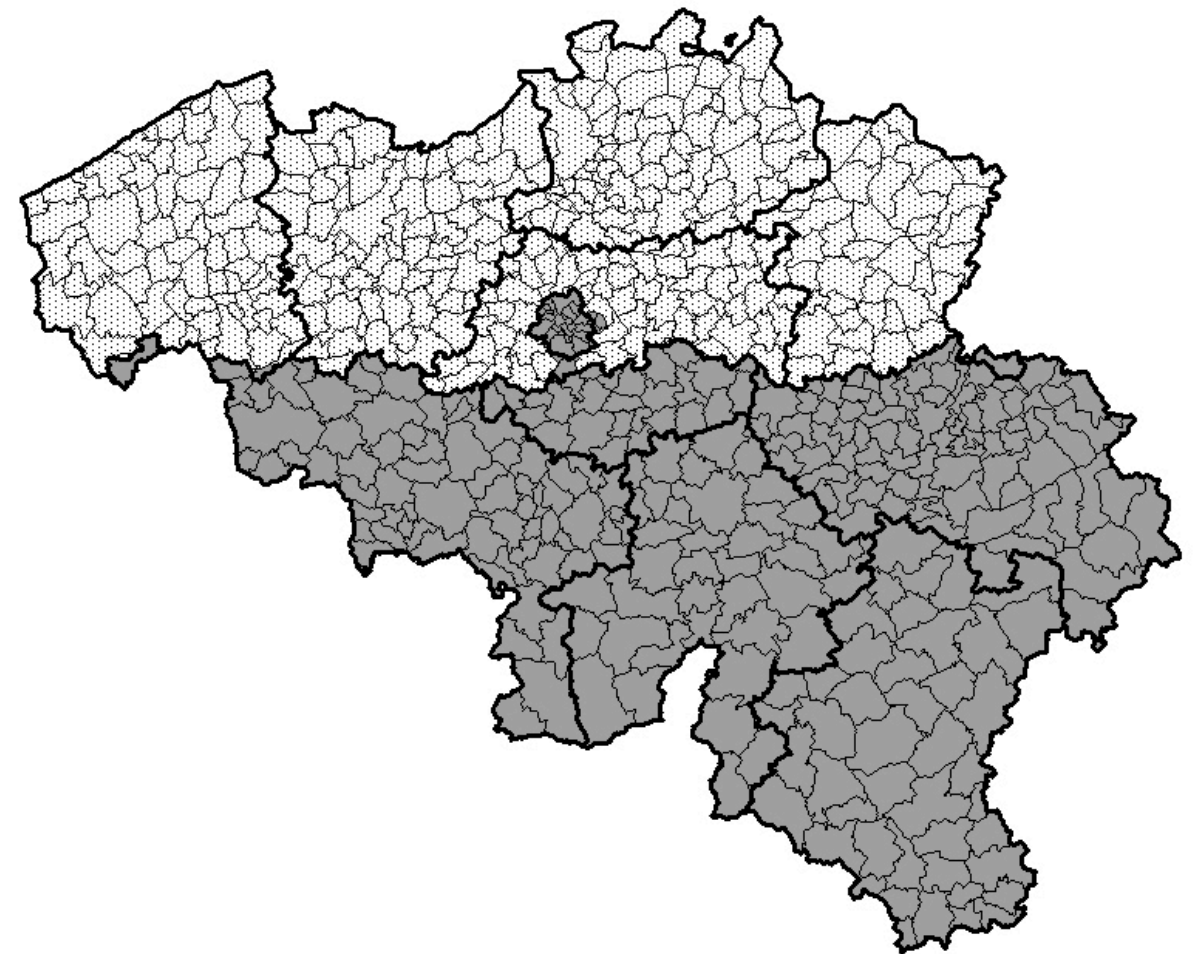
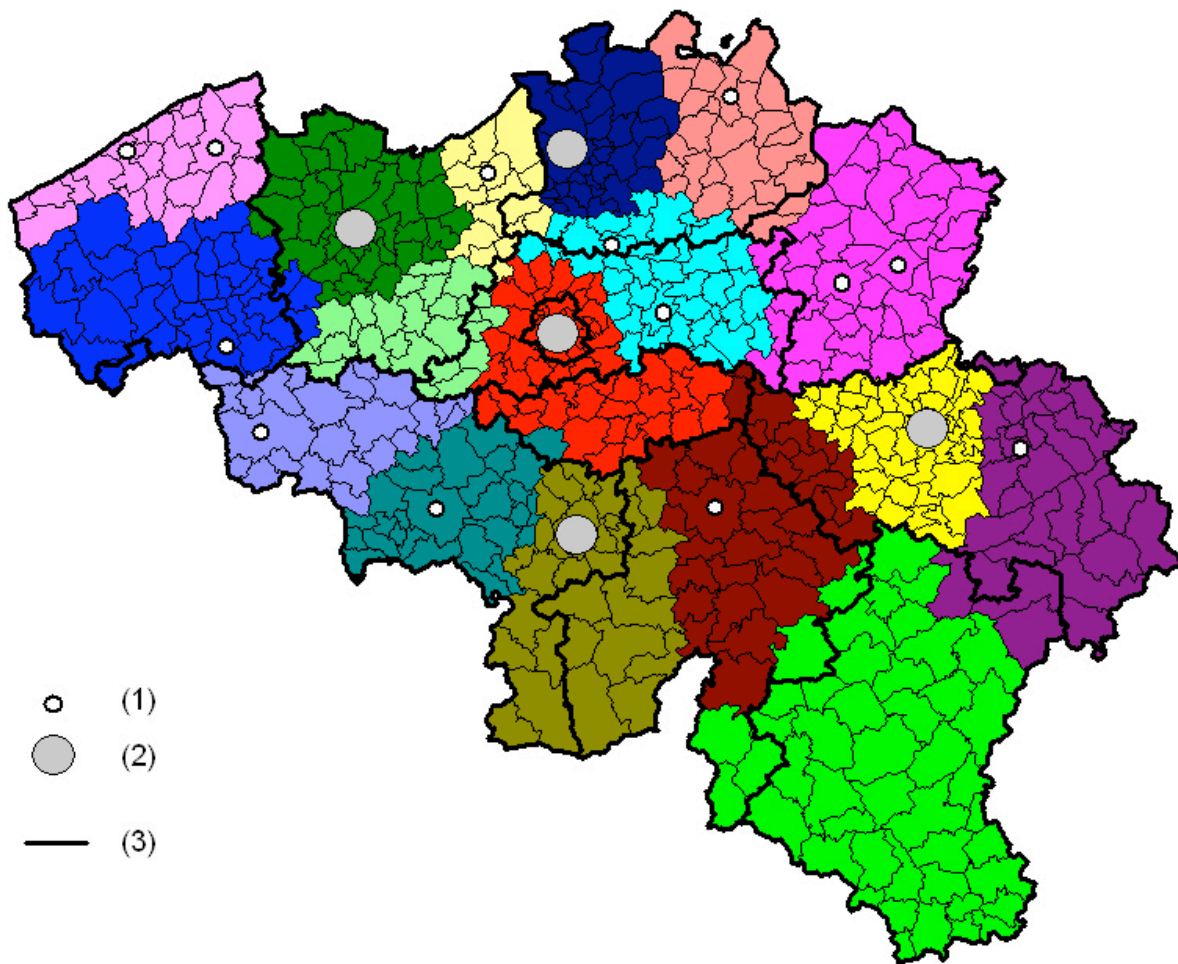
- Connections in the brain ?





