

# Community detection



# LEARNING OUTCOMES

UNDERSTAND WHAT **COMMUNITIES** ARE

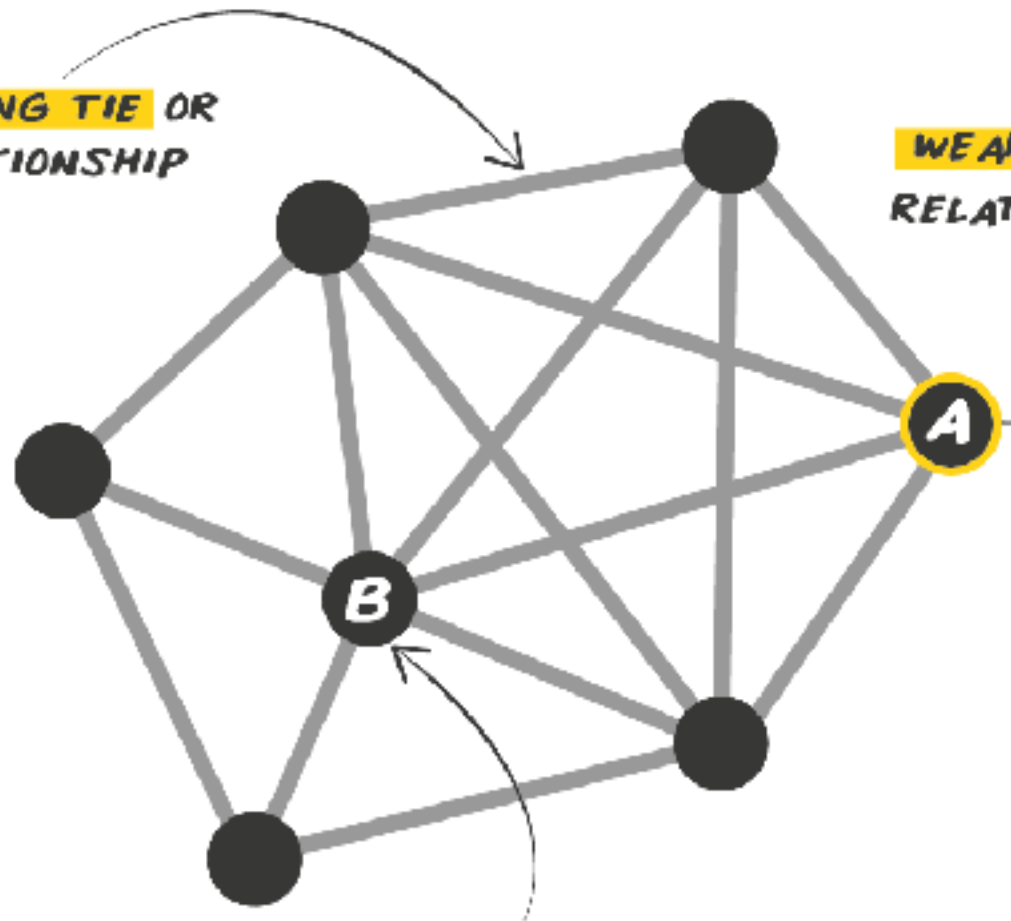
BE ABLE TO **DESCRIBE A NETWORK** IN TERMS OF COMMUNITIES

LEARN DIFFERENT TYPES OF COMMUNITY **CLASSIFICATIONS**



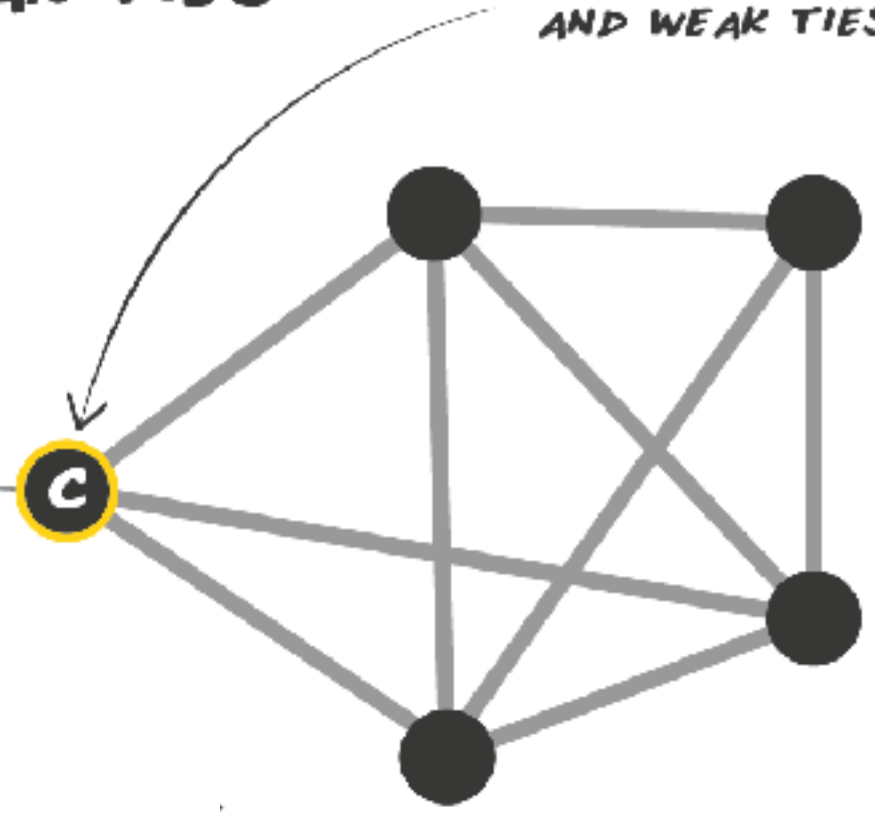
# GRANOVETTER'S STRENGTH OF WEAK TIES

**STRONG TIE OR RELATIONSHIP**



**WEAK TIE OR RELATIONSHIP**

IT'S VALUABLE TO HAVE A COMBINATION OF STRONG AND WEAK TIES



EVEN THOUGH B HAS MORE TIES THAN A, ALL THOSE TIES LIKELY HAVE THE SAME INFORMATION BECAUSE THEY ALL KNOW EACH OTHER WELL

FOR EXAMPLE, A CAN SHARE INFORMATION WITH C THAT C WOULDN'T GET FROM ANYONE ELSE IN THEIR GROUP, AND VICE VERSA.



LINCOLN EIGHTH  
**ALL VALLEY KARATE CHAMPIONSHIP**

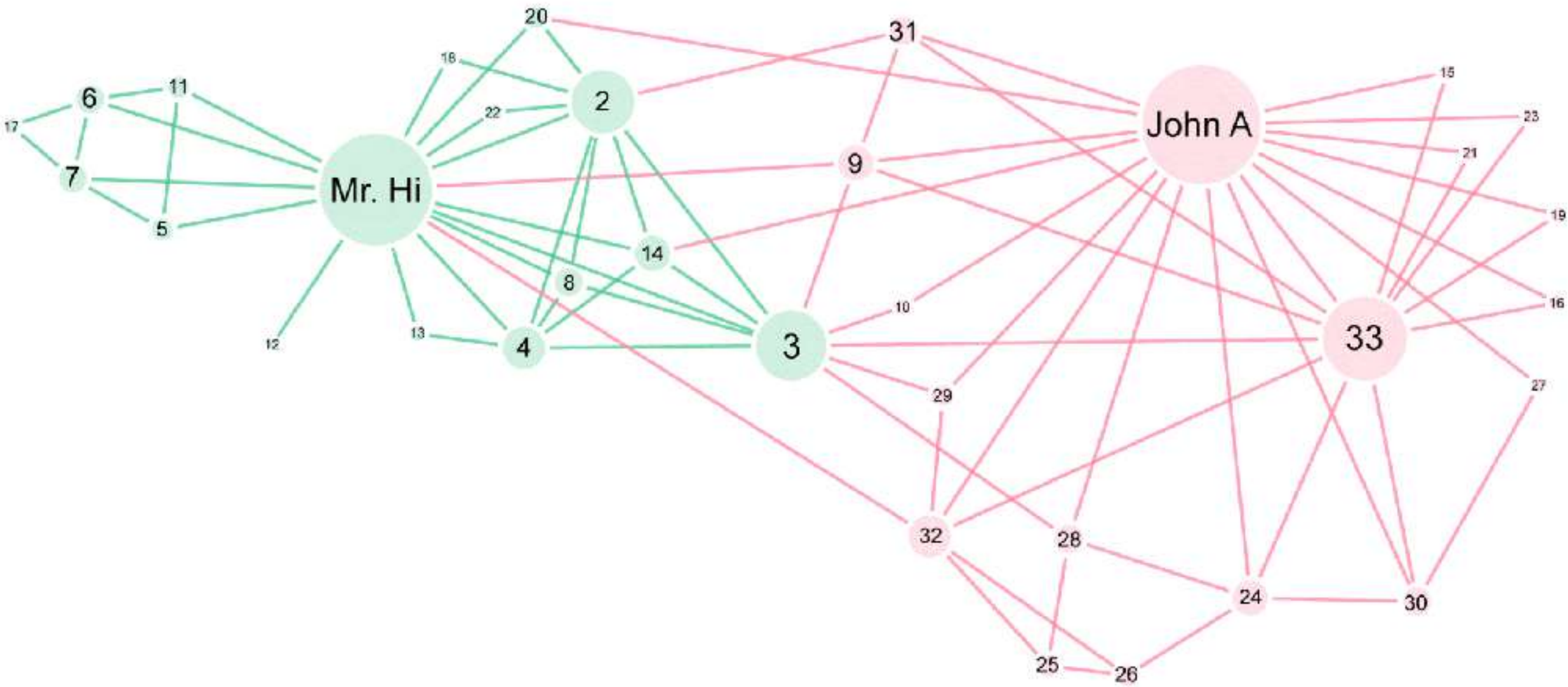
SEMI-FINALS

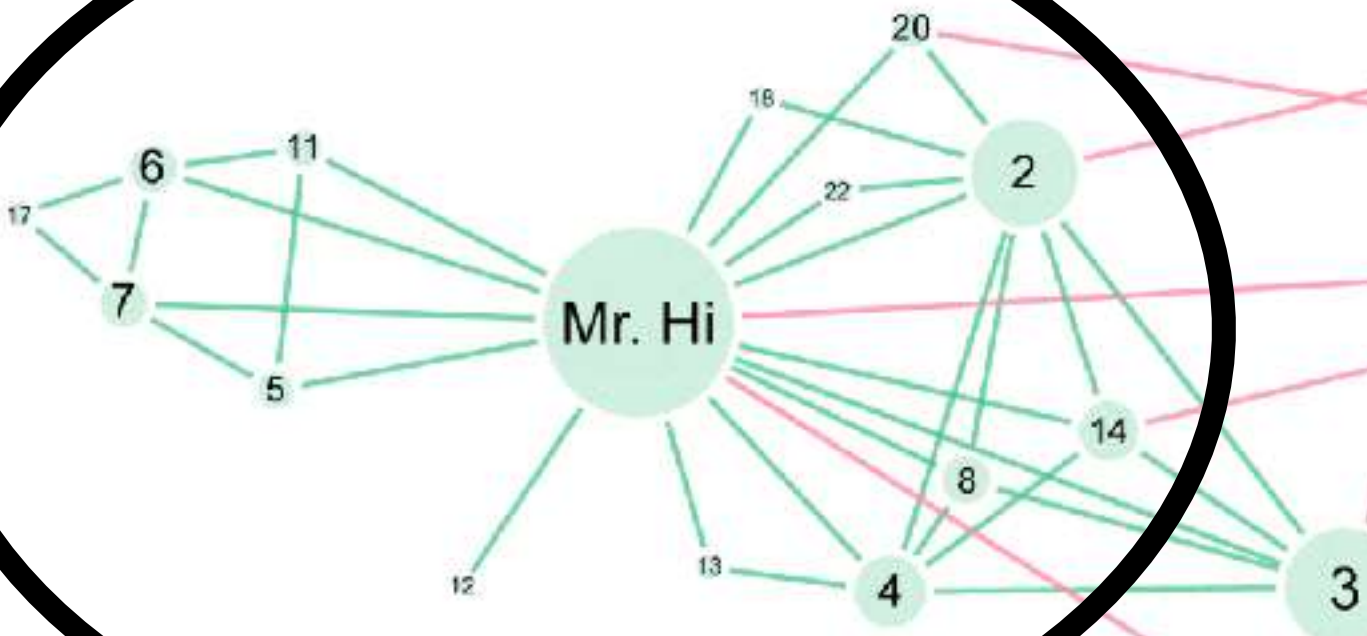
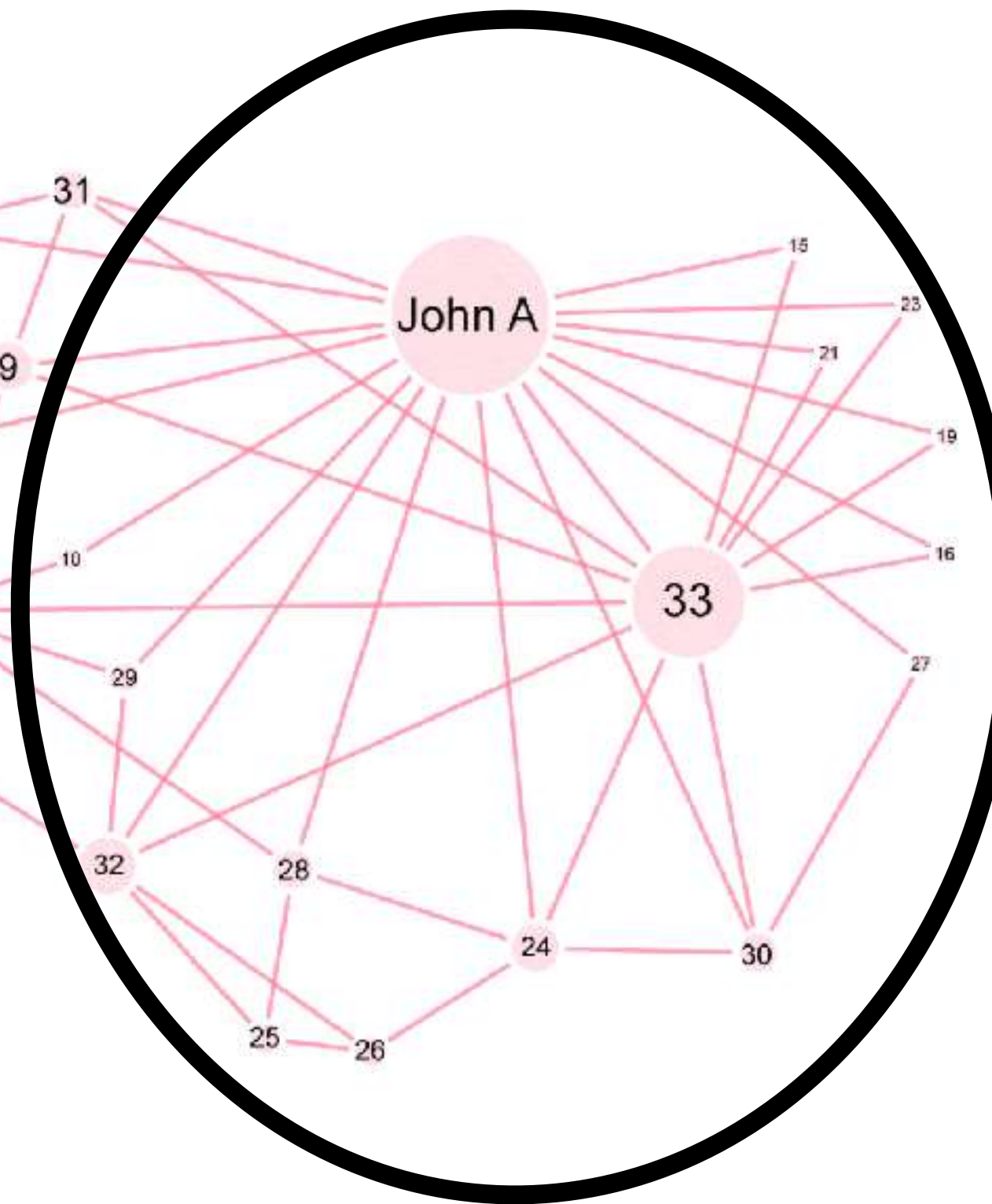
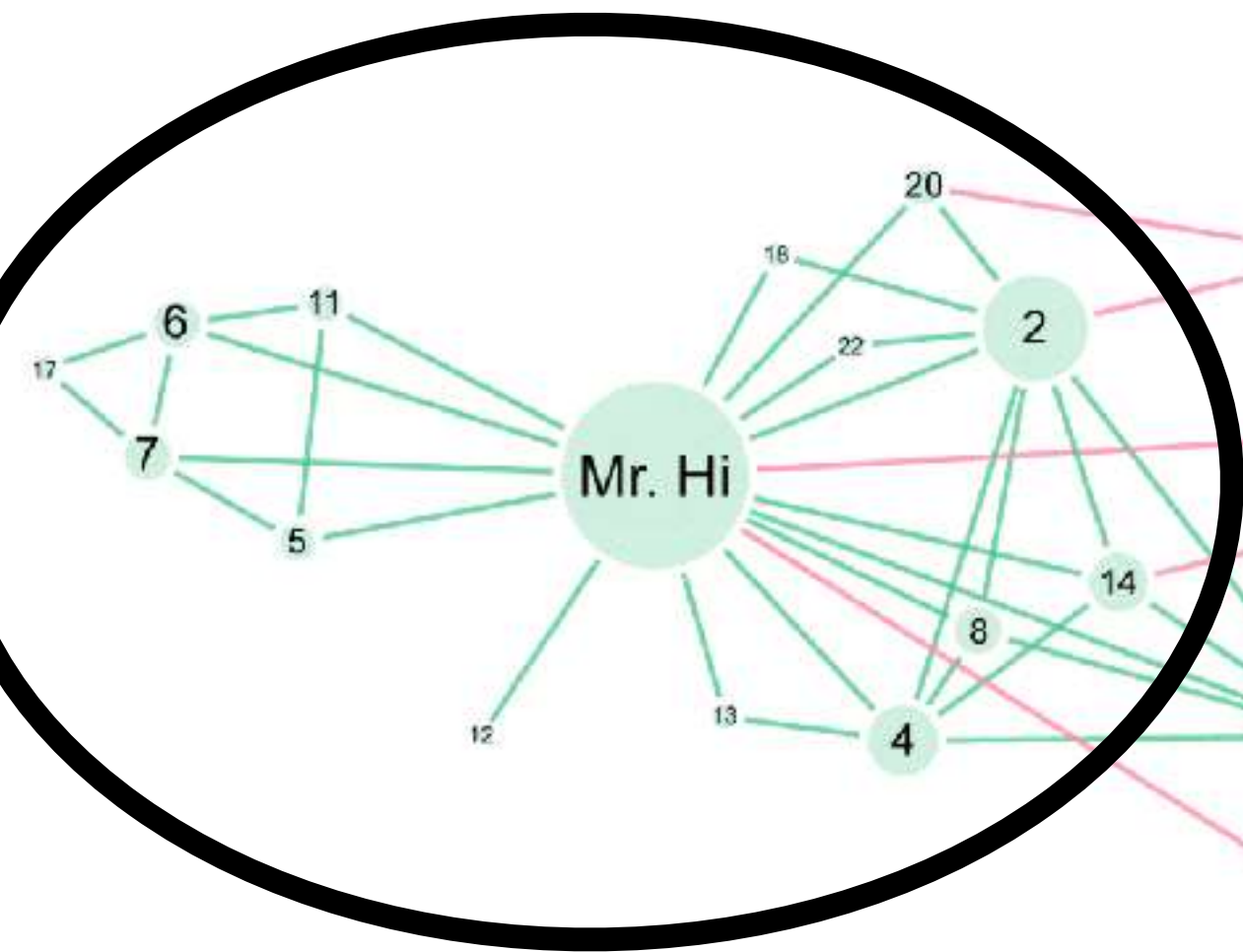
FINALS