# COMMUNITY STRUCTURE IN REAL GRAPHS

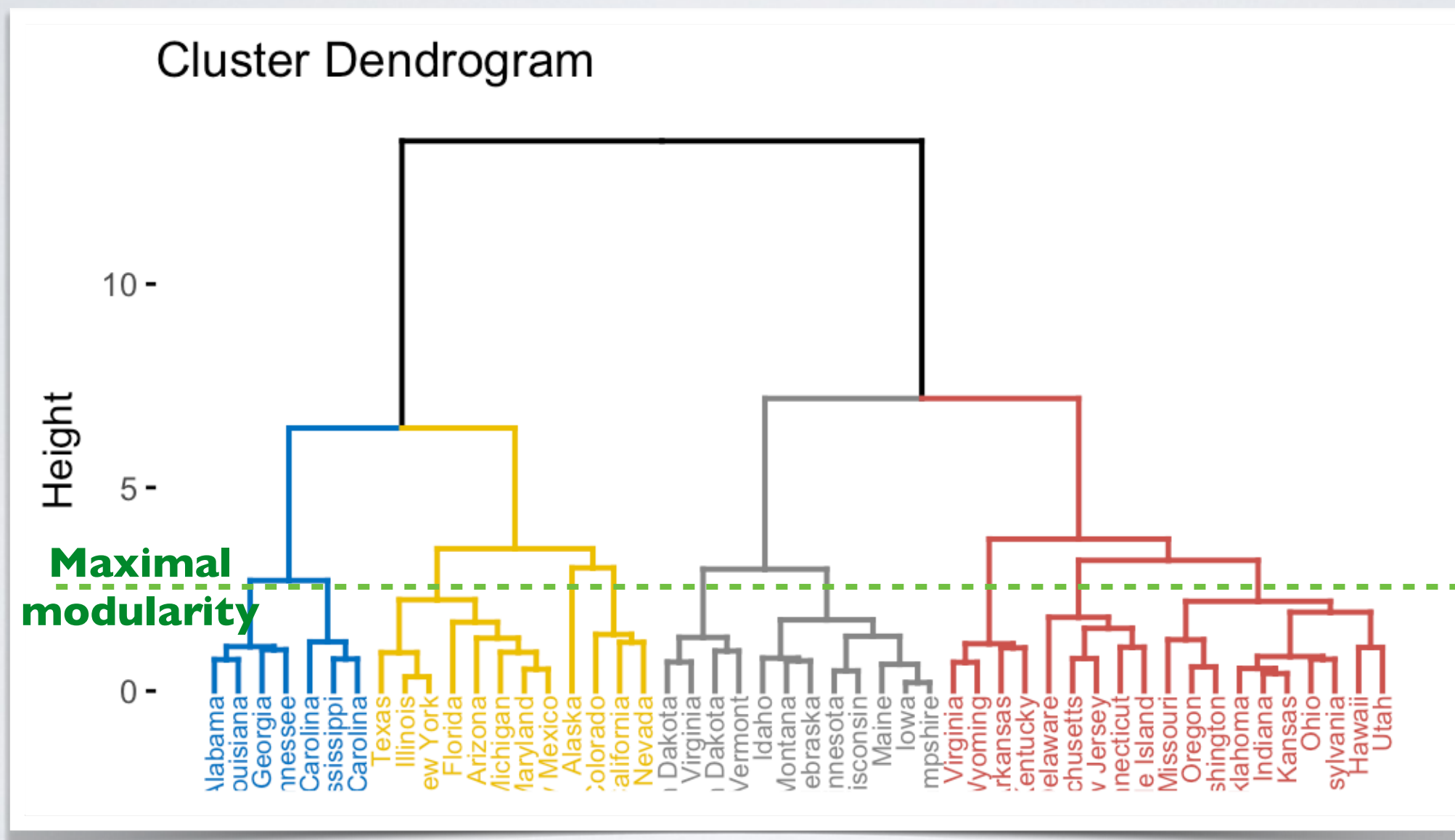
- Phone call communications in Belgium ?



# FIRST METHOD BY GIRVAN & NEWMAN

- 1) Compute the betweenness of all edges
- 2) Remove the edge of highest betweenness
- 3) Repeat until all edges have been removed
  - Connected components are communities
- => It is called a *divisive* method
- => What you obtain is a dendrogram
- How to cut this dendrogram at the *best* level ?

# FIRST METHOD BY GIRVAN & NEWMAN





# FIRST METHOD BY GIRVAN & NEWMAN

- Introduction of the **Modularity**
- The modularity is computed for a partition of a graph
  - (each node belongs to one and only one community)
- It compares :
  - The **observed** *fraction of edges inside communities*
  - To the **expected** *fraction of edges inside communities* in a random network

# MODULARITY

$$Q = \frac{1}{(2m)} \sum_{vw} \left[ A_{vw} - \frac{k_v k_w}{(2m)} \right] \delta(c_v, c_w)$$

Original formulation

# MODULARITY

$$Q = \frac{1}{(2m)} \sum_{vw} \left[ A_{vw} - \frac{k_v k_w}{(2m)} \right] \delta(c_v, c_w)$$

Sum over all pairs of nodes

# MODULARITY

$$Q = \frac{1}{(2m)} \sum_{vw} \left[ A_{vw} - \frac{k_v k_w}{(2m)} \right] \delta(c_v, c_w)$$

| if in same community



# MODULARITY

$$Q = \frac{1}{(2m)} \sum_{vw} \left[ A_{vw} - \frac{k_v k_w}{(2m)} \right] \delta(c_v, c_w)$$

| if there is an edge between them

# MODULARITY

$$Q = \frac{1}{(2m)} \sum_{vw} \left[ A_{vw} - \frac{k_v k_w}{(2m)} \right] \delta(c_v, c_w)$$

Probability of an edge in  
a configuration model  
(Edges at random, keeping degrees)

# MODULARITY

Can also be defined  
as a sum by community

$$Q = \frac{1}{L} \sum_{i=1}^{|C|} (L_i - \frac{1}{2} K_i^2)$$

with  $L_i = L(H(c_i))$  the number of edges inside community  $i$  and  $K_i = \sum_{u \in c_i} k_u$  the sum of degrees of nodes in community  $i$ .

# MODULARITY

- Modularity compares the observed network to a **null model**
  - Usually the configuration model
    - Multi-edges and loops are allowed
  - Other models could be used, such as ER random graphs.
- Natural extension to weighted/multi-edge networks



# FIRST METHOD BY GIRVAN & NEWMAN

- Back to the method:
  - Create a dendrogram by removing edges
  - Cut the dendrogram at the best level using modularity
- => In the end, your objective is... to optimize the Modularity, right ?
- Why not optimizing it directly !

# MODULARITY MAXIMIZATION

- From 2004 to 2008: The golden age of Modularity
- Scores of methods proposed to maximize it
  - ▶ Graph spectral approaches
  - ▶ Meta-heuristics approaches (simulated annealing, multi-agent...)
  - ▶ Local/Global approaches...
- => 2008: the Louvain algorithm

# LOUVAIN ALGORITHM

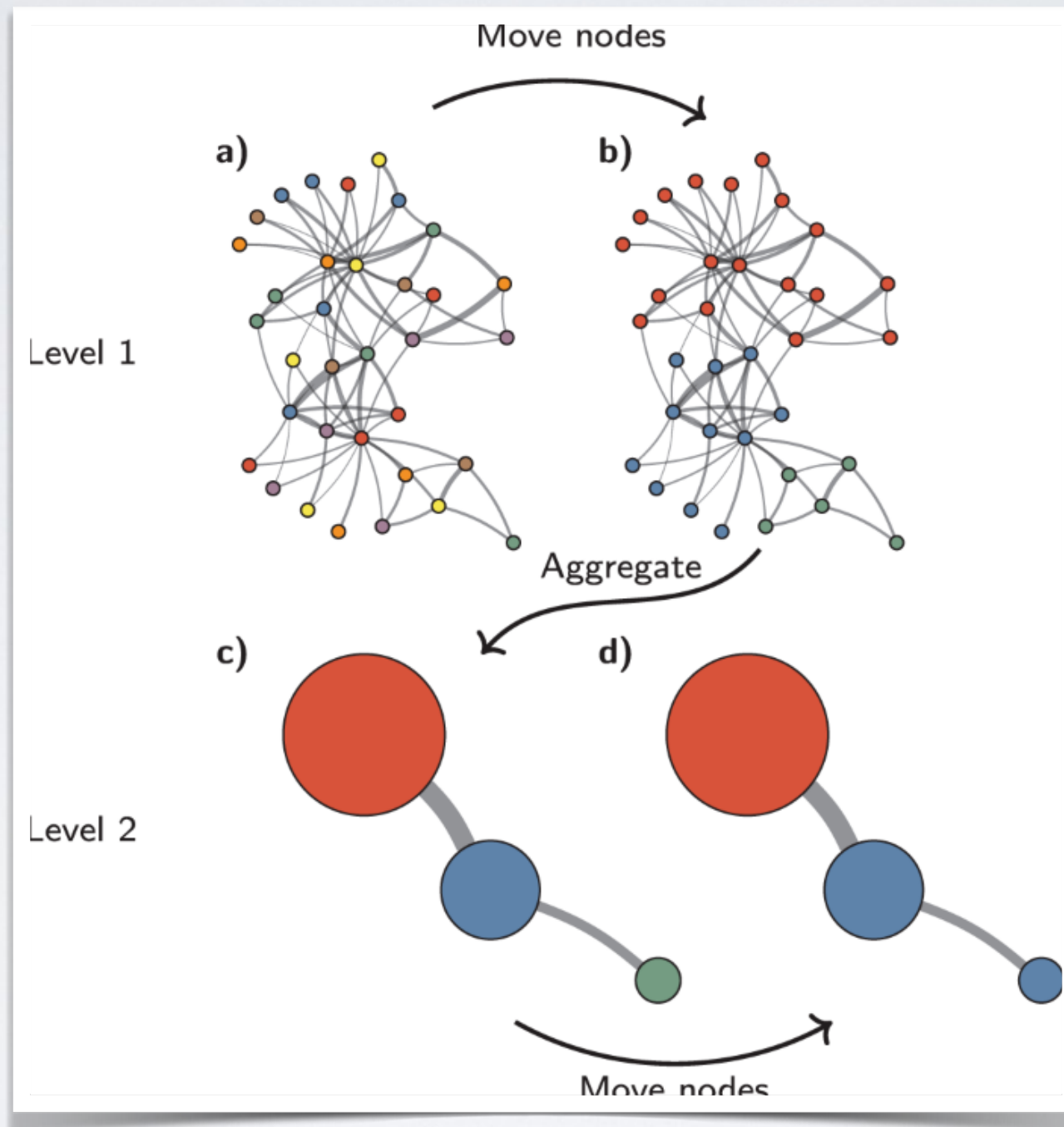
- Simple, greedy approach
  - Easy to implement
  - Fast
- Yields a hierarchical community structure
- Beat state of the art on all aspects (when introduced)
  - Speed
  - Max modularity obtained
  - Do not fall in some traps (see later)