ALL VALLEY  
CHAMPION



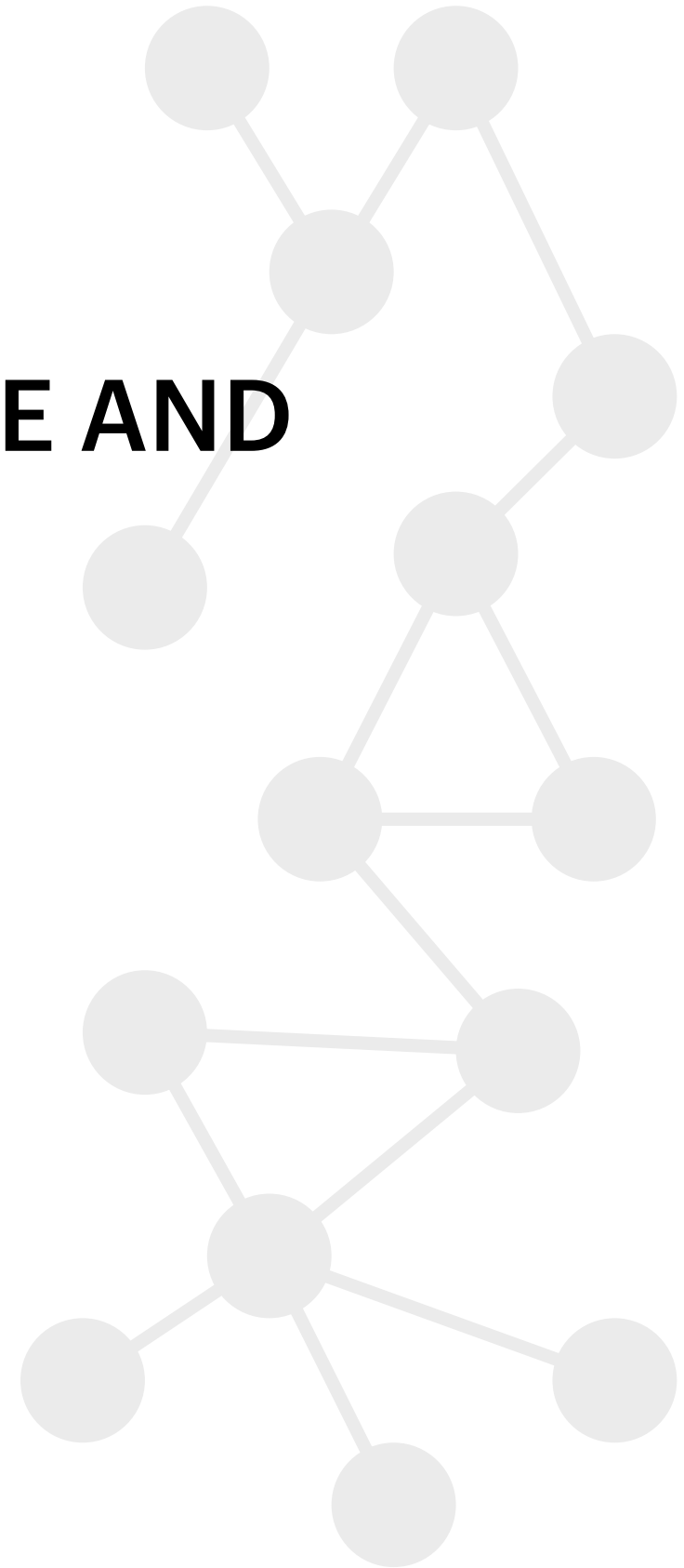






# DEFINITIONS

**INTERNAL AND EXTERNAL DEGREE:  
THE NUMBER OF NEIGHBOURS INSIDE AND  
OUTSIDE THE COMMUNITY**



# DEFINITIONS

**INTERNAL AND EXTERNAL DEGREE:  
THE NUMBER OF NEIGHBOURS INSIDE AND  
OUTSIDE THE COMMUNITY**

$$k_i = k_i^{int} + k_i^{ext}$$

**i is called internal node of community **c** if**

$$k_i^{ext} = 0 \text{ And } k_i^{int} > 0$$

**i is called boundary node of community **c** if**

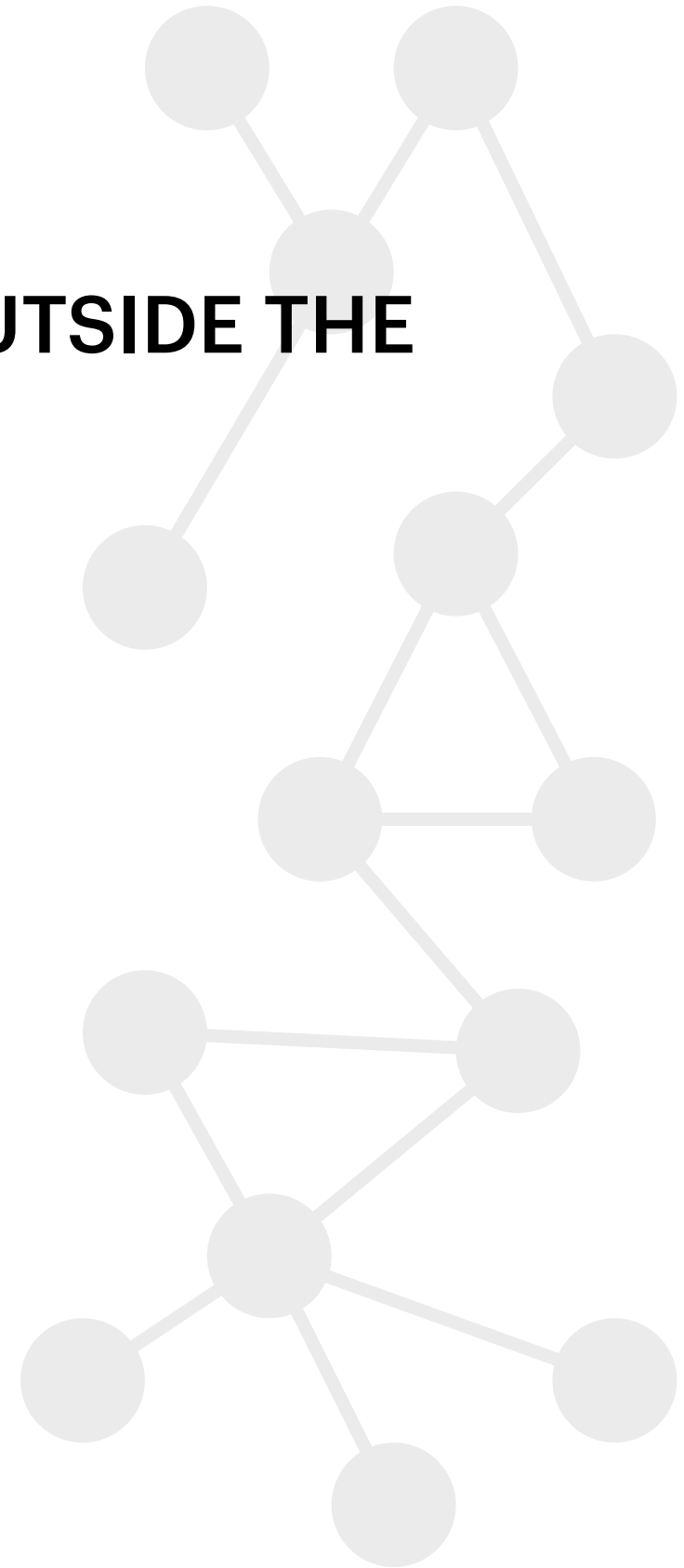
$$k_i^{ext} > 0 \text{ And } k_i^{int} > 0$$





# DEFINITIONS

**INTERNAL AND EXTERNAL DEGREE:**  
THE NUMBER OF NEIGHBOURS INSIDE AND OUTSIDE THE  
COMMUNITY



# DEFINITIONS

## **INTERNAL AND EXTERNAL DEGREE:**

**THE NUMBER OF NEIGHBOURS INSIDE AND OUTSIDE THE COMMUNITY**

## **NUMBER OF INTERNAL LINKS:**

**THE NUMBER OF LINKS BETWEEN NODES WITHIN THE COMMUNITY**



# DEFINITIONS

## **INTERNAL AND EXTERNAL DEGREE:**

**THE NUMBER OF NEIGHBOURS INSIDE AND OUTSIDE THE COMMUNITY**

## **NUMBER OF INTERNAL LINKS:**

**THE NUMBER OF LINKS BETWEEN NODES WITHIN THE COMMUNITY**

## **COMMUNITY DEGREE:**

**THE SUM OF DEGREE OF ALL THE NODES IN THE COMMUNITY**

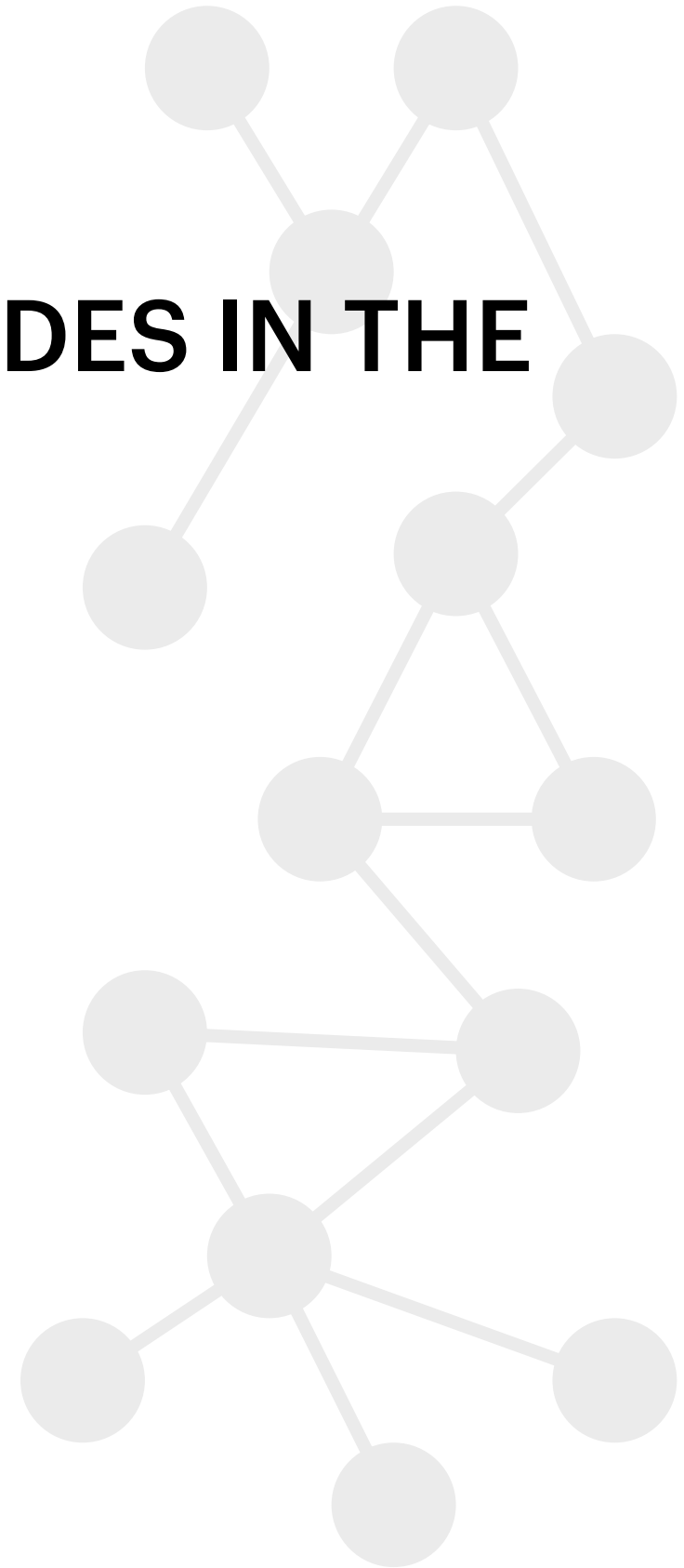


# DEFINITIONS

**COMMUNITY DEGREE:**

**THE SUM OF DEGREE OF ALL THE NODES IN THE  
COMMUNITY**

$$k_C = \sum_{i \in C} k_i$$





# DEFINITIONS

## **INTERNAL AND EXTERNAL DEGREE:**

THE NUMBER OF NEIGHBOURS INSIDE AND OUTSIDE THE COMMUNITY

## **NUMBER OF INTERNAL LINKS:**

THE NUMBER OF LINKS BETWEEN NODES WITHIN THE COMMUNITY

## **COMMUNITY DEGREE:**

THE SUM OF DEGREE OF ALL THE NODES IN THE COMMUNITY

## **INTERNAL LINK DENSITY:**

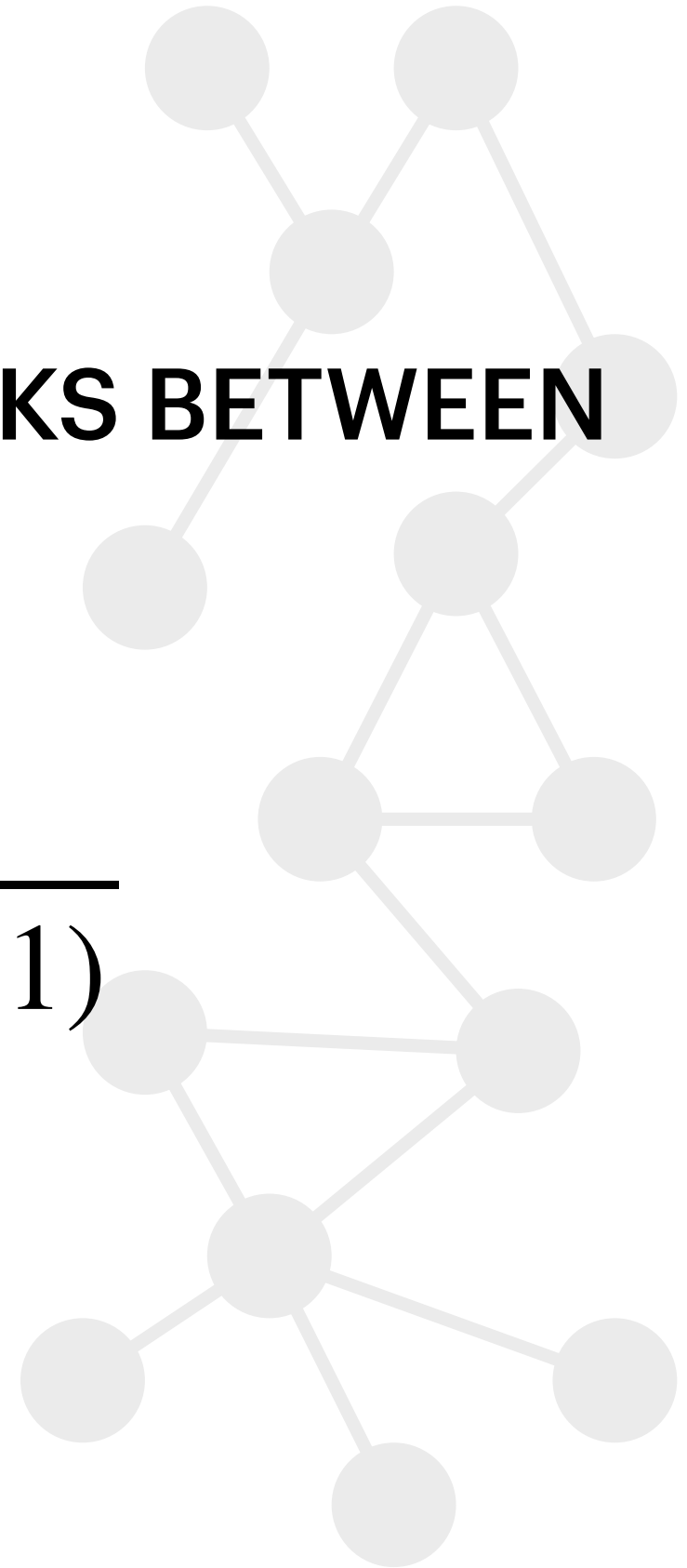
DENSITY THAT CONSIDERS ONLY LINKS BETWEEN MEMBERS OF THE COMMUNITY



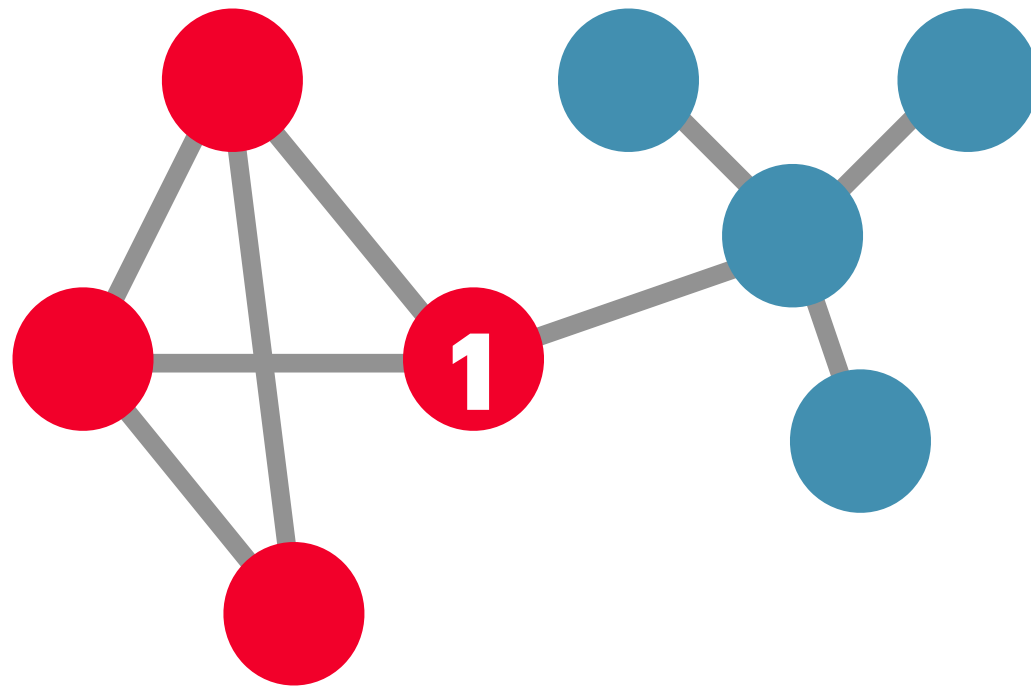
# DEFINITIONS

**INTERNAL LINK DENSITY:**  
DENSITY THAT CONSIDERS ONLY LINKS BETWEEN  
MEMBERS OF THE COMMUNITY

$$\delta_C^{int} = \frac{L_C}{L_C^{max}} = \frac{2L_C}{N_C(N_C - 1)}$$

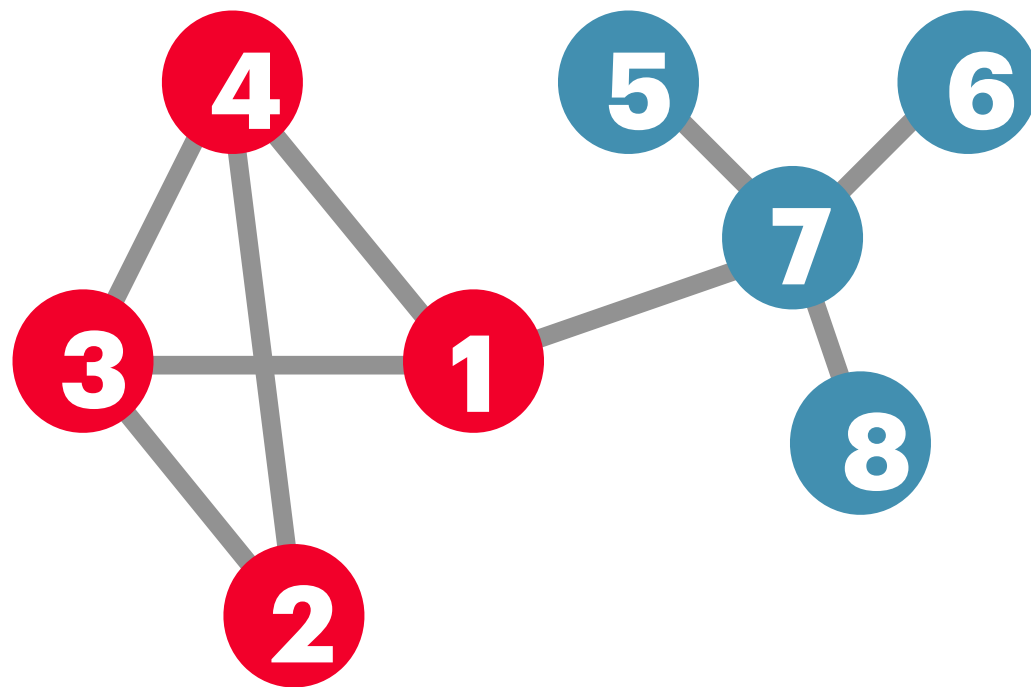


# DEFINITIONS



?  $k_1^{ext}$ ,  $k_1^{int}$ ,  $\delta_{red}^{int}$ ,  $k_{blue}$

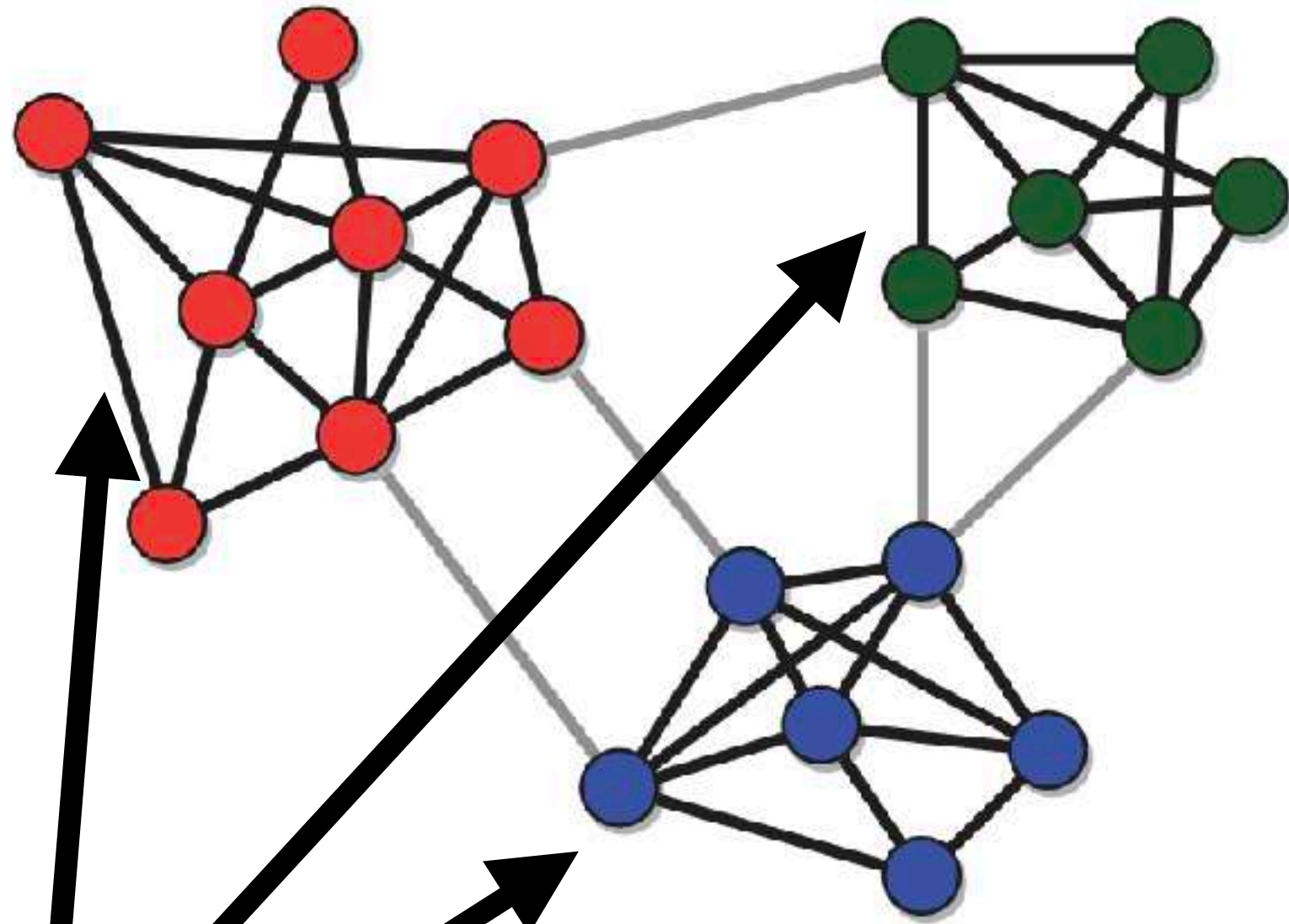
# DEFINITIONS



**Which are the boundary nodes?**

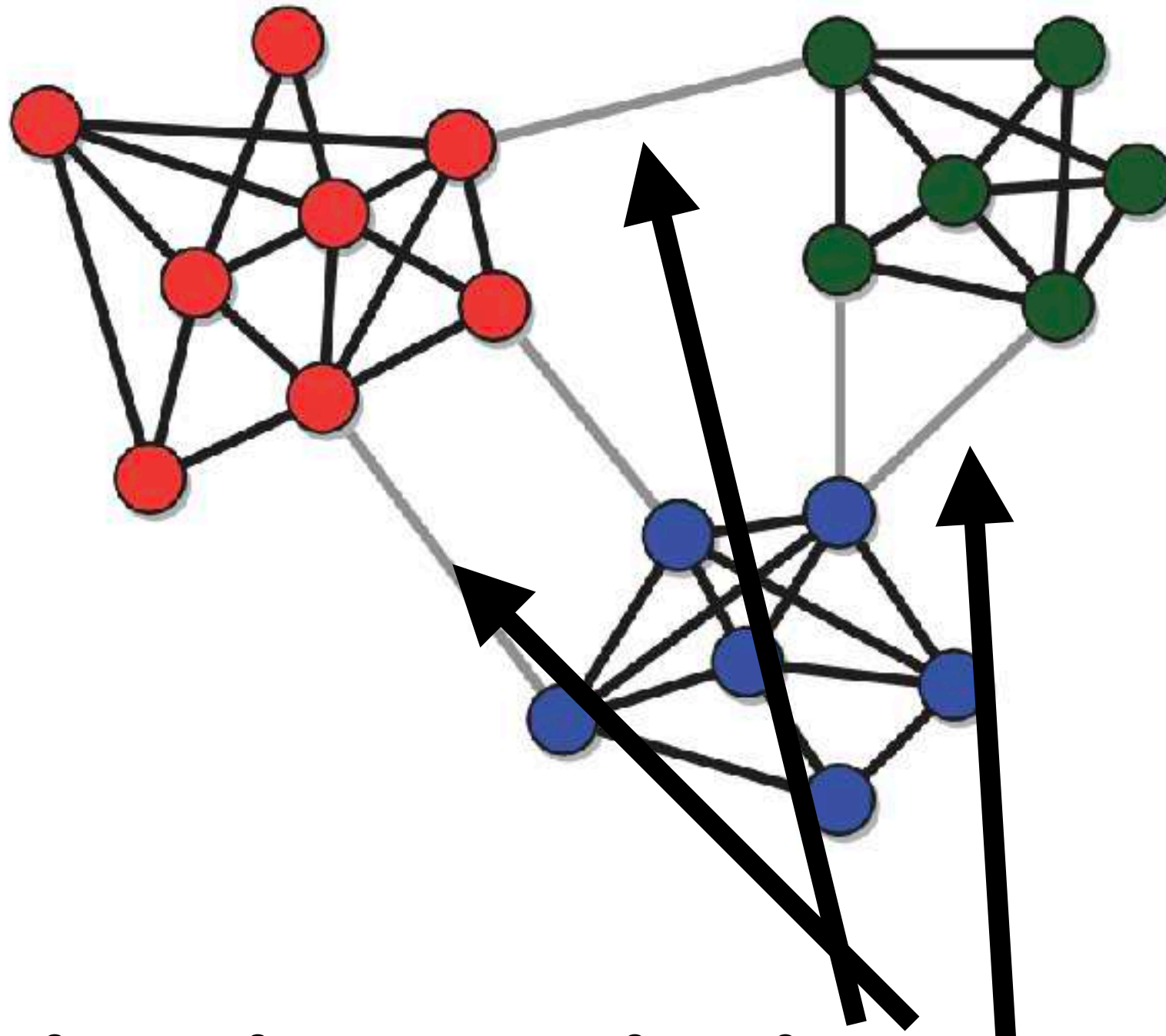


# DEFINITIONS



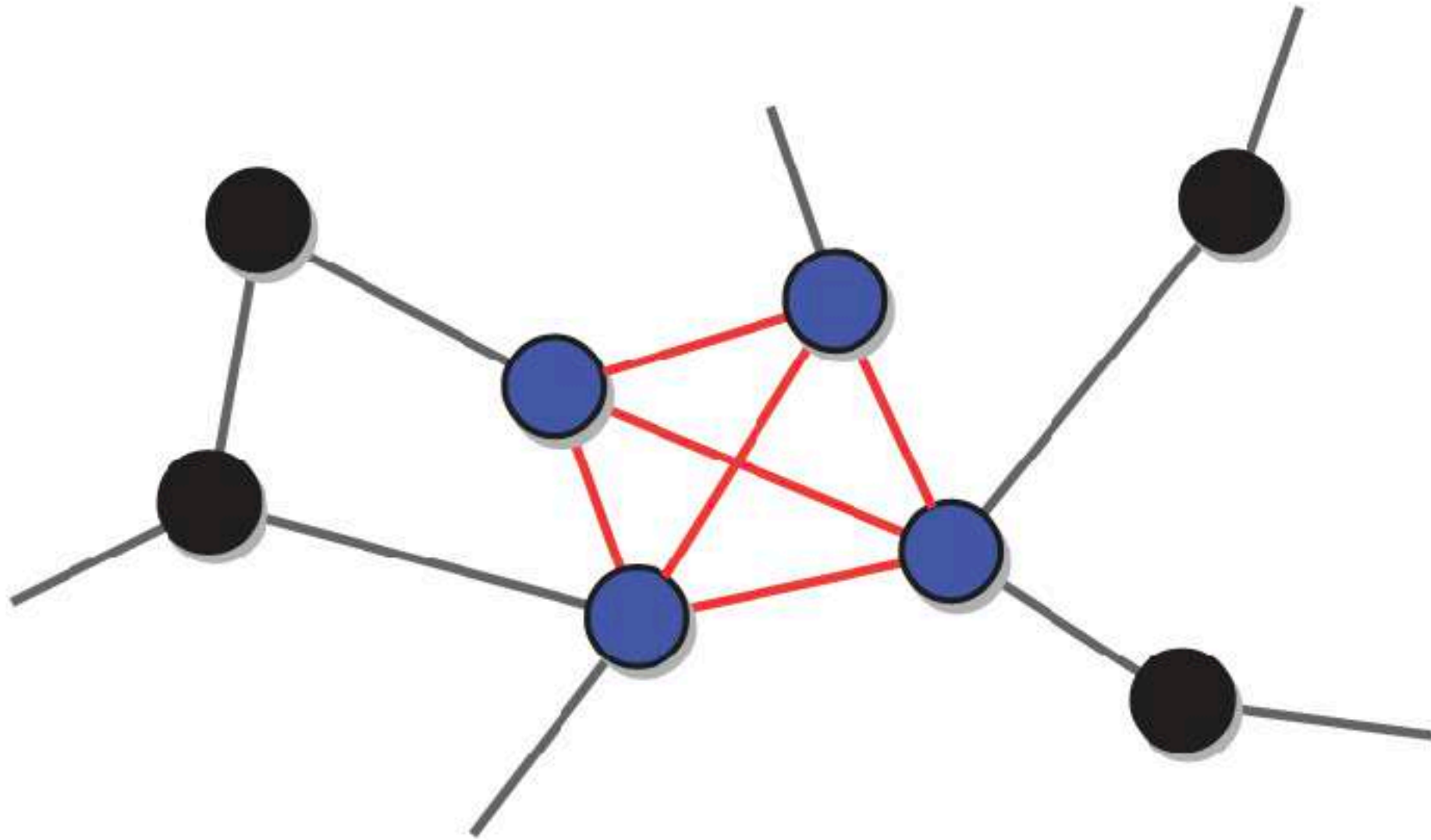
**high cohesion, high separation**

# DEFINITIONS



**high cohesion, high separation**

# DEFINITIONS

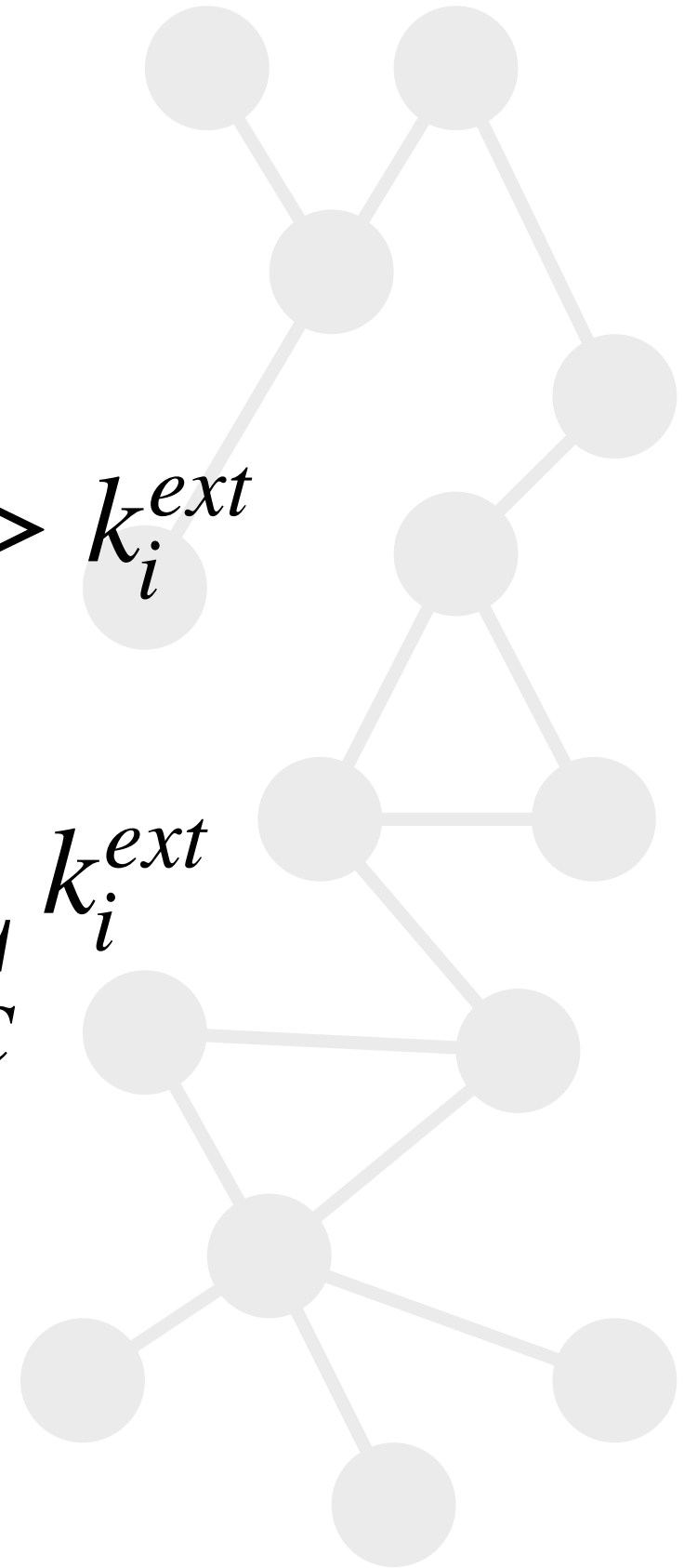


**clique** (a fully connected subgraph)

# DEFINITIONS

**Strong community:**  $\forall i \in C : k_i^{int} > k_i^{ext}$

**Weak community:**  $\sum_{i \in C} k_i^{int} > \sum_{i \in C} k_i^{ext}$

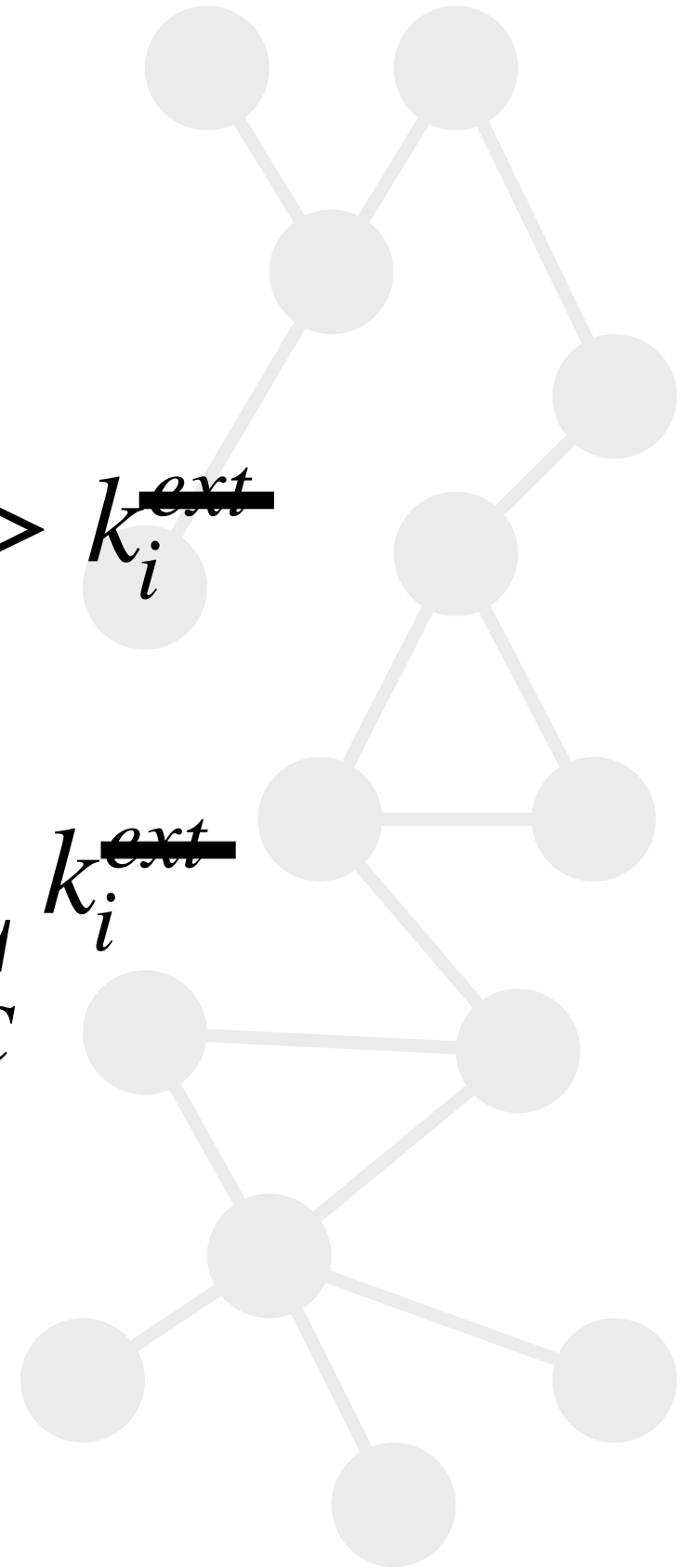


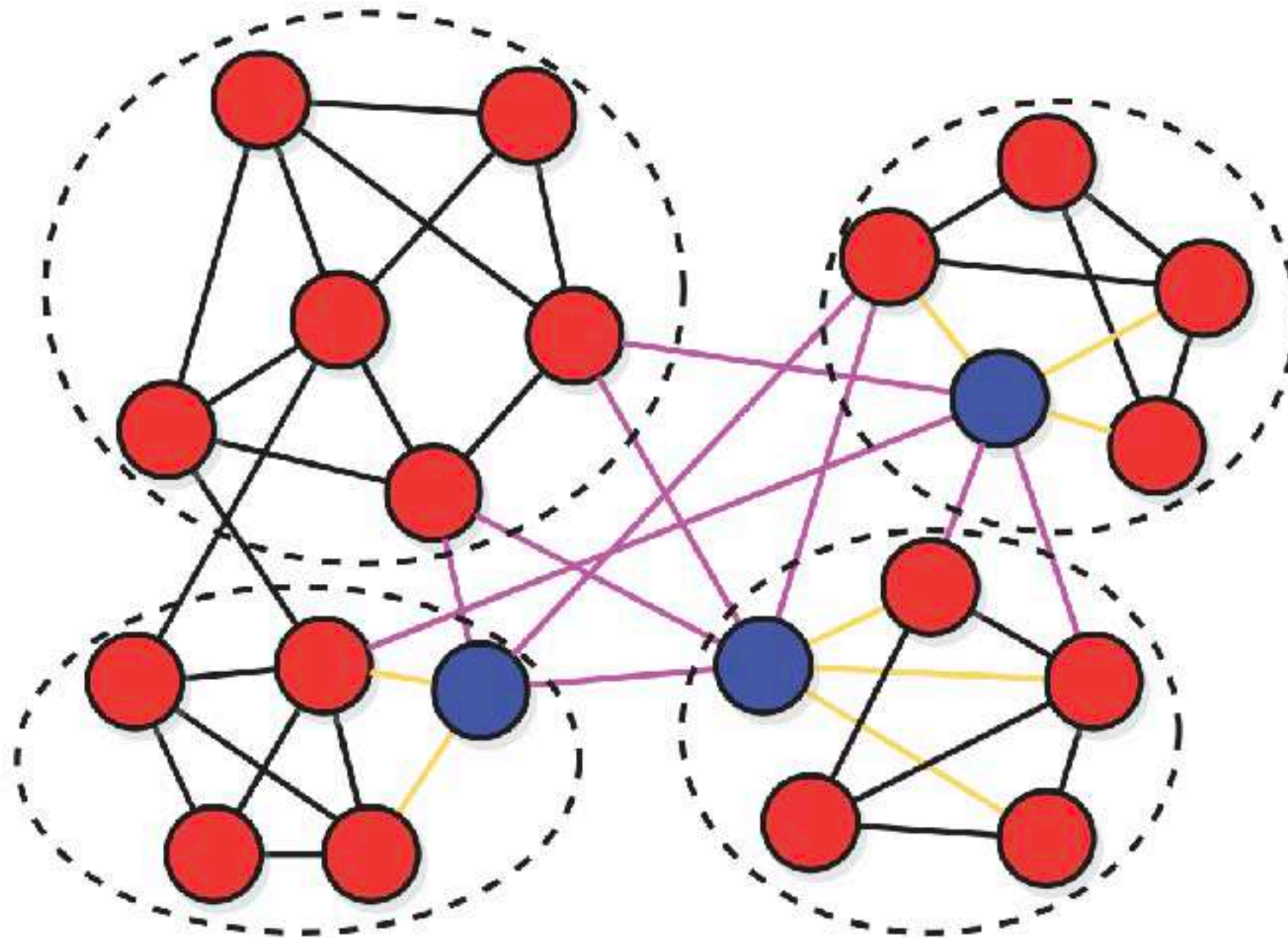


# DEFINITIONS

**Strong community:**  $\forall i \in C : k_i^{int} > k_i^{ext}$

**Weak community:**  $\sum_{i \in C} k_i^{int} > \sum_{i \in C} k_i^{ext}$





Strong and weak communities. The four subnetworks enclosed in the dashed contours are weak communities according to both definitions we have given. They are also strong communities according to the less stringent definition, as the internal degree of each node exceeds the number of links joining the node with those of every other community. However, three of the subnetworks are not strong communities in the more stringent sense, because some nodes (in blue) have external degree larger than their internal degree (the internal and external links of these nodes are colored in yellow and magenta, respectively). Adapted from Fortunato and Hric (2016).