# LOUVAIN ALGORITHM

- Each node start in its own community
- Repeat until convergence
  - FOR each node:
    - FOR each neighbor:
      - if adding node to its community increase modularity, do it
- When converged, create an *induced network*
  - Each community becomes a node
  - Edge weight is the sum of weights of edges between them
- Trick: Modularity is computed *by community*
  - Global Modularity = sum of modularities of each community



# LOUVAIN ALGORITHM



# ALTERNATIVES

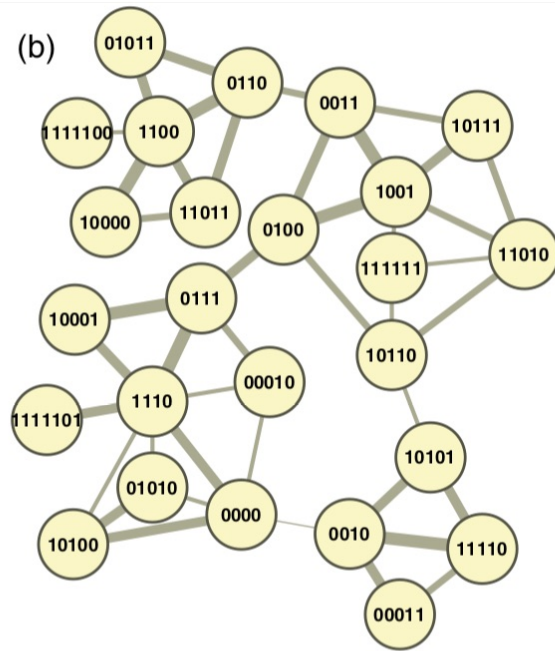
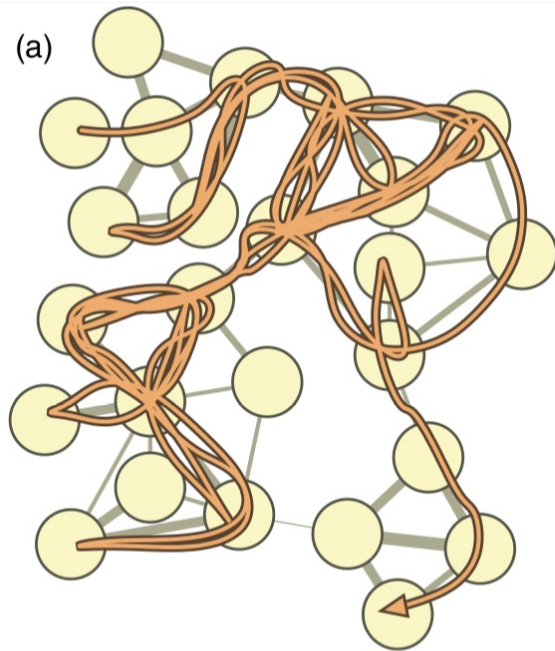
- Most serious alternatives
  - Infomap (based on information theory —compression)
  - Stochastic block models (bayesian inference)
- These methods have a clear definition of what are good communities. Theoretically grounded

# INFOMAP

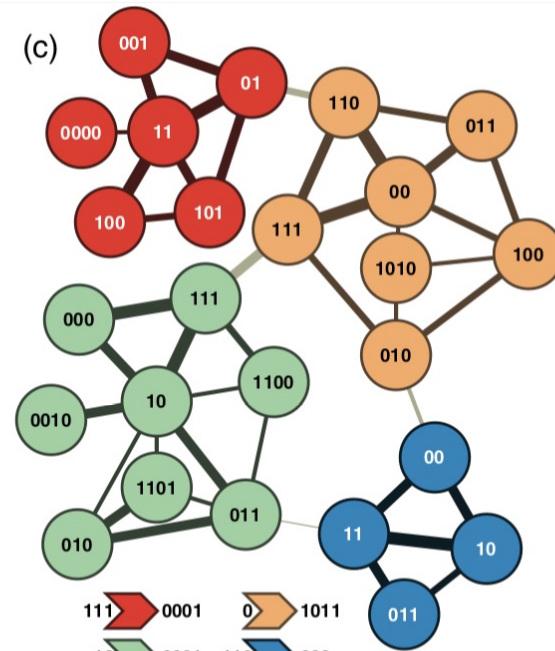
- [Rosvall & Bergstrom 2009]
- Find the partition minimizing the *description* of any *random walk* on the network
- We want to *compress* the description of random walks