# PARTITIONS

**A PARTITION IS A DIVISION OF THE NETWORK IN COMMUNITIES**



# PARTITIONS

**SUPPOSE YOU HAVE A NETWORK  $G$  WITH 10 NODES  
1,2,...,10**

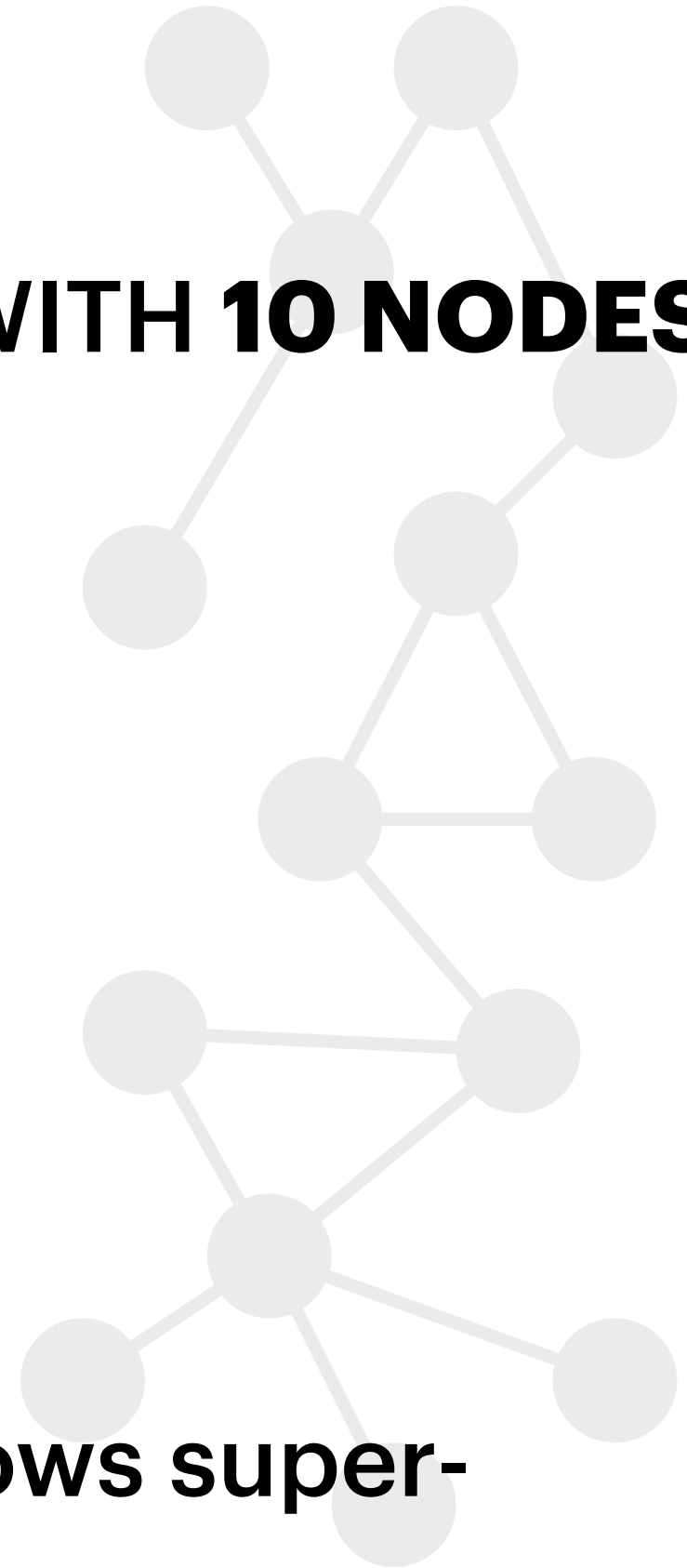
**{1,2,...,10}**

**{1} {2} {3} ... {10}**

**{1,2} {3,6,9} {5,8,10} {7,4}**

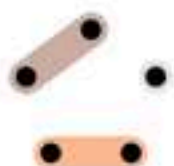
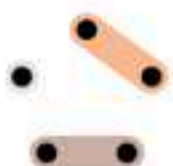
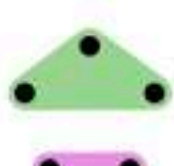
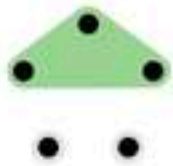
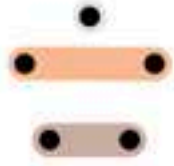
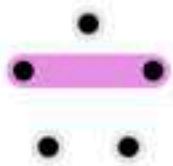
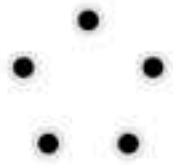
**THESE ARE ALL VALID PARTITIONS**

**The number of possible partitions grows super-exponentially**

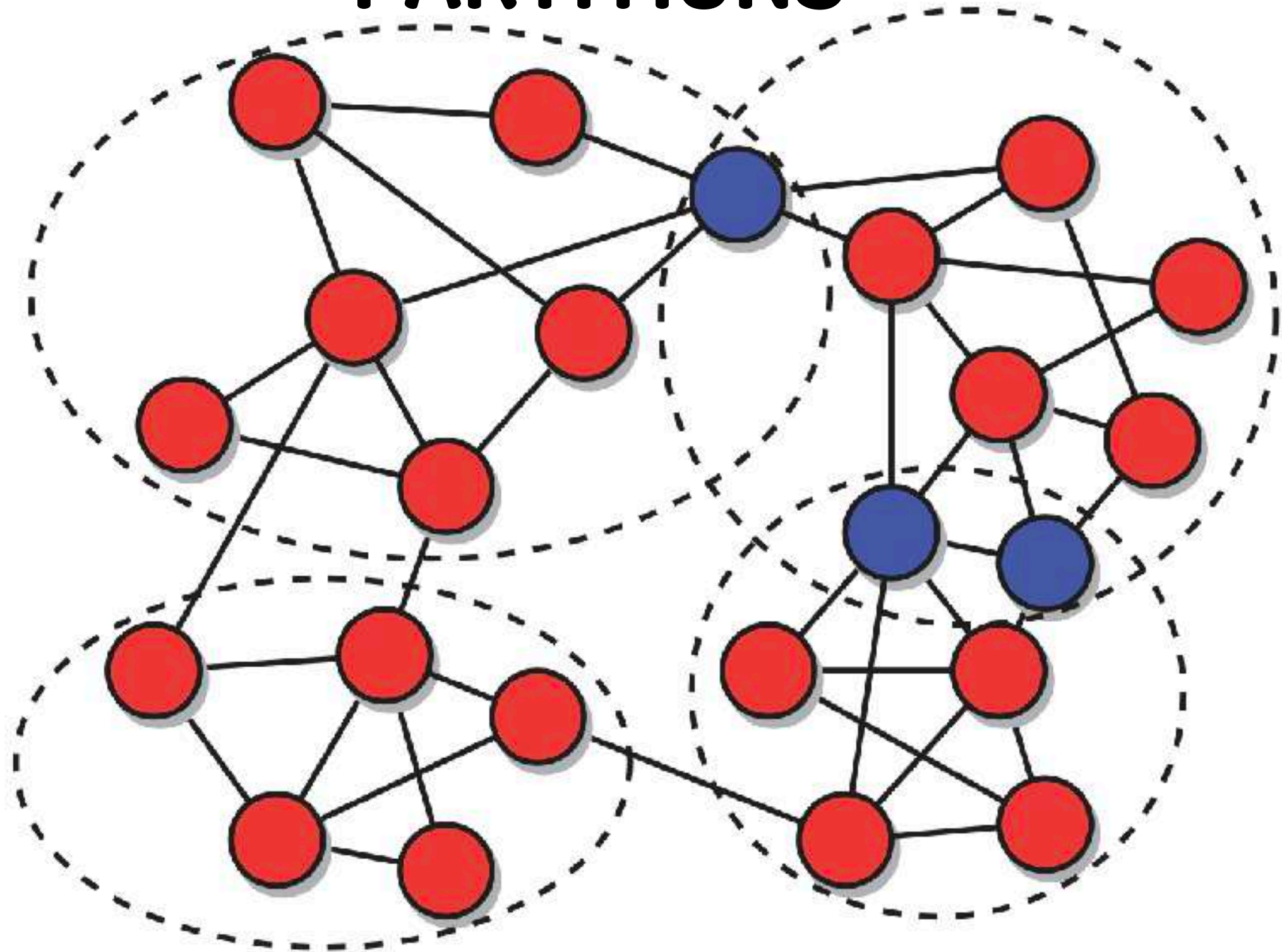




# PARTITIONS



# PARTITIONS

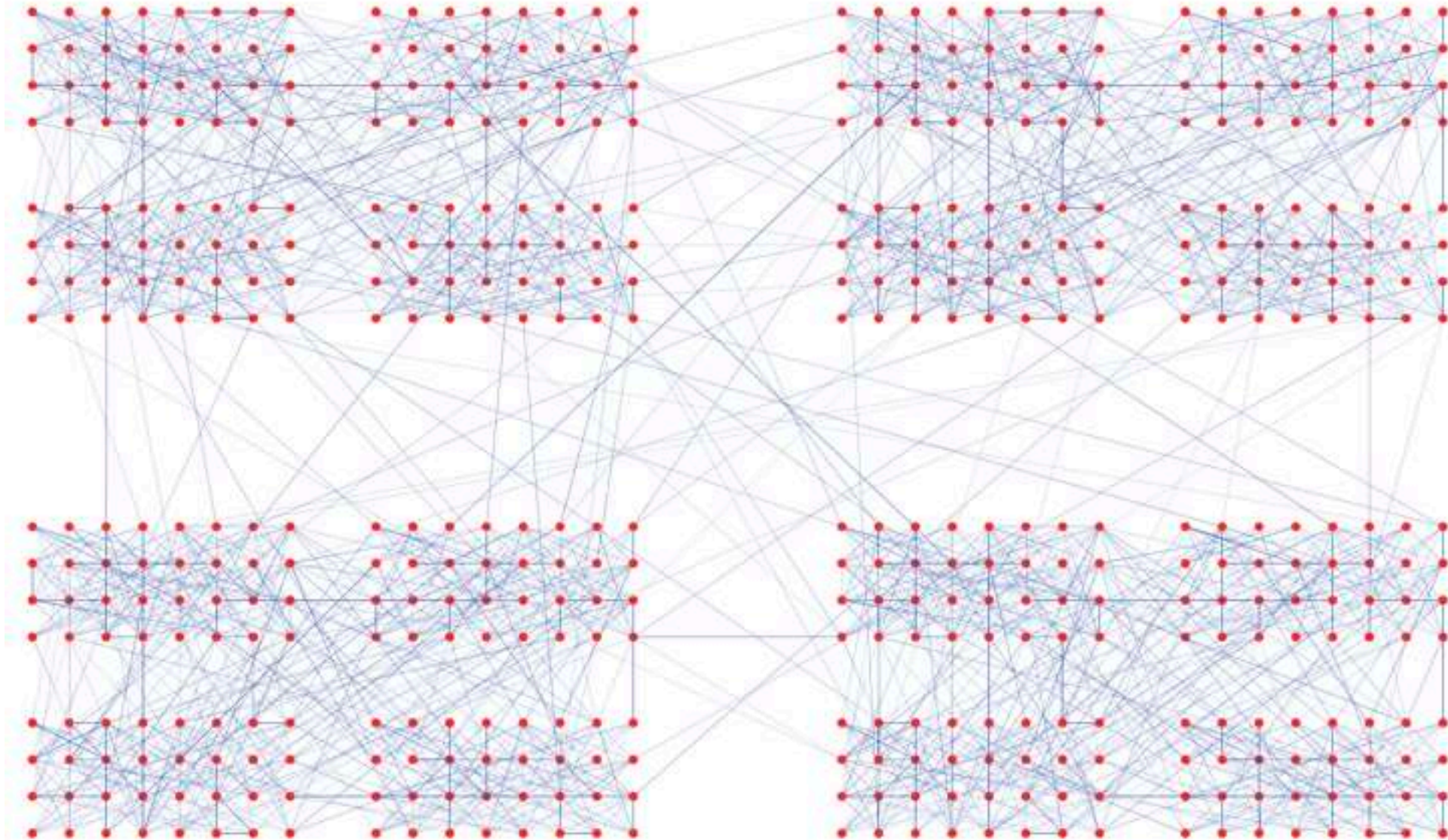


**COMMUNITIES CAN OVERLAP**

(You are part of different communities, think about it)



# PARTITIONS

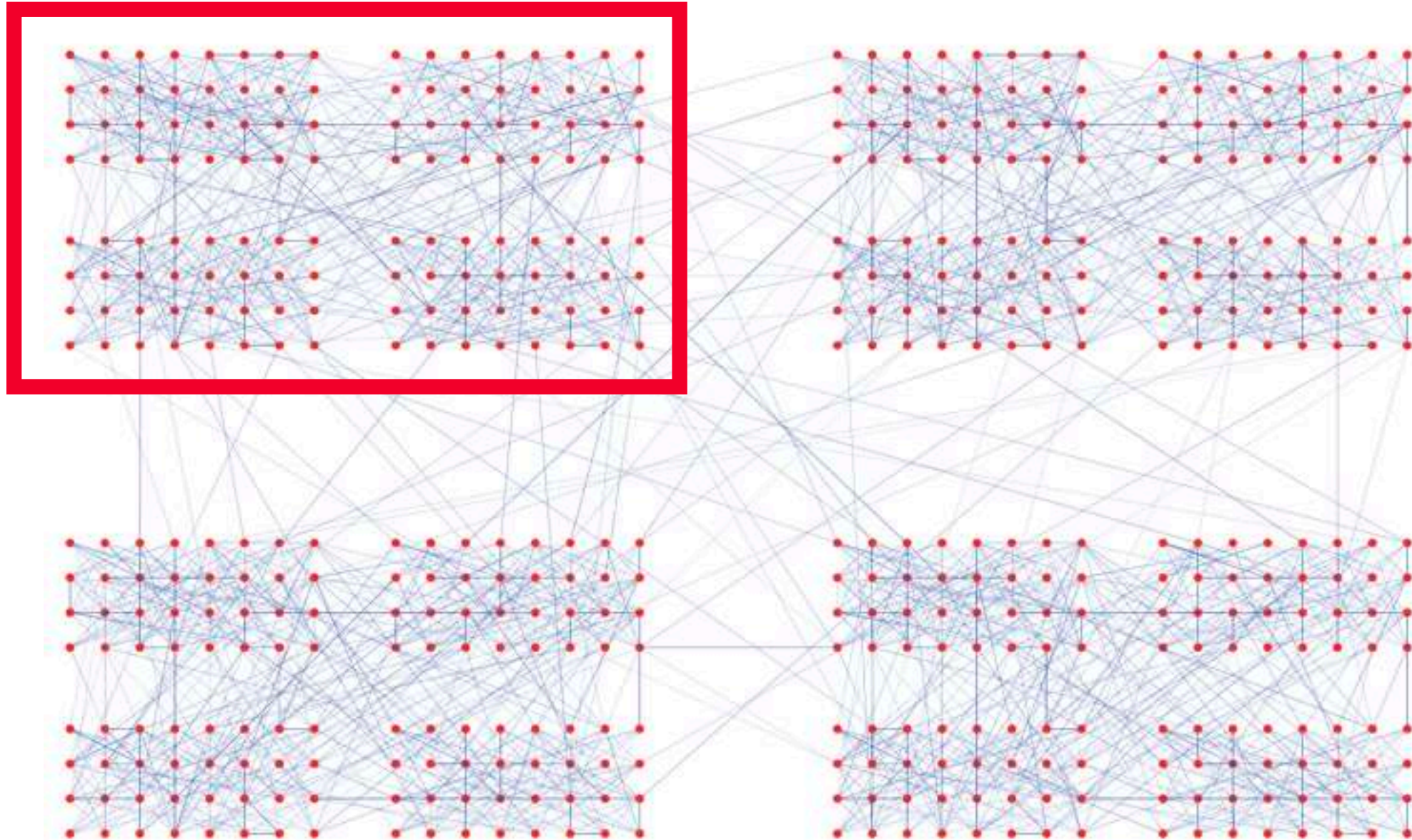


**COMMUNITIES CAN BE HIERARCHICAL**

**(There might be communities within communities)**



# PARTITIONS

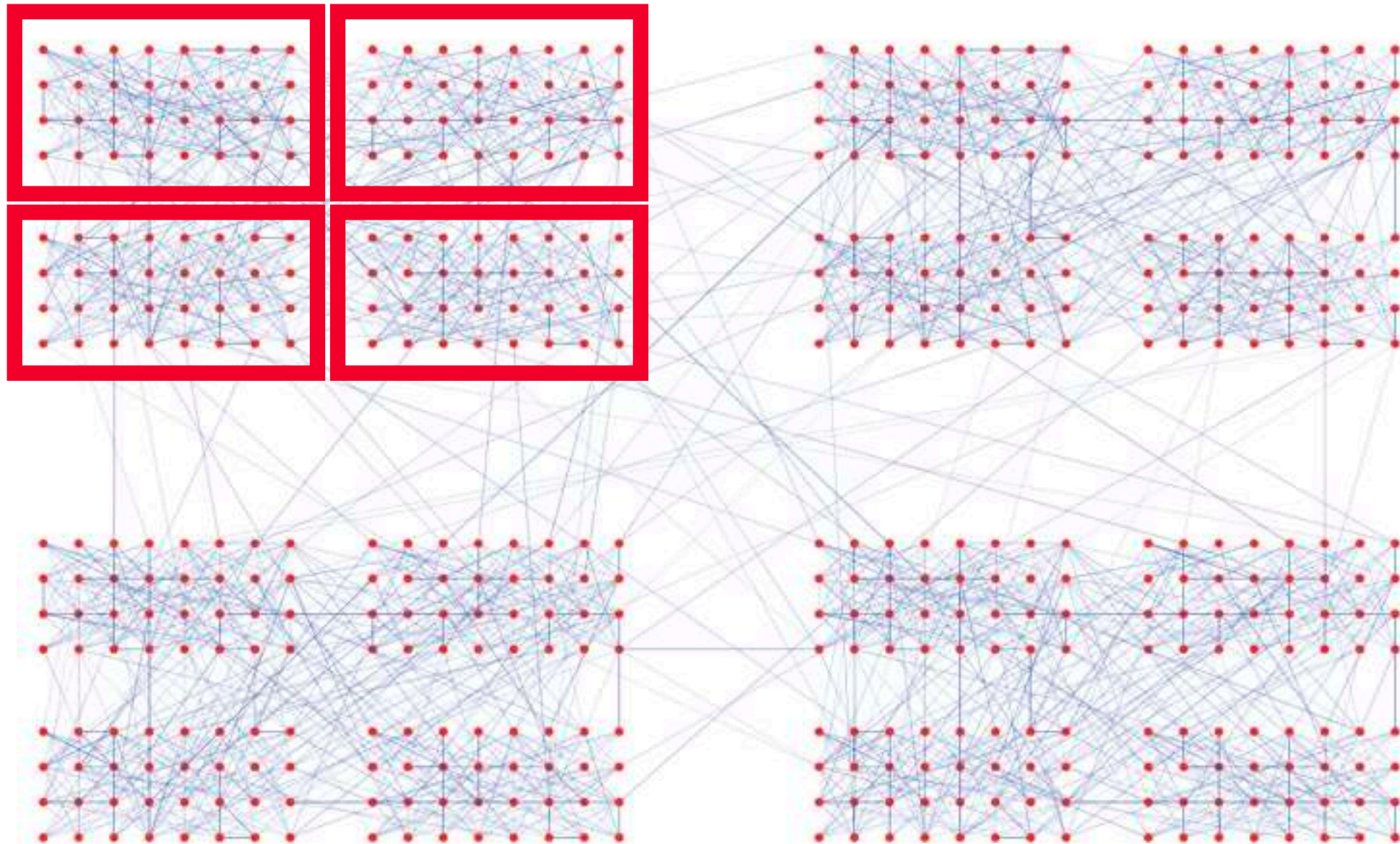


**COMMUNITIES CAN BE HIERARCHICAL**

**(There might be communities within communities)**



# PARTITIONS

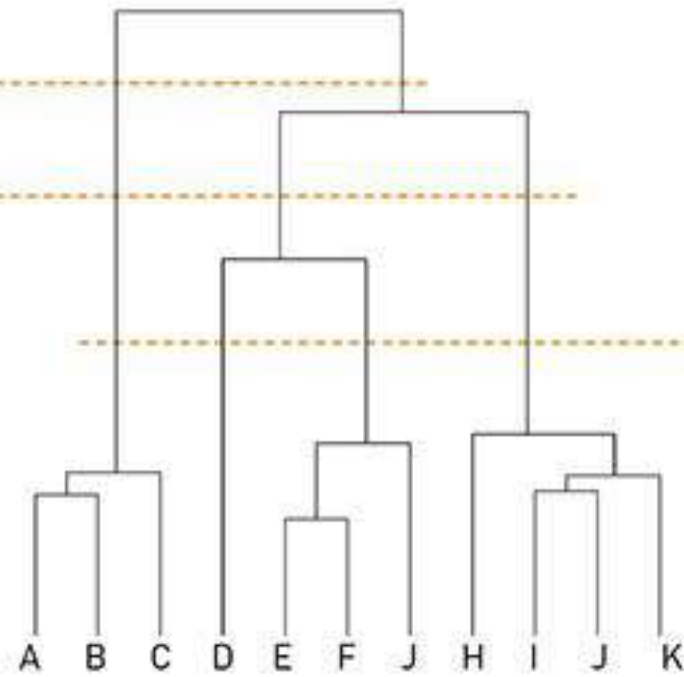


**COMMUNITIES CAN BE HIERARCHICAL**

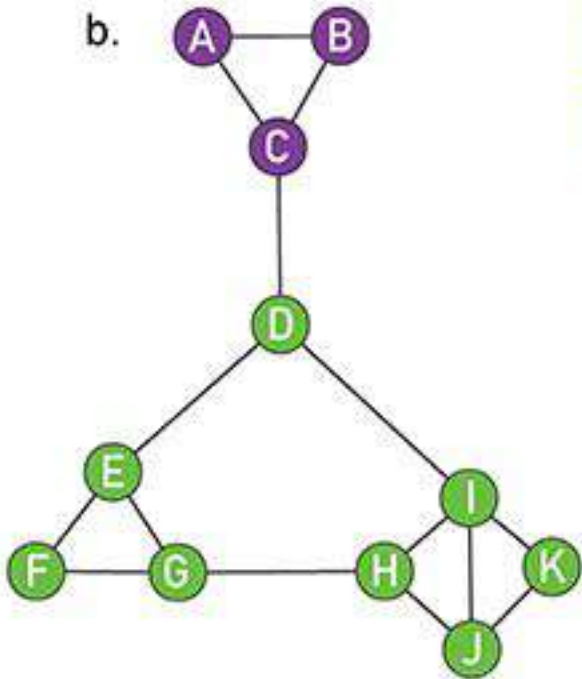
**(There might be communities within communities)**

# PARTITIONS

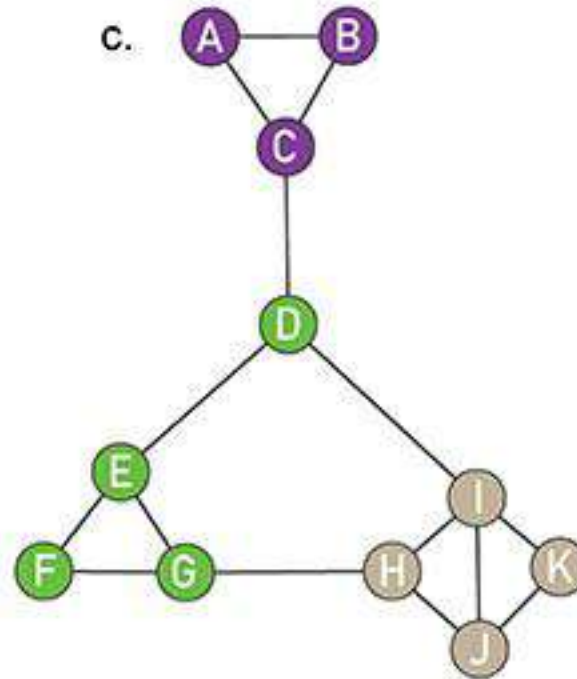
a.



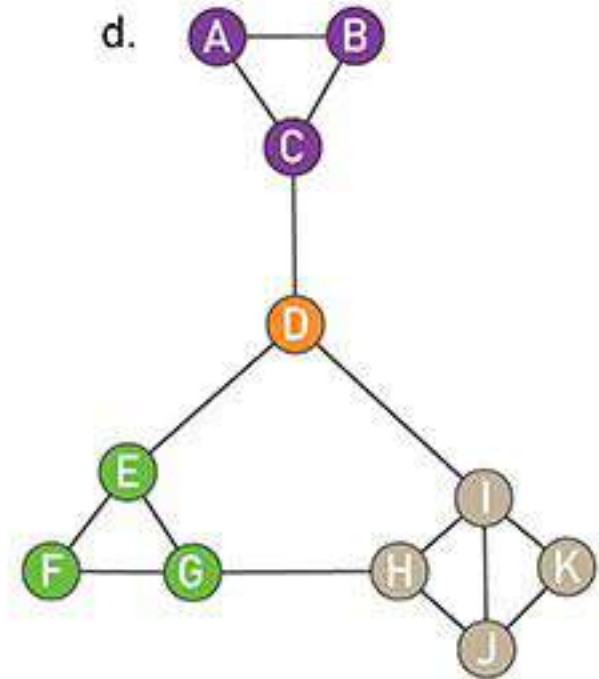
b.



c.



d.



# Dendrogram





# **EXERCISE**

**MAKE SOME EXAMPLES  
OF SOCIAL AND FINANCIAL  
NETWORKS WITH COMMUNITIES**

# PART I RECAP

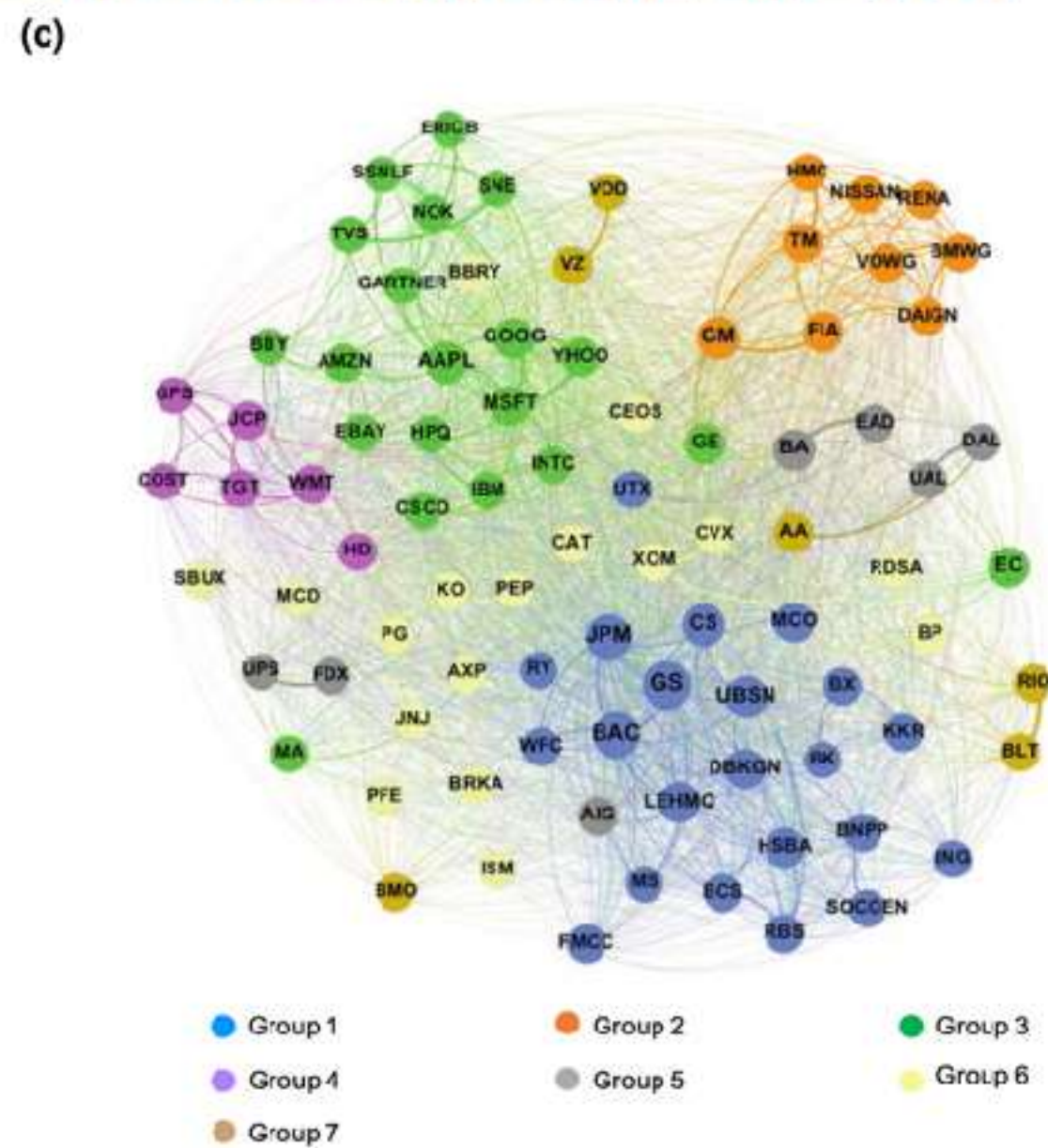
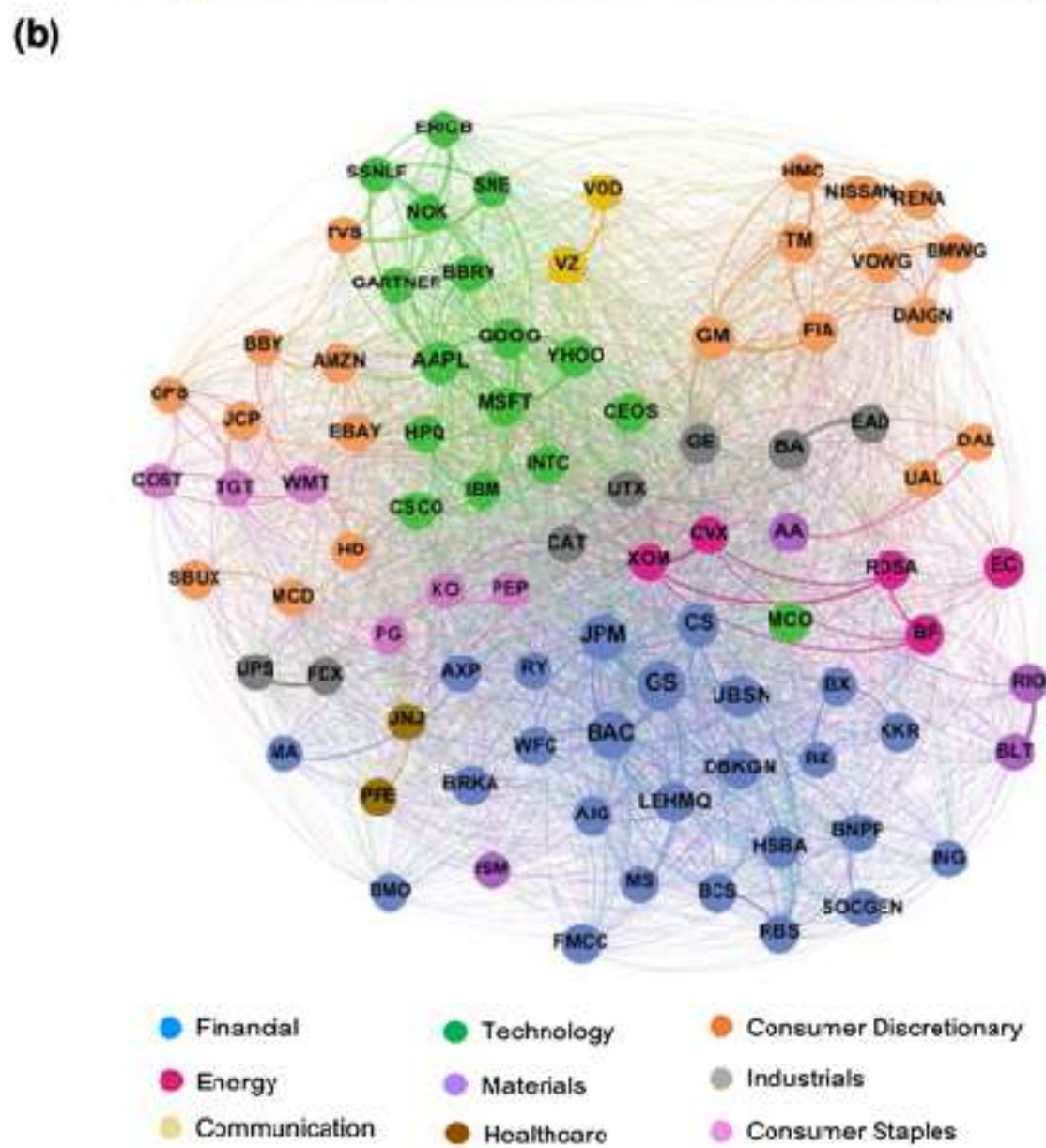
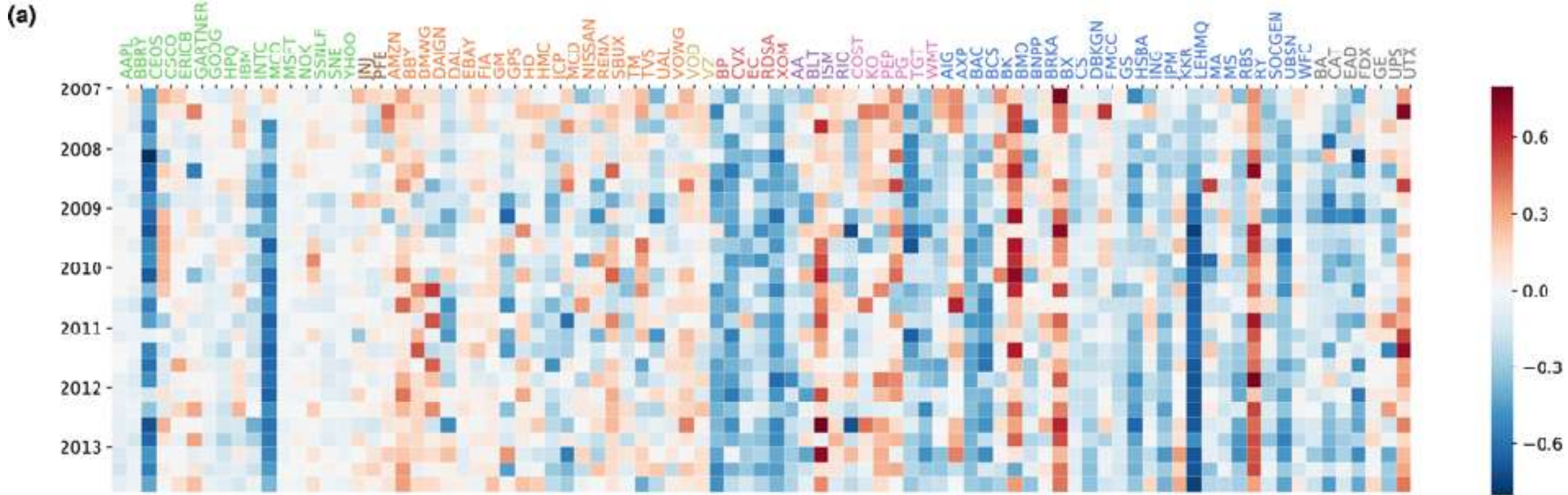
We saw what communities are and how they are **defined**

We explored some **examples**

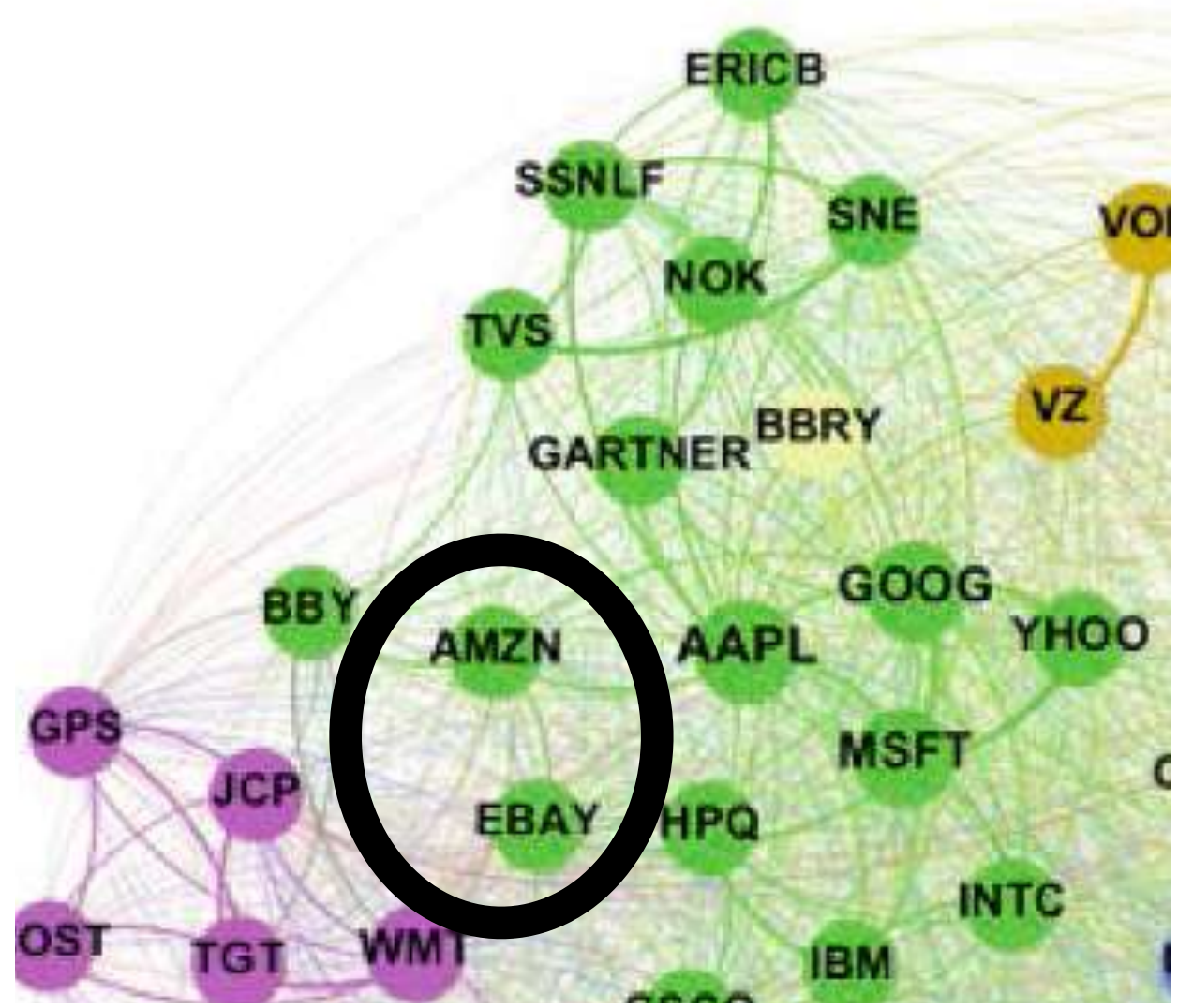
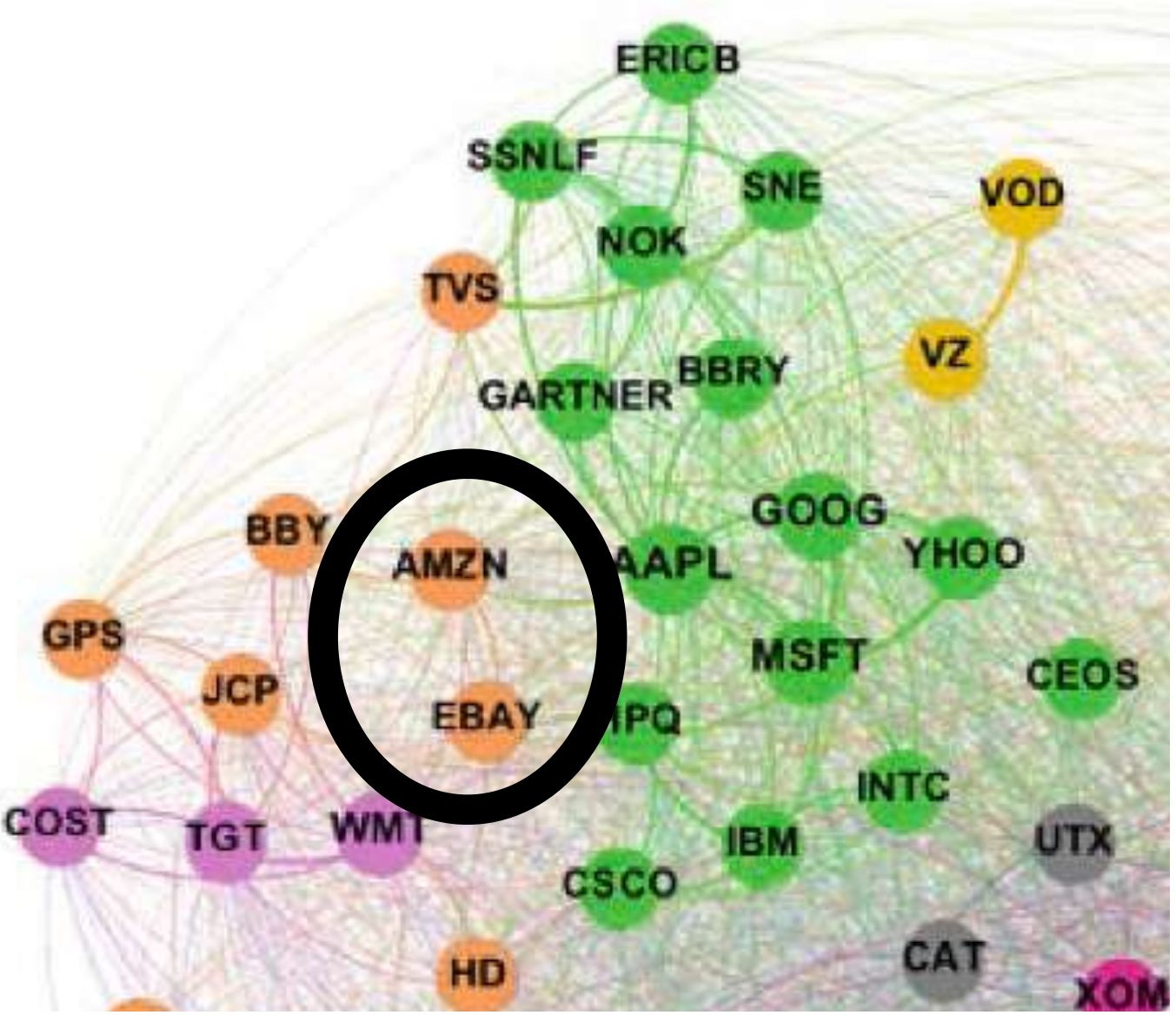
We now have all the tools to learn about **community detection**











- |   |   |   |
|---|---|---|
| <span style="color: blue;">●</span> Financial       | <span style="color: green;">●</span> Technology | <span style="color: orange;">●</span> <u>Consumer Discretionary</u> |
| <span style="color: pink;">●</span> Energy          | <span style="color: purple;">●</span> Materials | <span style="color: grey;">●</span> Industrials                     |
| <span style="color: yellow;">●</span> Communication | <span style="color: brown;">●</span> Healthcare | <span style="color: magenta;">●</span> Consumer Staples             |

# COMMUNITY DETECTION

The task of **finding communities** in a network  
We now have all the tools to learn about **community detection**