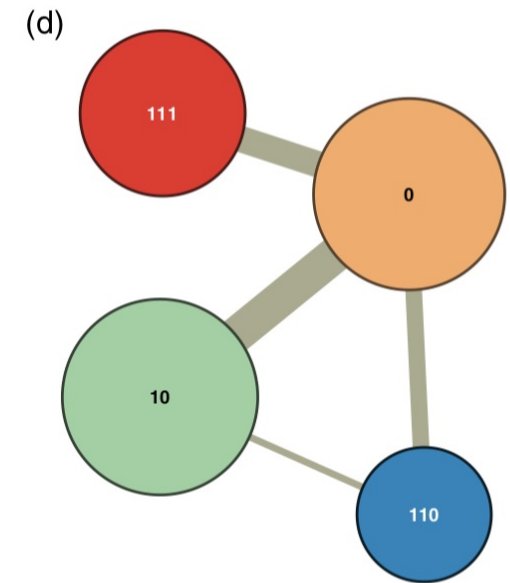
# INFOMAP



1111100 1100 0110 11011 10000 11011 0110 0011 10111 1001  
 0011 1001 0100 0111 10001 1110 0111 10001 0111 1110 0000  
 1110 10001 0111 1110 0111 1110 1111101 1110 0000 10100 0000  
 1110 10001 0111 0100 10110 11010 10111 1001 0100 1001 10111  
 1001 0100 1001 0100 0011 0100 0011 0110 11011 0110 0011 0100  
 1001 10111 0011 0100 0111 10001 1110 10001 0111 0100 10110  
 111111 10110 10101 11110 00011



111 0000 11 01 101 100 101 01 0001 0 110 011 00 110 00 111  
 1011 10 111 000 10 111 000 111 10 011 10 000 111 10 111 10  
 0010 10 011 010 011 10 000 111 0001 0 111 010 100 011 00 111  
 00 011 00 111 00 111 110 111 110 1011 111 01 101 01 0001 0 110  
 111 00 011 110 111 1011 10 111 000 10 000 111 0001 0 111 010  
 1010 010 1011 110 00 10 011



111 0000 11 01 101 100 101 01 0001 0 110 011 00 110 00 111  
 1011 10 111 000 10 111 000 111 10 011 10 000 111 10 111 10  
 0010 10 011 010 011 10 000 111 0001 0 111 010 100 011 00 111  
 00 011 00 111 00 111 110 111 110 1011 111 01 101 01 0001 0 110  
 111 00 011 110 111 1011 10 111 000 10 000 111 0001 0 111 010  
 1010 010 1011 110 00 10 011

Random  
walk

Description  
Without  
Communities

With communities

**Huffman coding:** short codes for frequent items

**Prefix free:** no code is a prefix of another one (avoid fix length/separators)



# The Infomap method

## Finding the optimal partition M:

- Shannon's source coding theorem (Shannon's entropy)

for a probability distribution  $P = \{p_i\}$  such that  $\sum_i p_i = 1$ , the lower limit of the per-step code-length is

$$L(\mathcal{P}) = H(\mathcal{P}) \equiv - \sum_i p_i \log p_i.$$

- Minimise the expected description length of the random walk

Sum of Shannon entropies of multiple codebooks weighted by the rate of usage

probability of between modules movements of a RW, i.e. the rate of usage of the index codebook

probability of within modules movements of a RW, i.e. the rate of usage of the module codebook

$$L(\mathbf{M}) = q H(\mathcal{Q}) + \sum_{i=1}^m p_i H(\mathcal{P}^i)$$

Expected decryption length of partition M

Entropy of movement between modules, i.e. the frequency weighted average length of codewords

Entropy of movement inside modules, i.e. the frequency weighted average length of codewords in the module codebook

## Algorithm

1. Compute the fraction of time each node is visited by the random walker ([Power-method on adjacency matrix](#))
2. Explore the space of possible partitions ([deterministic greedy search algorithm - similar to Louvain but here we join nodes if they decrease the description length](#))
3. Refine the results with simulated annealing ([heat-bath algorithm](#))

# INFOMAP

- To sum up:
  - Infomap defines a *quality function* for a partition different than modularity
  - Any algorithm can be used to optimize it (like Modularity)
- Advantage:
  - Infomap can recognize random networks (no communities)