# COMMUNITY DETECTION

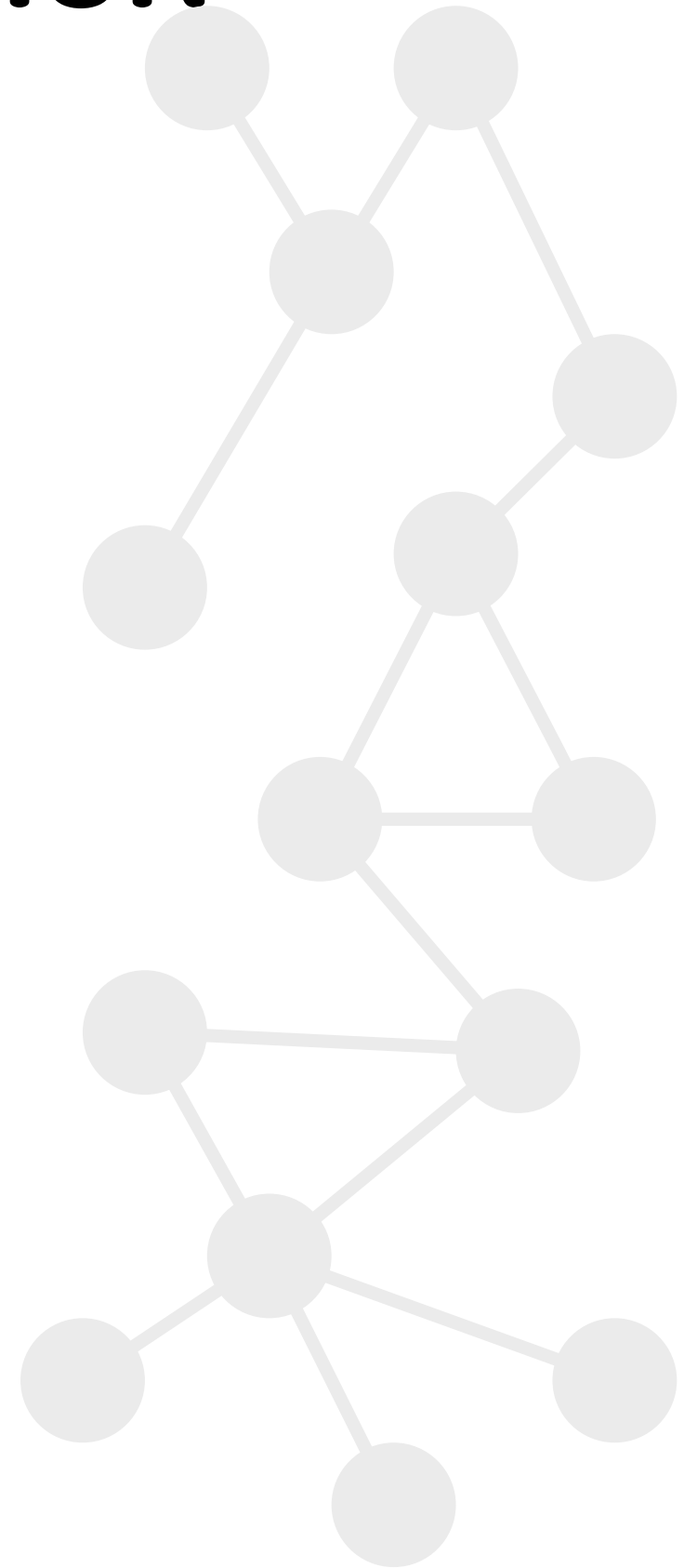
## FOUR APPROACHES

Bridge removal

Modularity maximisation

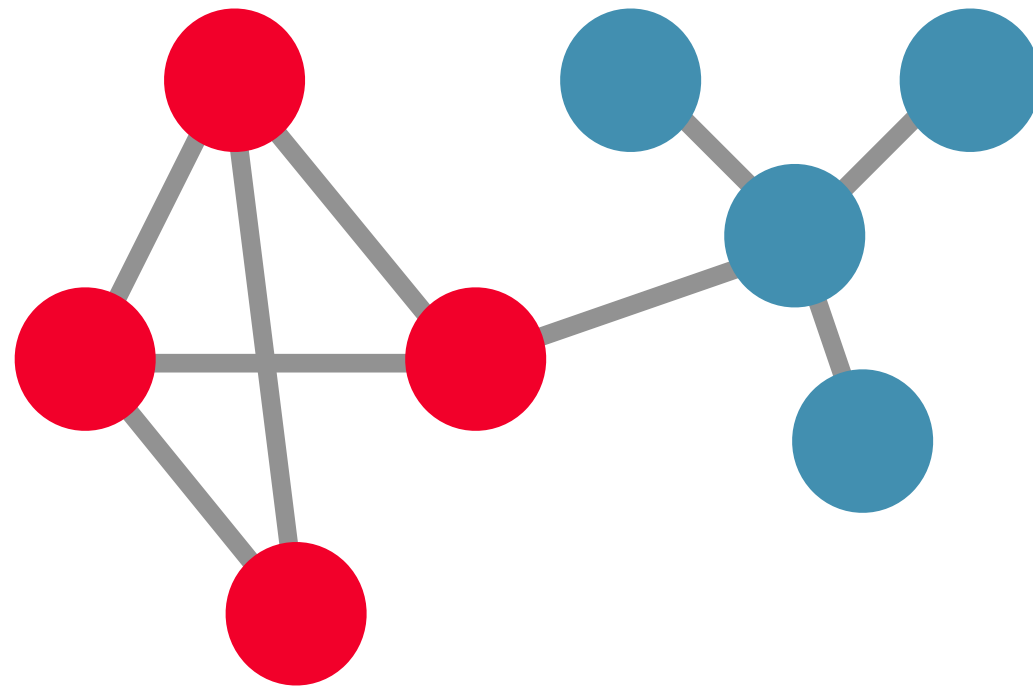
Label propagation

Stochastic block modelling





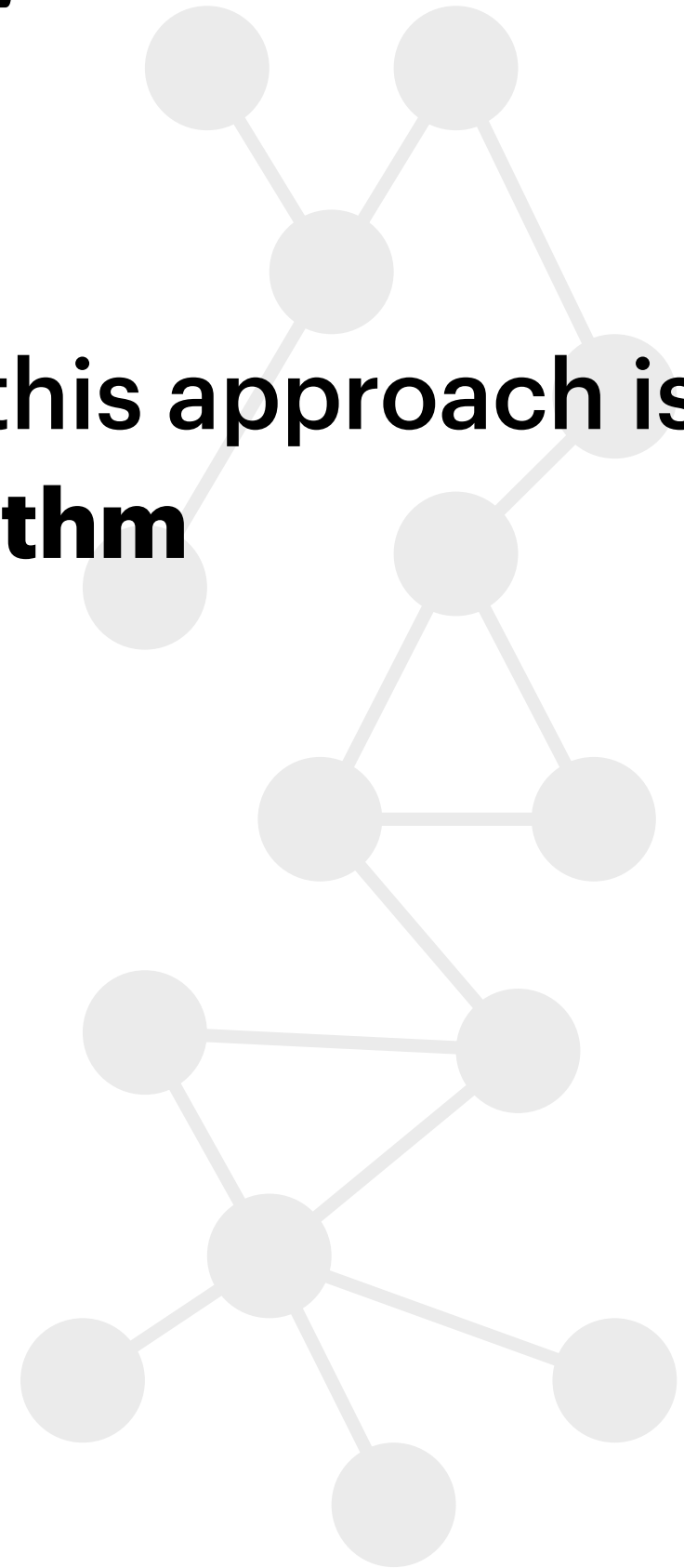
# BRIDGE REMOVAL



**A bridge is a link whose removal breaks the network into two parts**

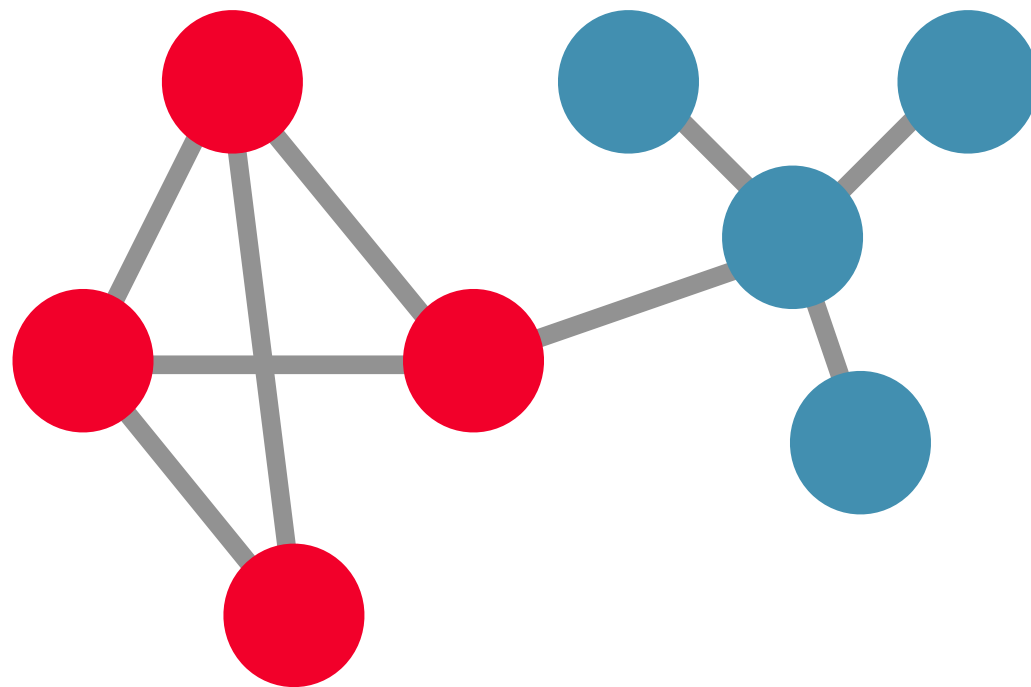
# BRIDGE REMOVAL

The most famous algorithm based on this approach is the **Girvan-Newman algorithm**

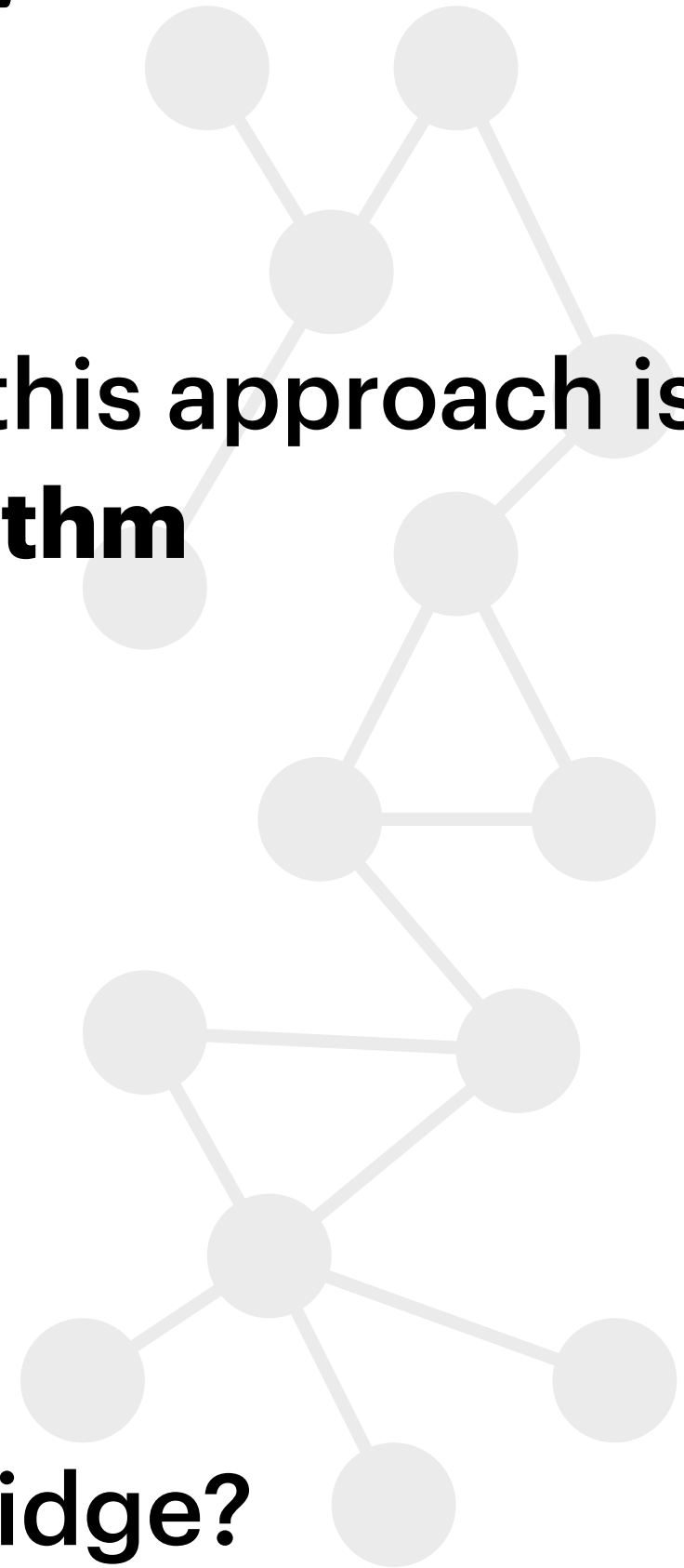


# BRIDGE REMOVAL

The most famous algorithm based on this approach is the **Girvan-Newman algorithm**



How do we find a bridge?



# BRIDGE REMOVAL

The most famous algorithm based on this approach is the **Girvan-Newman algorithm**

1 - compute link **betweenness** for all the links

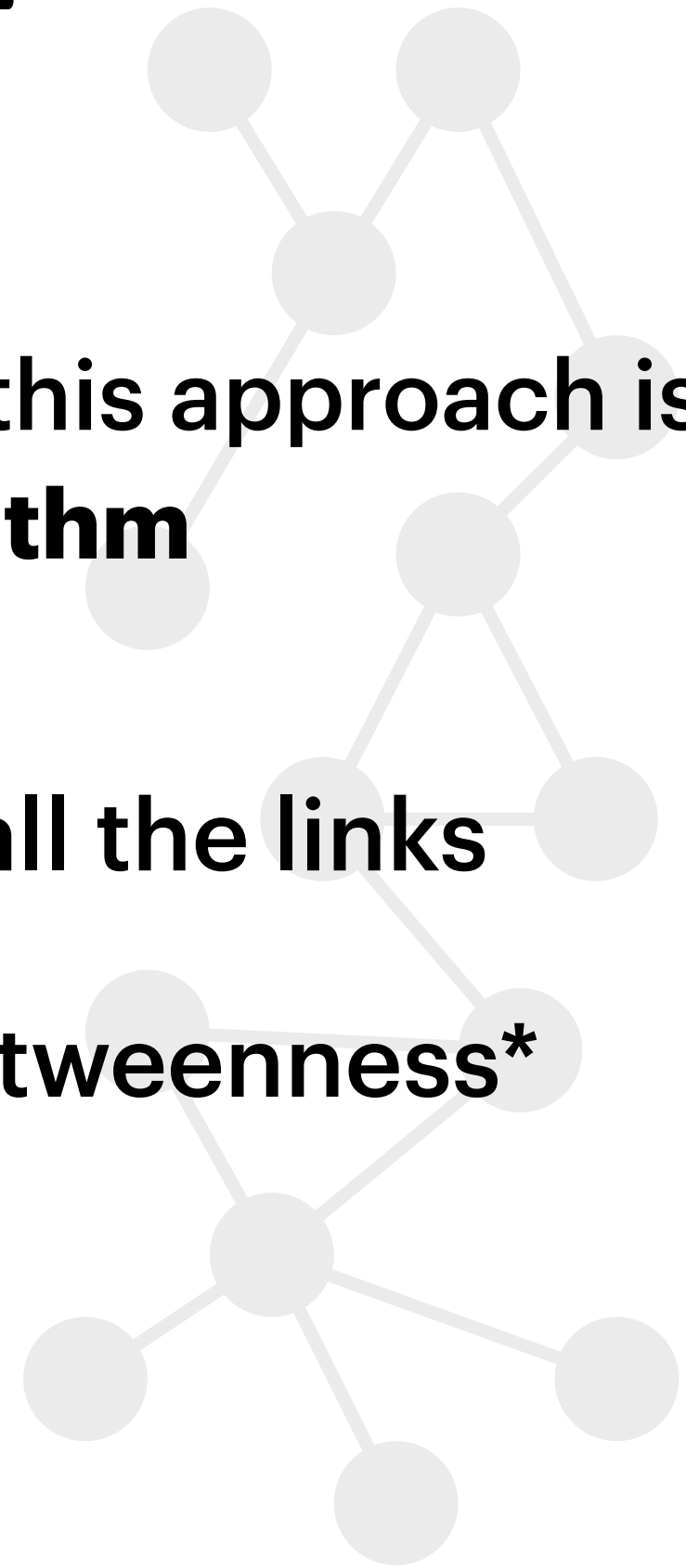


# BRIDGE REMOVAL

The most famous algorithm based on this approach is the **Girvan-Newman algorithm**

- 1 - compute link **betweenness** for all the links
- 2 - **remove** the link with highest betweenness\*

\*in case of a tie, pick a random one among those with highest betweenness

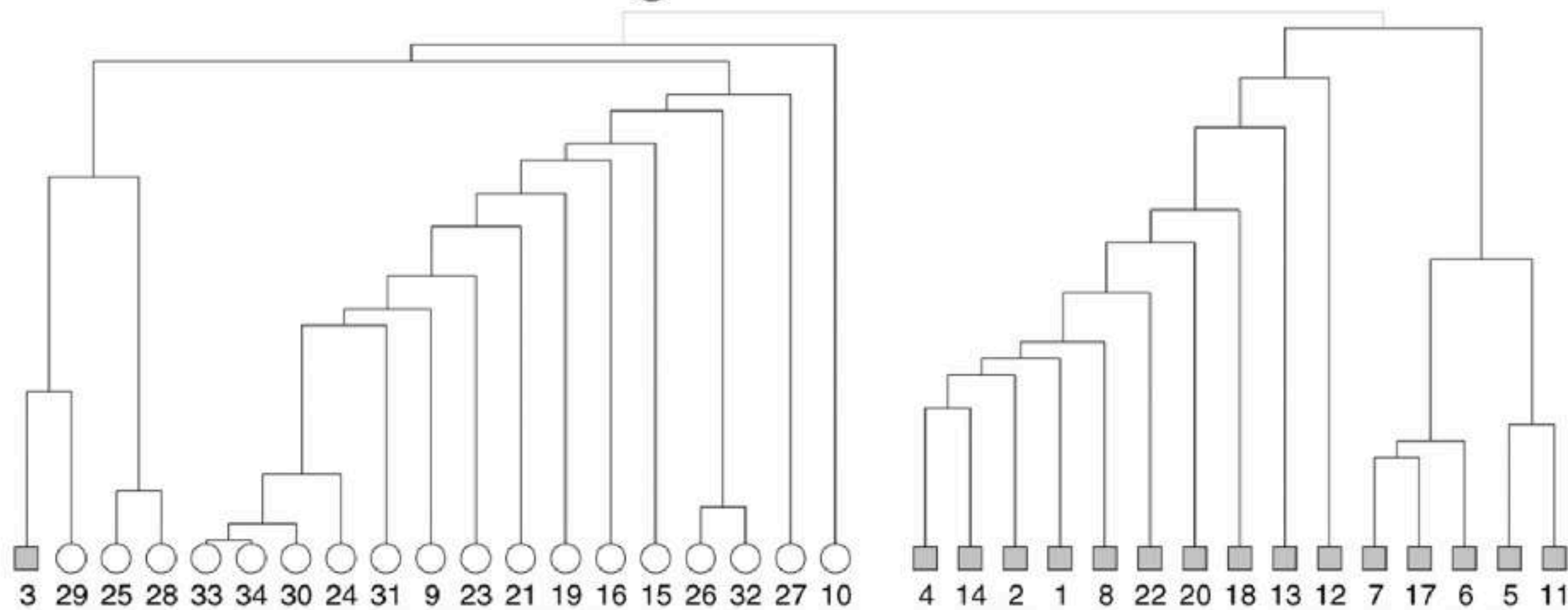
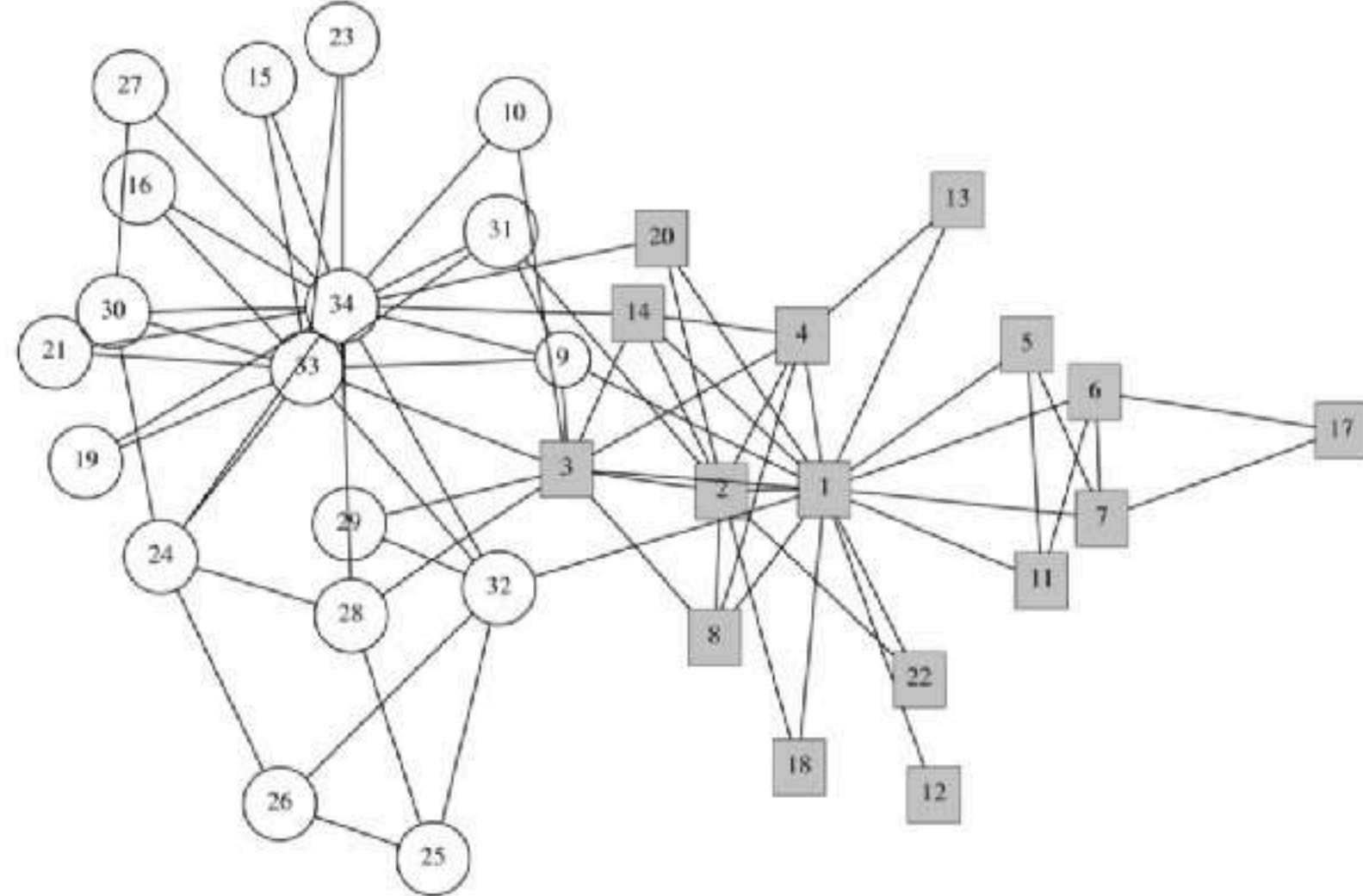