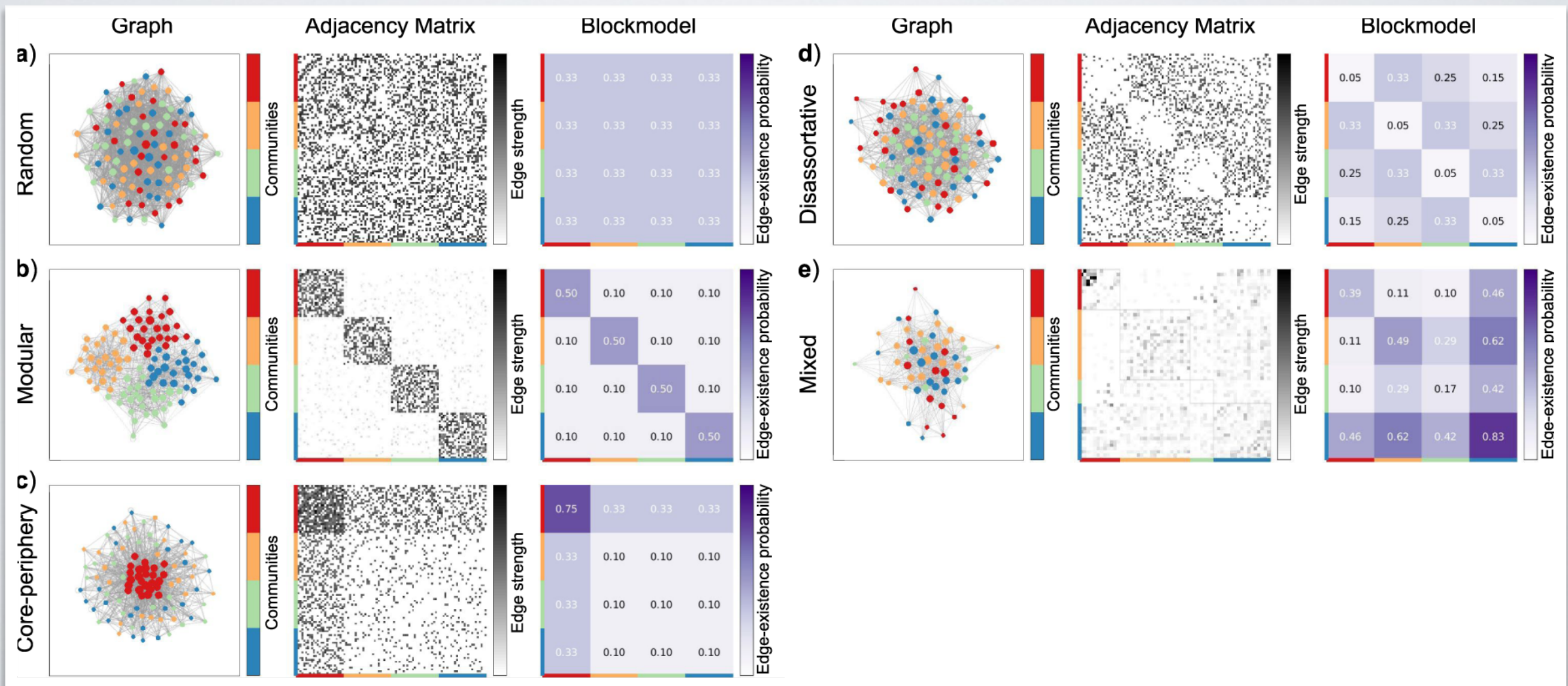
# STOCHASTIC BLOCK MODELS

- Stochastic Block Models (SBM) are based on statistical models of networks
- They are in fact more general than usual communities.
- The model is:
  - Each node belongs to 1 and only 1 community
  - To each pair of communities, there is an associated density (probability of each edge to exist)



# STOCHASTIC BLOCK MODELS

- SBM can represent different things:
  - Associative SBM: density inside nodes of a same communities  $\gg$  density of pairs belonging to different communities.





# STOCHASTIC BLOCK MODELS

- General idea of SBM community detection:
  - Specify the desired number of cluster
  - Find parameters to optimize the maximum likelihood
    - Principle: The best parameters are those that allow to generate the observed network with the highest probability
- Main weakness of this approach
  - Number of clusters must be specified (avoid trivial solution)
- MDL (Minimum Description Length) approaches exist to automatically find the number of blocks

# EVALUATION OF COMMUNITY STRUCTURE

# EVALUATION

- Similar to clustering:
  - ▶ Intrinsic/Internal evaluation
    - Partition quality function
    - Individual Community quality function
  - ▶ Comparison of observed communities and expected communities
    - Synthetic networks with community structure
    - Real networks with Ground Truth

# INTRINSIC EVALUATION



# INTRINSIC EVALUATION

- Partition quality function
  - Already defined: **Modularity**, **graph compression**, etc.

- Quality function for individual community

- Internal Clustering Coefficient

- Conductance:  $\frac{|E_{out}|}{|E_{out}| + |E_{in}|}$

- Fraction of external edges

$|E_{in}|, |E_{out}|$ :  
# of links to nodes inside  
(respectively, outside) the  
community

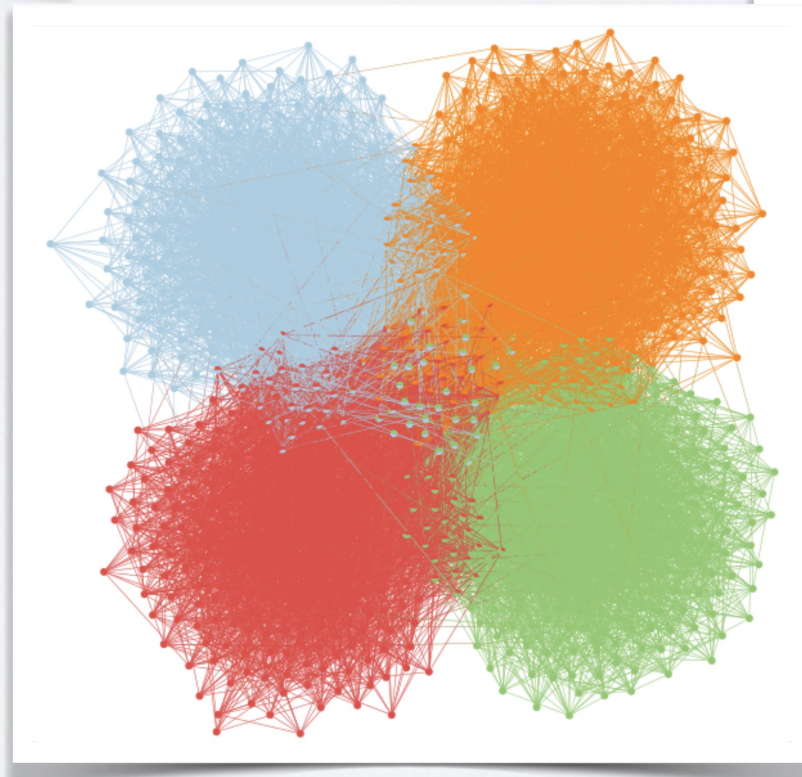
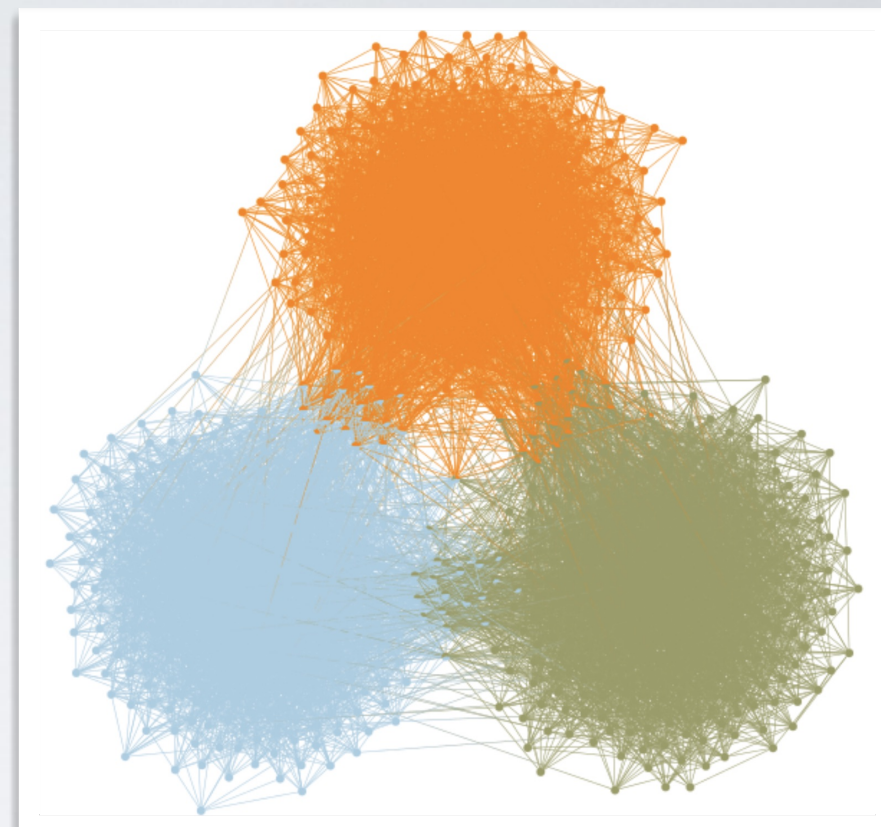
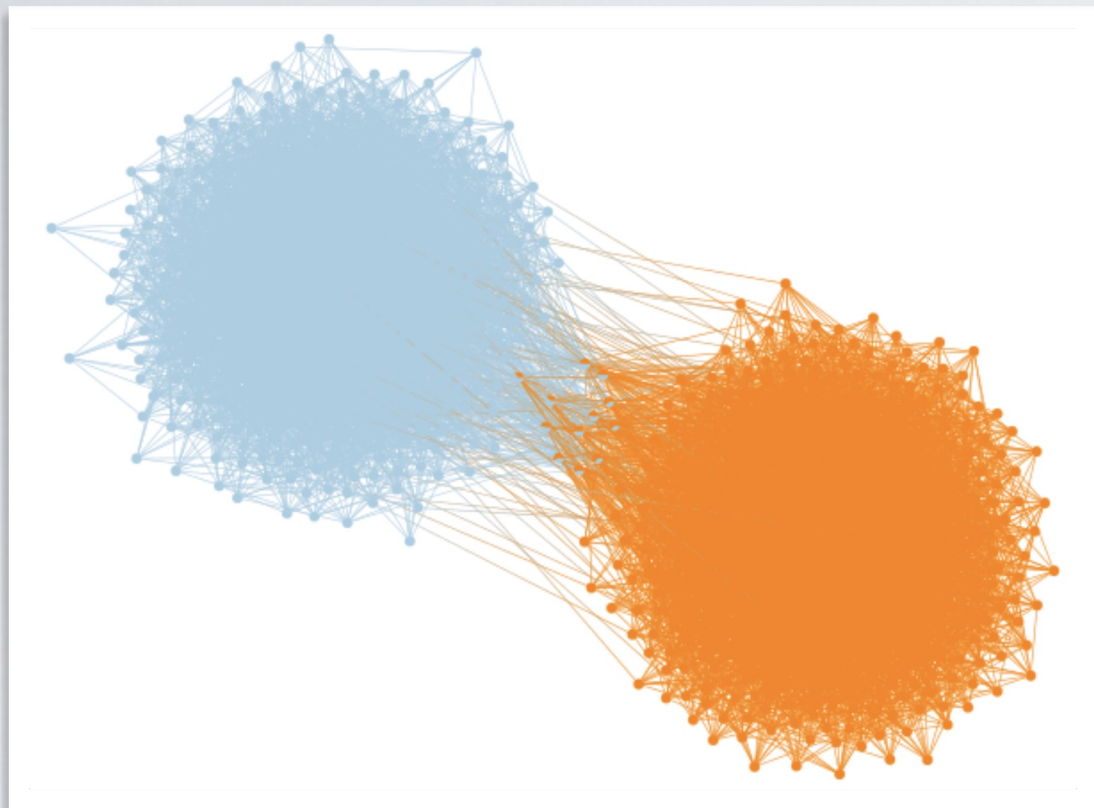
COMPARISON WITH  
GROUND TRUTH

# SYNTHETIC NETWORKS

- Planted Partition models:
  - ▶ Another name for SBM with manually chosen parameters
    - Assign degrees to nodes
    - Assign nodes to communities
    - Assign density to pairs of communities
    - Attribute randomly edges
  - ▶ Problem: how to choose parameters?
    - Either oversimplifying (all nodes same degrees, all communities same #nodes, all intern densities equals...)
    - Or ad-hoc process (sample values from distributions)