# BRIDGE REMOVAL

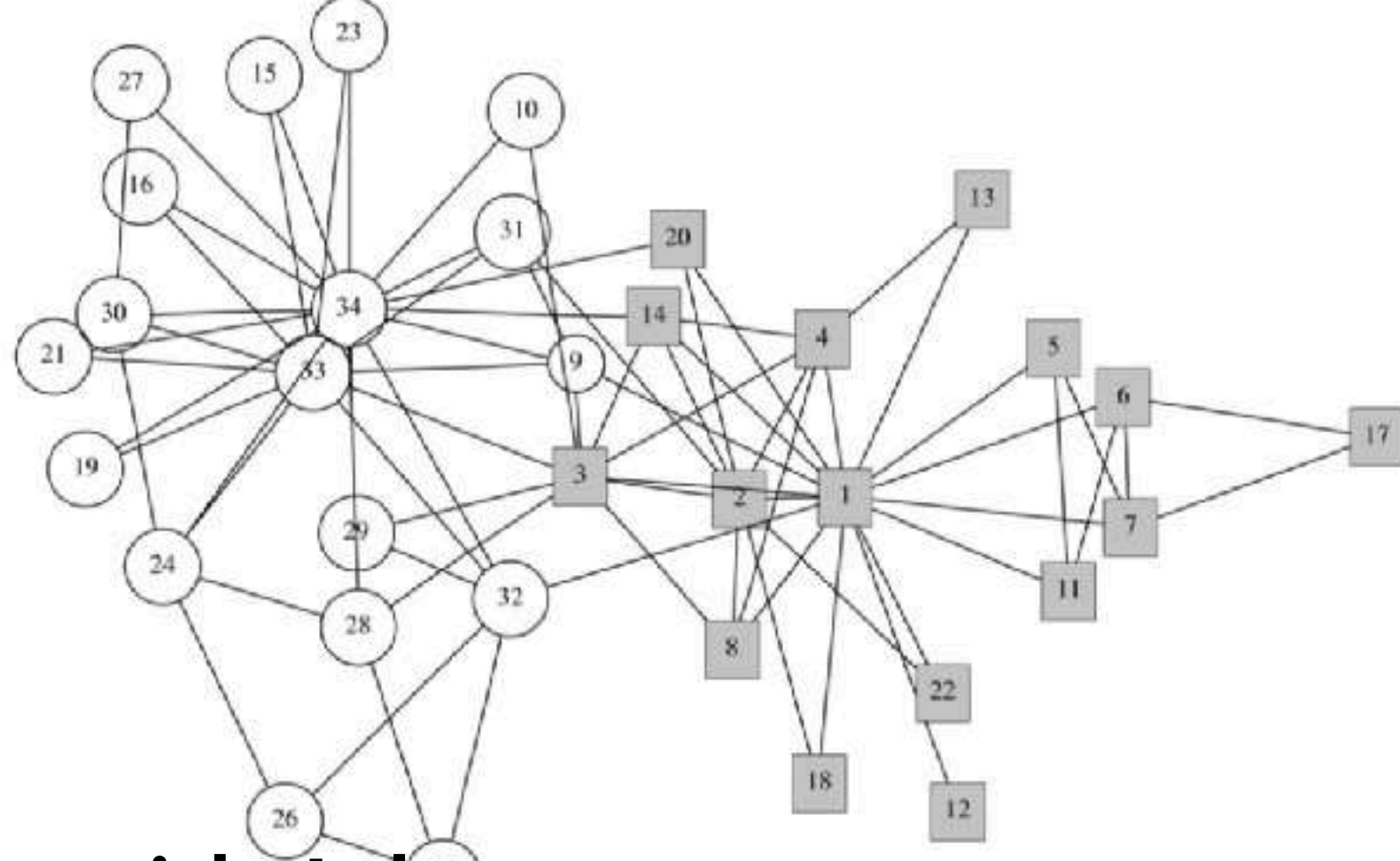
The most famous algorithm based on this approach is the **Girvan-Newman algorithm**

- 1 - compute link **betweenness** for all the links
- 2 - **remove** the link with highest betweenness\*
- 3 - **repeat** 1 and 2 until you have no links left

\*in case of a tie, pick a random one among those with

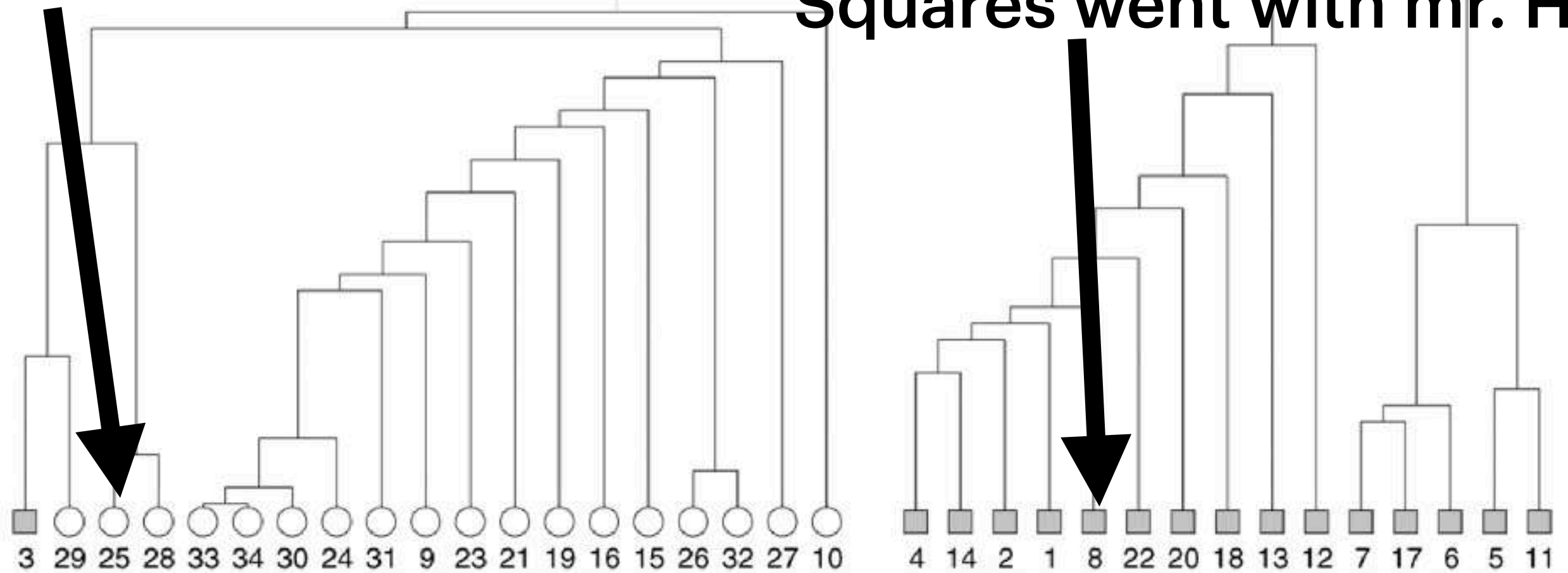






**Circles went with John a**

**Squares went with mr. Hi**



# FINAL VERDICT



**GREAT FIRST ATTEMPT, BUT COMPUTING LINK  
BETWEENNESS FOR LARGE NETWORKS THAT MANY  
TIMES IS IMPOSSIBLE**



# MODULARITY MAXIMISATION

**MAIN IDEA: WE CALCULATE HOW GOOD A COMMUNITY IS VS RANDOM BASELINE**



# MODULARITY MAXIMISATION

**MAIN IDEA: WE CALCULATE HOW GOOD A COMMUNITY IS VS RANDOM BASELINE**

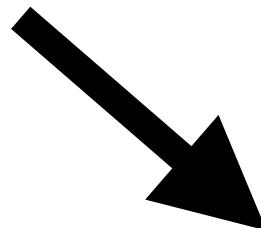
Originally introduced to know **where to cut** the dendrogram in Girvan-Newman



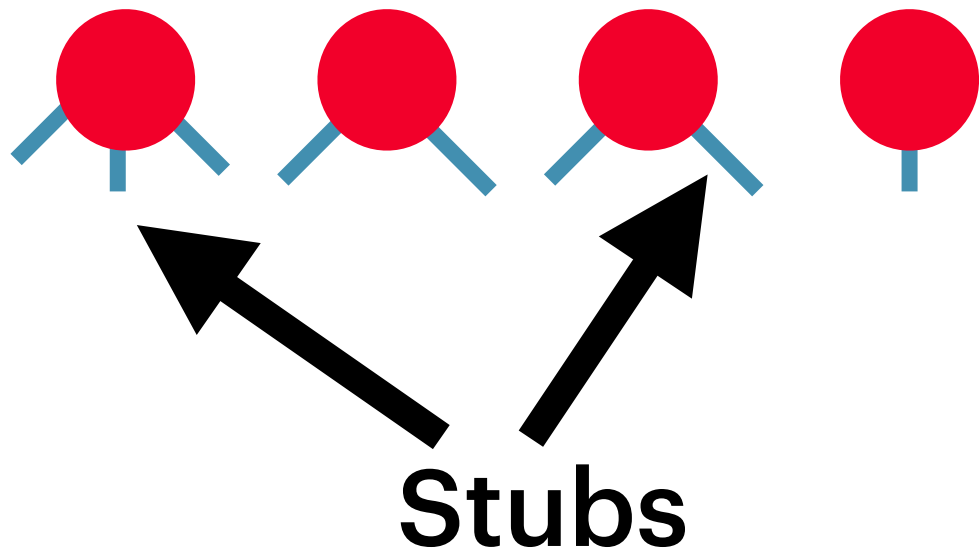


# MODULARITY MAXIMISATION

Difference between links in  $c$   
and expected links in  $c$  with  
configuration model


$$Q = \frac{1}{L} \sum_c \left( L_c - \frac{k_c^2}{4L} \right)$$

# MODULARITY MAXIMISATION

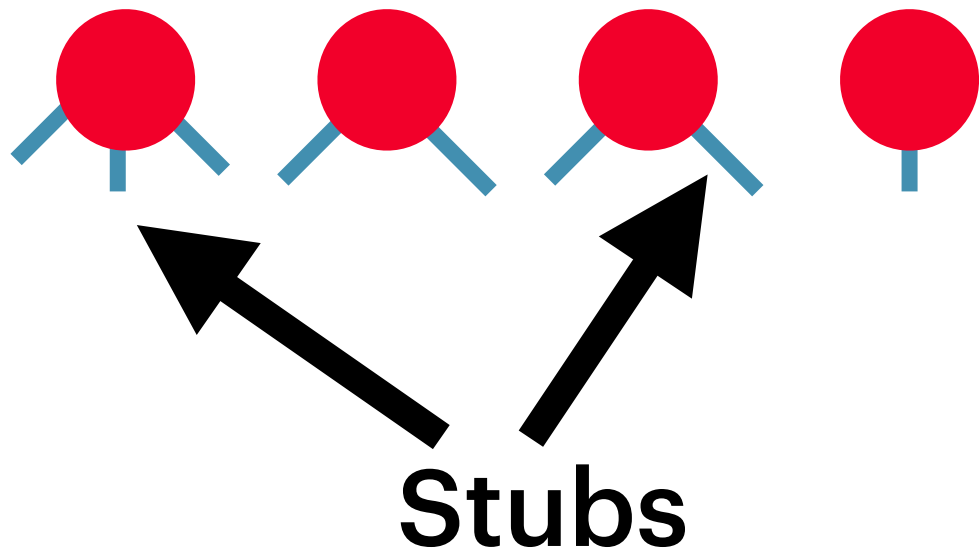


$$Q = \frac{1}{L} \sum_c \left( L_c - \frac{k_c^2}{4L} \right)$$

$\frac{k_c}{2L}$  Is the probability of randomly choosing **one stub** in the community



# MODULARITY MAXIMISATION



$$Q = \frac{1}{L} \sum_c \left( L_c - \frac{k_c^2}{4L} \right)$$

$\left( \frac{k_c}{2L} \right)^2$  Is the probability of randomly choosing **two stubs** in the community

# MODULARITY MAXIMISATION

There are  $L$  links in the network

# MODULARITY MAXIMISATION

There are  $L$  links in the network

Each link joins two stubs from community  $c$  with probability

$$\left( \frac{k_c}{2L} \right)^2$$

# MODULARITY MAXIMISATION

There are  $L$  links in the network

Each link joins two stubs from community  $c$  with probability

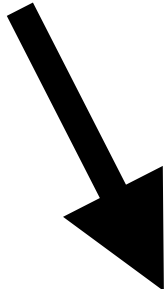
$$\left(\frac{k_c}{2L}\right)^2$$

Then, the expected number of links in the community is

$$L \left(\frac{k_c}{2L}\right)^2 = \frac{k_c^2}{4L}$$

# MODULARITY MAXIMISATION

Average


$$Q = \frac{1}{L} \sum_c \left( L_c - \frac{k_c^2}{4L} \right)$$



Difference between actual links in c and expected links in c

# MODULARITY MAXIMISATION

$$\text{Directed } Q_d = \frac{1}{L} \sum_c \left( L_c - \frac{k_C^{in} k_C^{out}}{L} \right)$$

$$\text{Weighted } Q_w = \frac{1}{W} \sum_c \left( W_c - \frac{s_C^2}{4W} \right)$$

$$\text{Weighted and directed } Q_{dw} = \frac{1}{W} \sum_c \left( W_c - \frac{s_C^{in} s_C^{out}}{W} \right)$$

# MODULARITY MAXIMISATION

Most famous algorithms: **Louvain, Leiden**

# MODULARITY MAXIMISATION

Most famous algorithms: **Louvain, Leiden**

- 1) start with no communities. Every nodes is moved to a community so that  $Q$  is maximised. Repeat until no modularity gain is possible

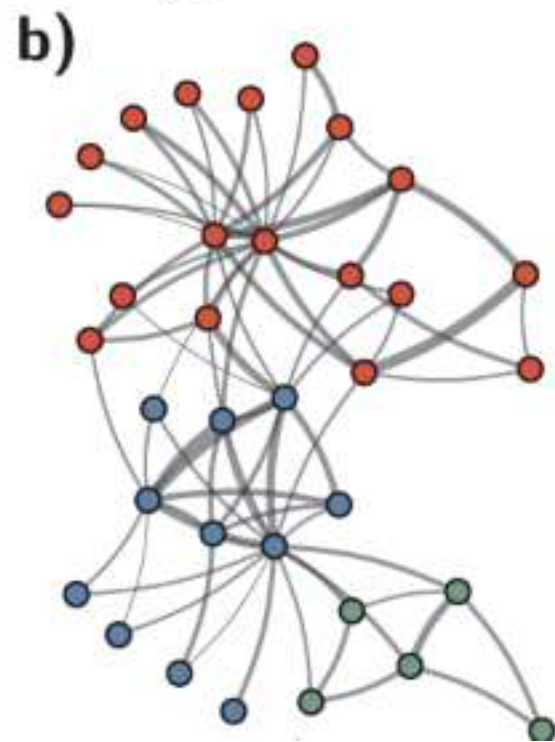
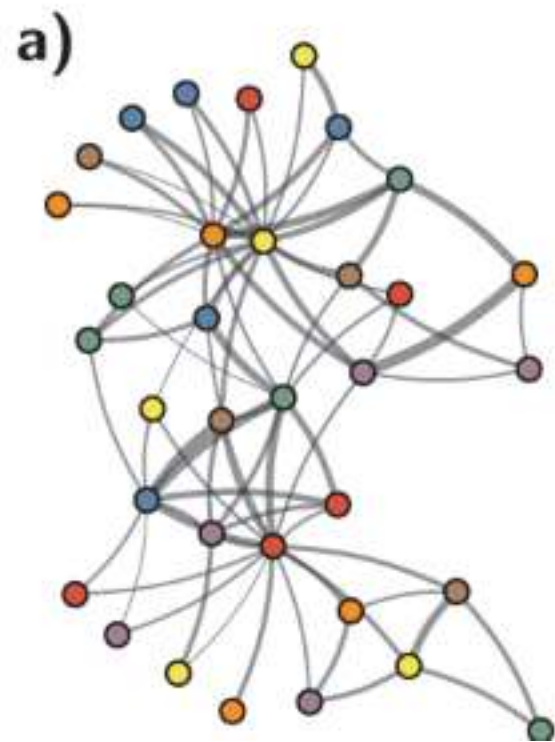


# MODULARITY MAXIMISATION

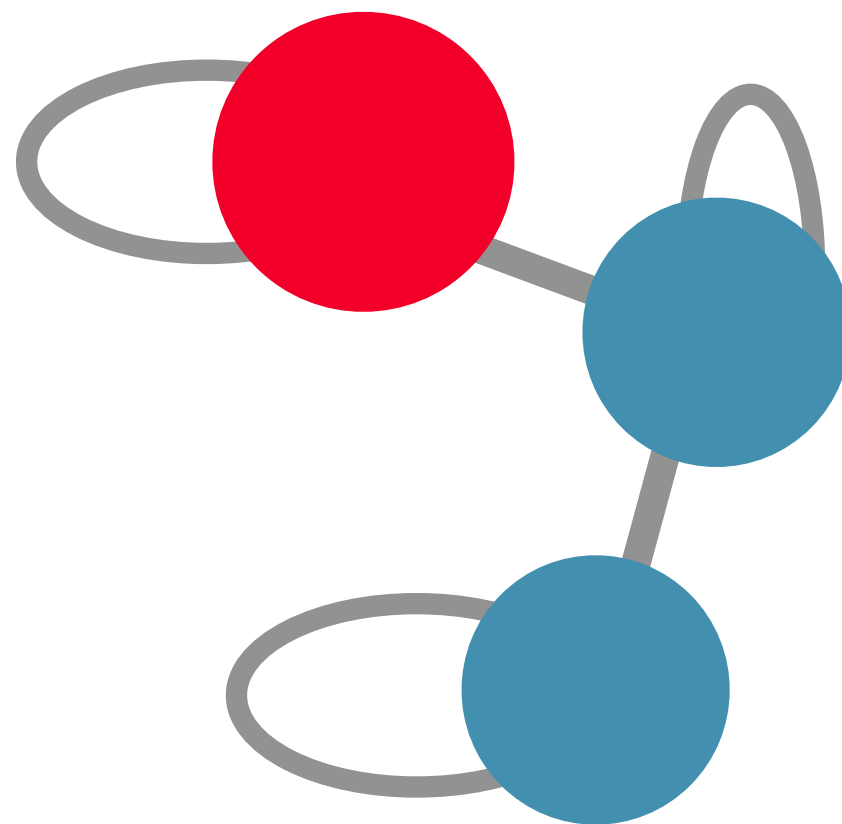
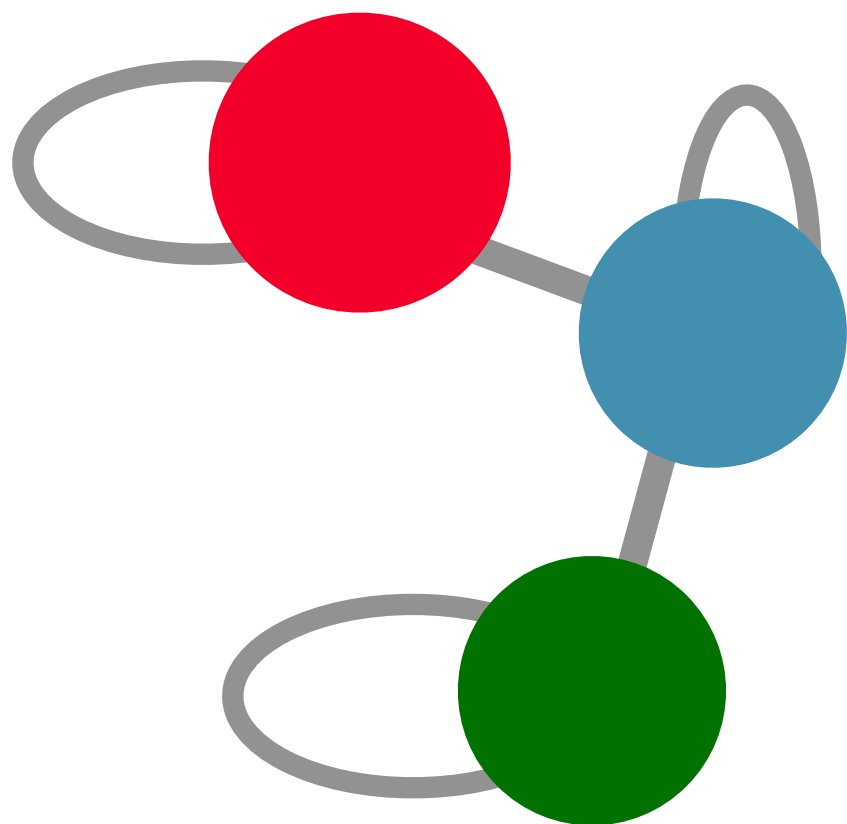
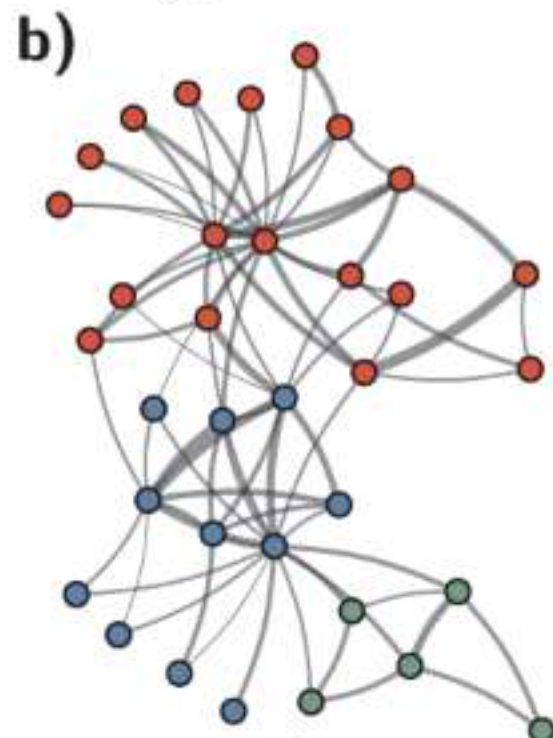
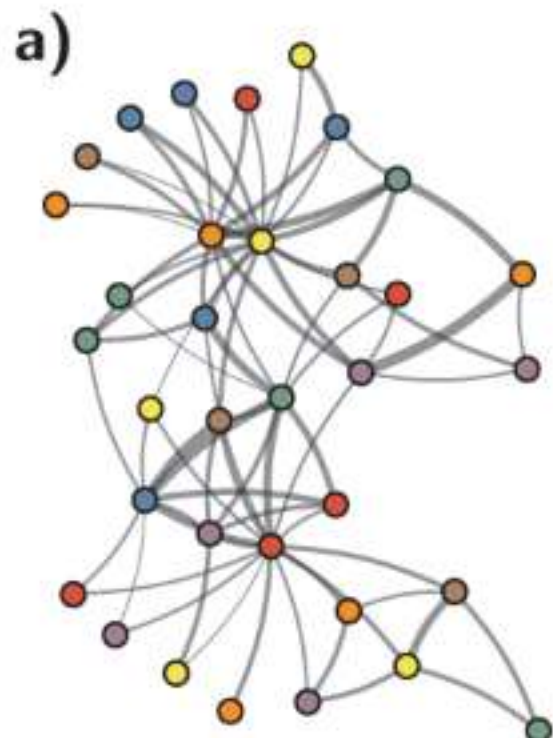
Most famous algorithms: **Louvain, Leiden**