# SYNTHETIC NETWORKS



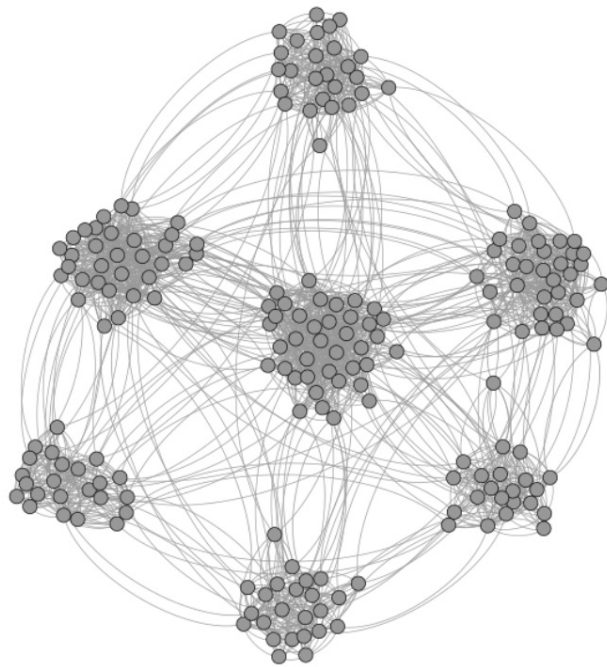


# SYNTHETIC NETWORKS

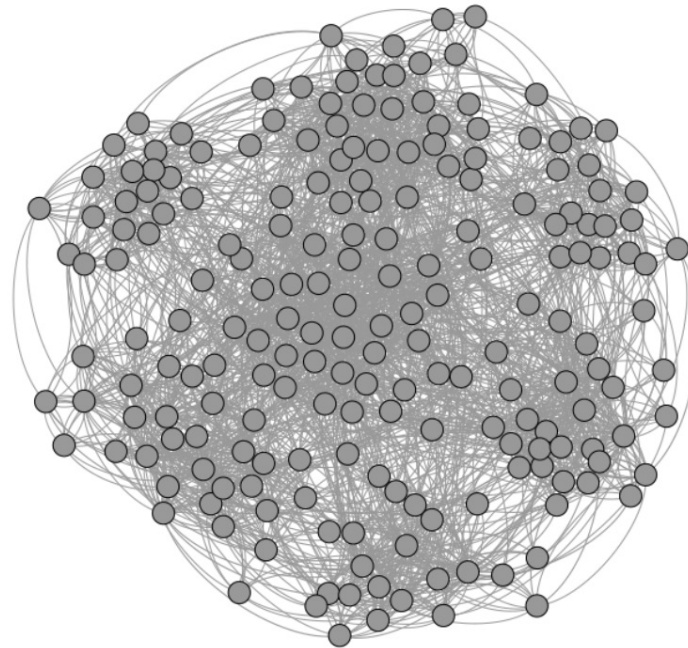
- LFR Benchmark [Lancichinetti 2008]
  - High level parameters:
    - Slope of the power law distribution of degrees/community sizes
    - Avg Degree, Avg community size
    - Mixing parameter: fraction of external edges of each node
  - Varying the mixing parameter makes community more or less well defined
- Currently the most popular

# SYNTHETIC NETWORKS

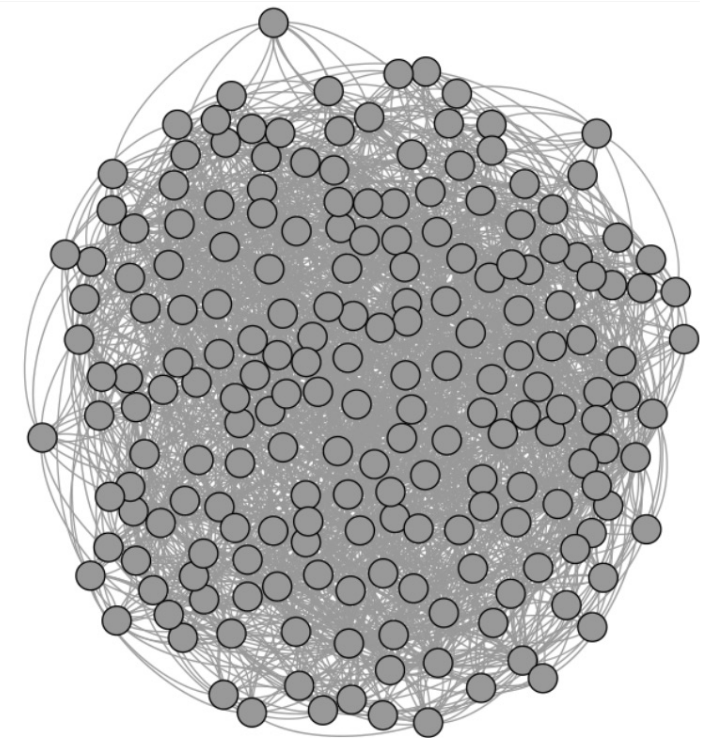
LFR Benchmark Networks with 200 Nodes



$\mu=0.1$   
#Edges=2206



$\mu=0.3$   
#Edges=2628



$\mu=0.5$   
#Edges=2462

# OTHER TYPES OF COMMUNITIES



# OVERLAPPING COMMUNITIES

- In real networks, communities are often overlapping
  - ▶ Some of your High-School friends might be also University Friends
  - ▶ A colleague might be a member of your family
  - ▶ ...
- Overlapping community detection is considered much harder
  - ▶ And is not well defined
- Difference between “attributes” and overlapping communities ?
  - ▶ Community of Women, Community of 17-19yo, Community of fans of...



# HIERARCHICAL COMMUNITIES