- 1) start with no communities. Every nodes is moved to a community so that  $Q$  is maximised. Repeat until no modularity gain is possible
- 2) the network becomes a weighted super-network, in which nodes are the communities of the original network, and weights are the number of links between communities (this includes self-loops)

Move nodes



Move nodes



# MODULARITY MAXIMISATION PROBLEMS

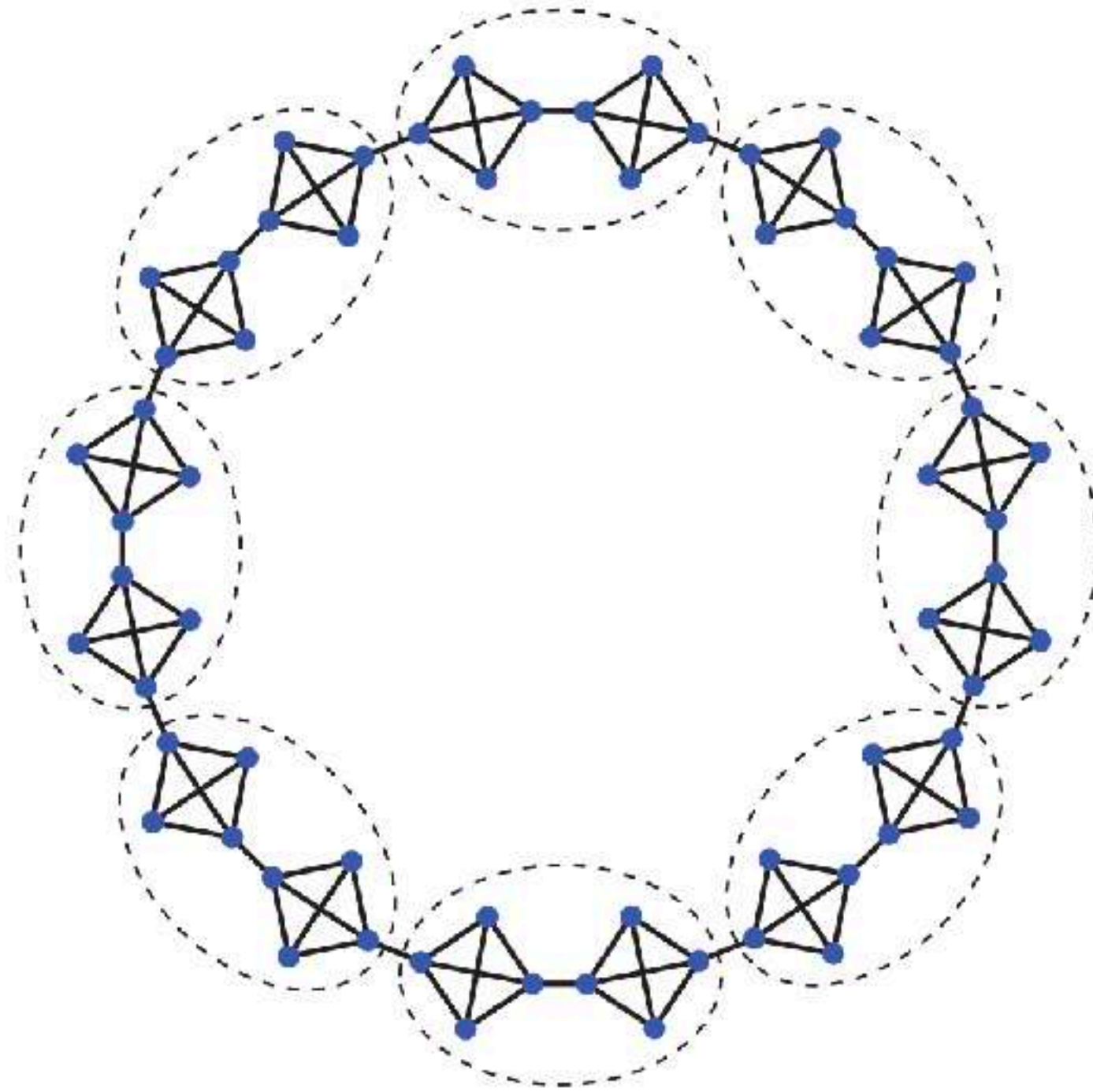
**Comparison:** On average Larger networks have larger modularity

**Uncertainty:** this approach can find positive modularity for random networks

**Resolution:** cannot find communities whose degree is smaller than

$$\sqrt{2L}$$

# MODULARITY MAXIMISATION

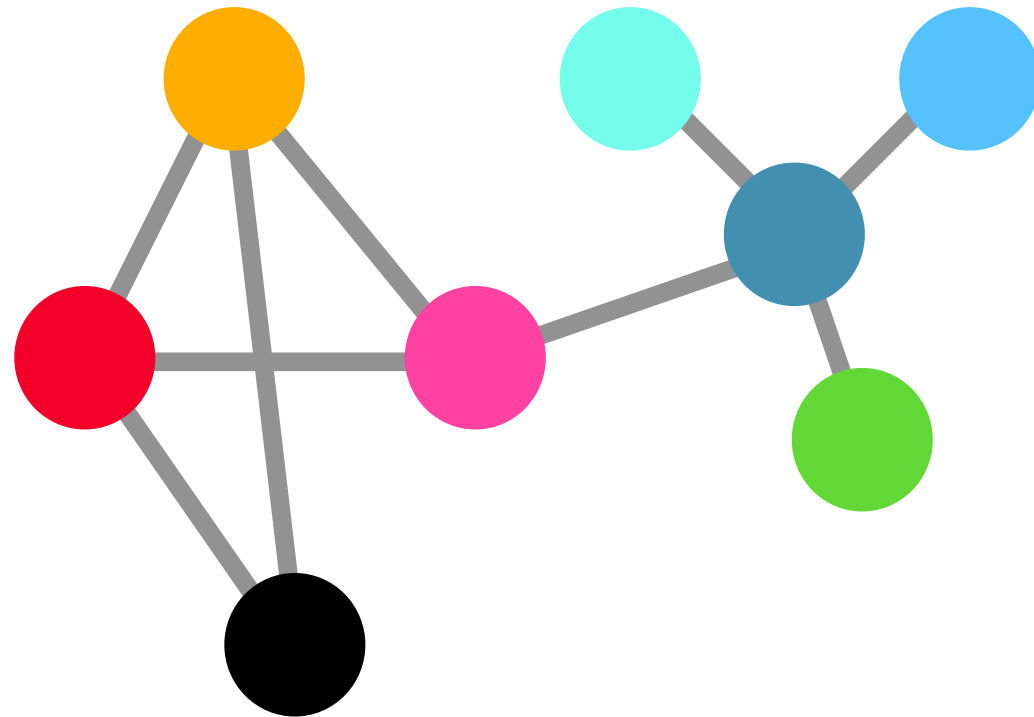


# LABEL PROPAGATION

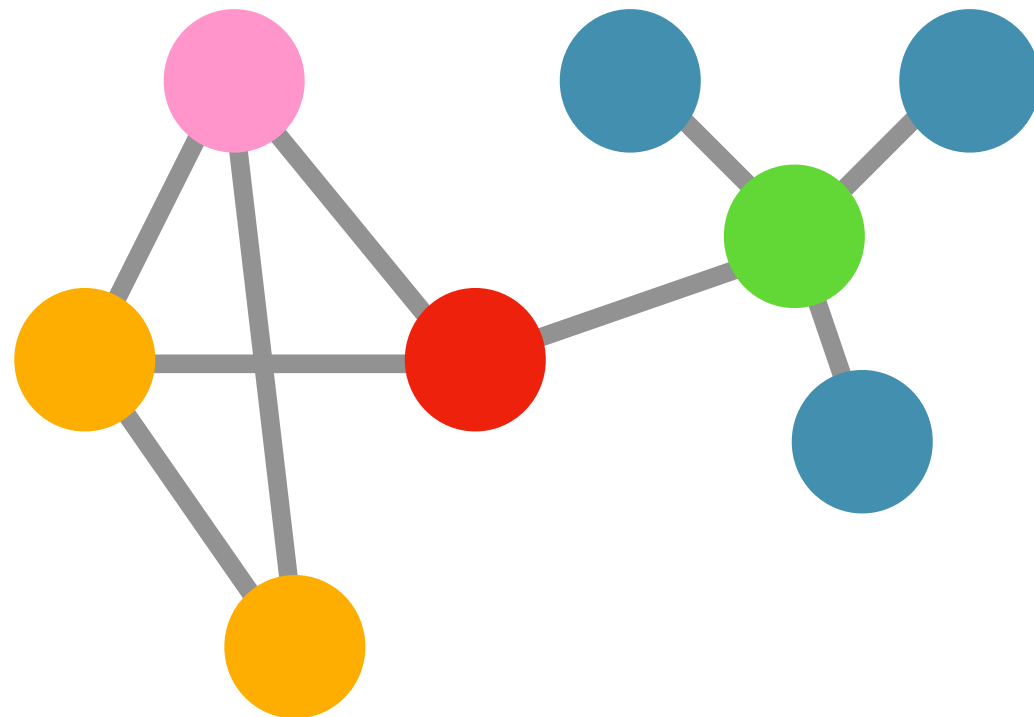
- 1) WE START WITH SINGLETONS**
- 2) ONE BY ONE, WITH RANDOM ORDER, NODES TAKE THE "LABEL" (IE COMMUNITY MEMBERSHIP) OF THE MAJORITY OF THEIR NEIGHBOURS**
- 3) WE REPEAT THIS UNTIL THE PARTITION IS STABLE (IE THERE ARE NO POSSIBLE CHANGES)**



# LABEL PROPAGATION

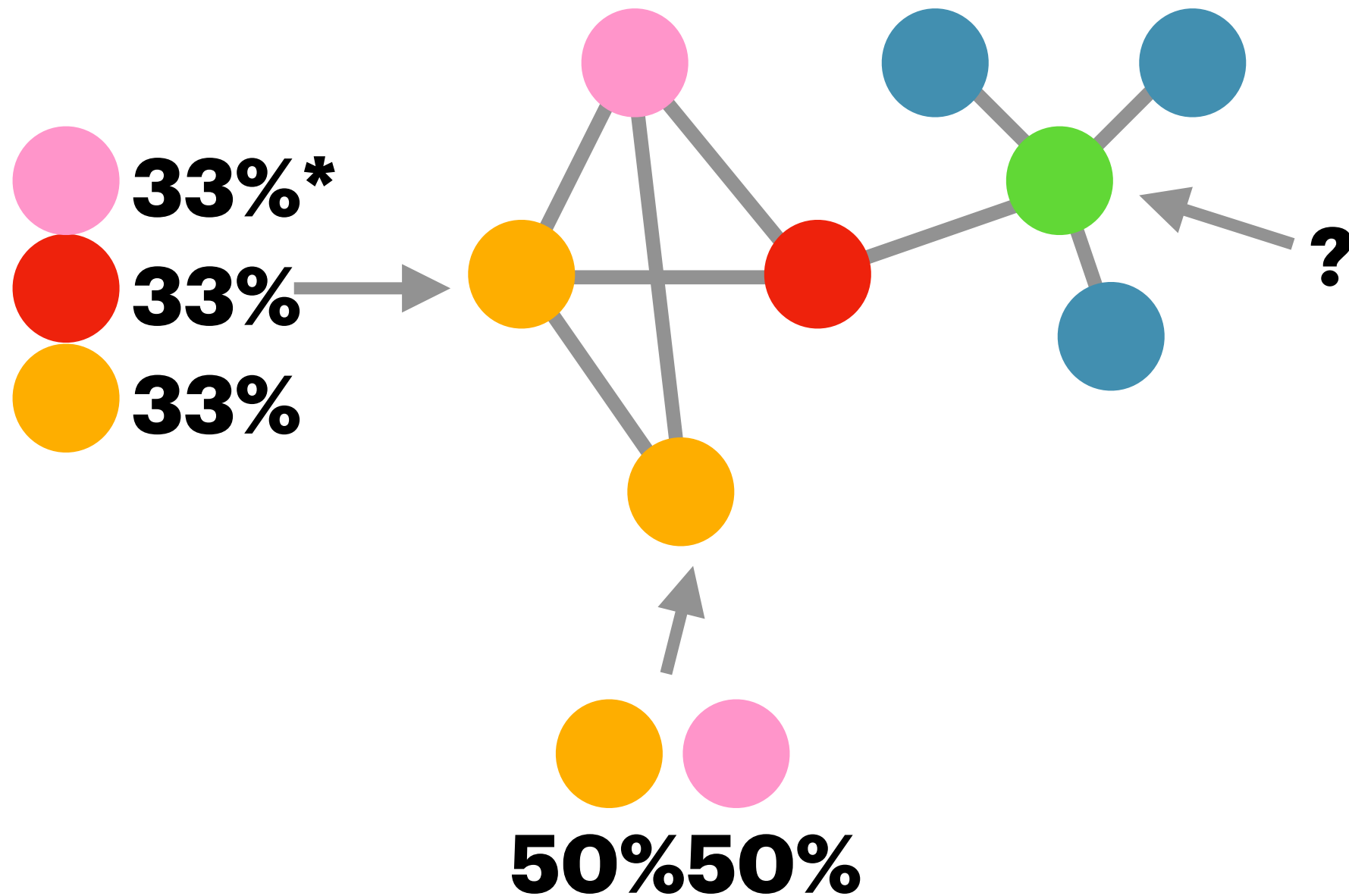


# LABEL PROPAGATION



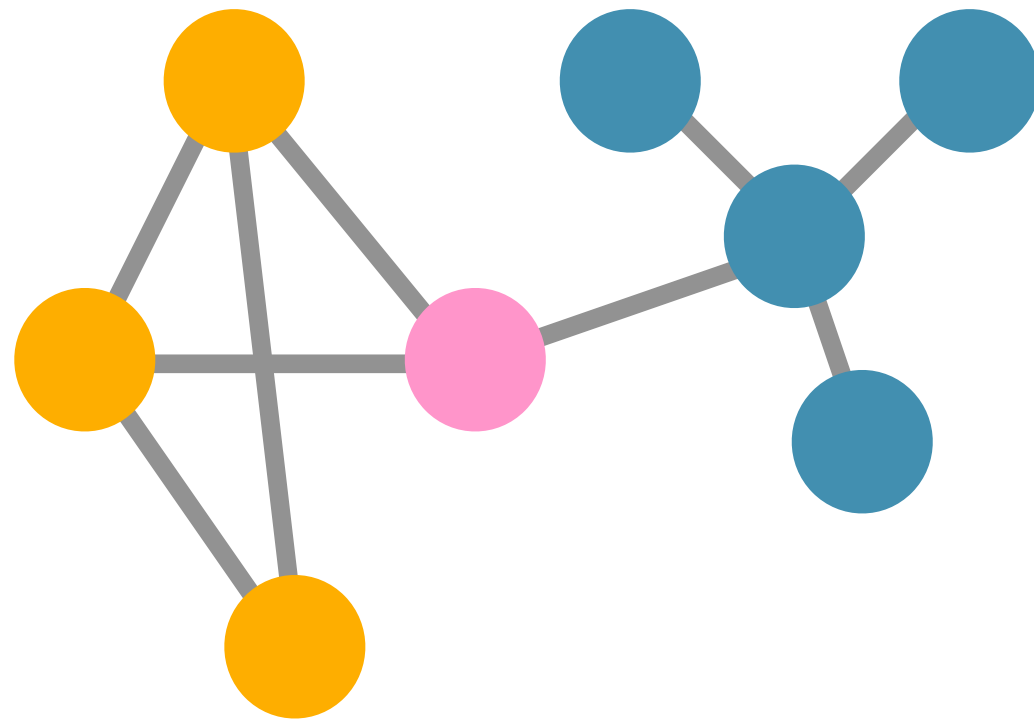


# LABEL PROPAGATION

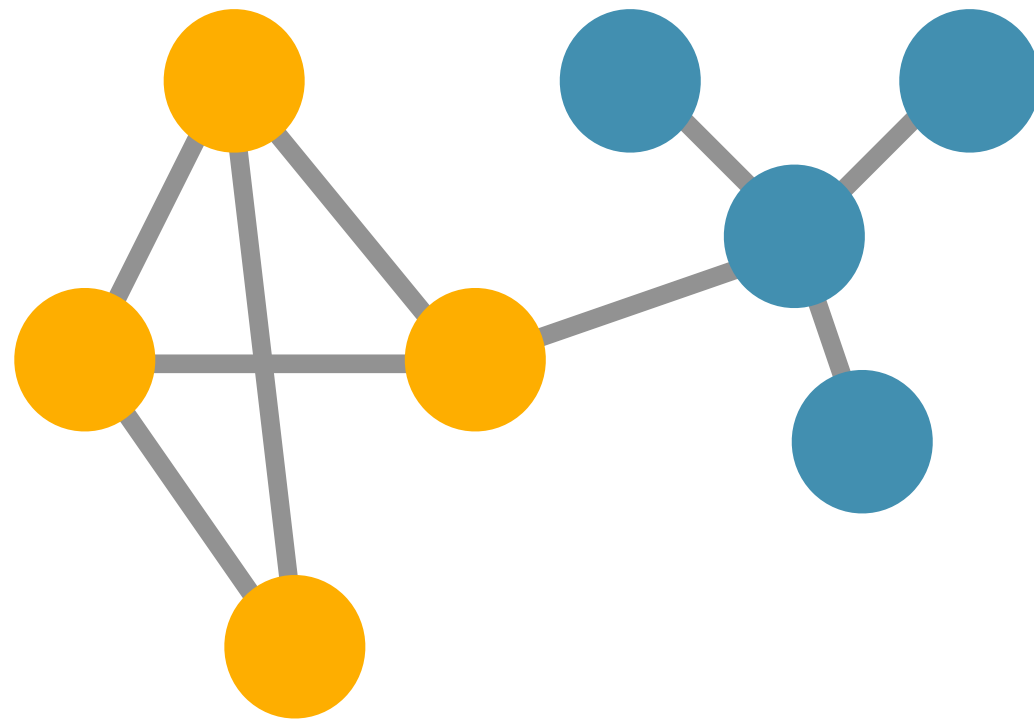


**\*Actually 1/3!!!**

# LABEL PROPAGATION



# LABEL PROPAGATION



# **LABEL PROPAGATION**

## **ISSUES**

**DIFFERENT RUNS FIND DIFFERENT COMMUNITIES  
NEEDS TO BE RUN MULTIPLE TIMES**

## **STRENGTHS**

**VERY FAST**

**IF SOME MEMBERSHIPS ARE KNOWN, THEY CAN BE  
USED TO INITIALISE THE NETWORK**



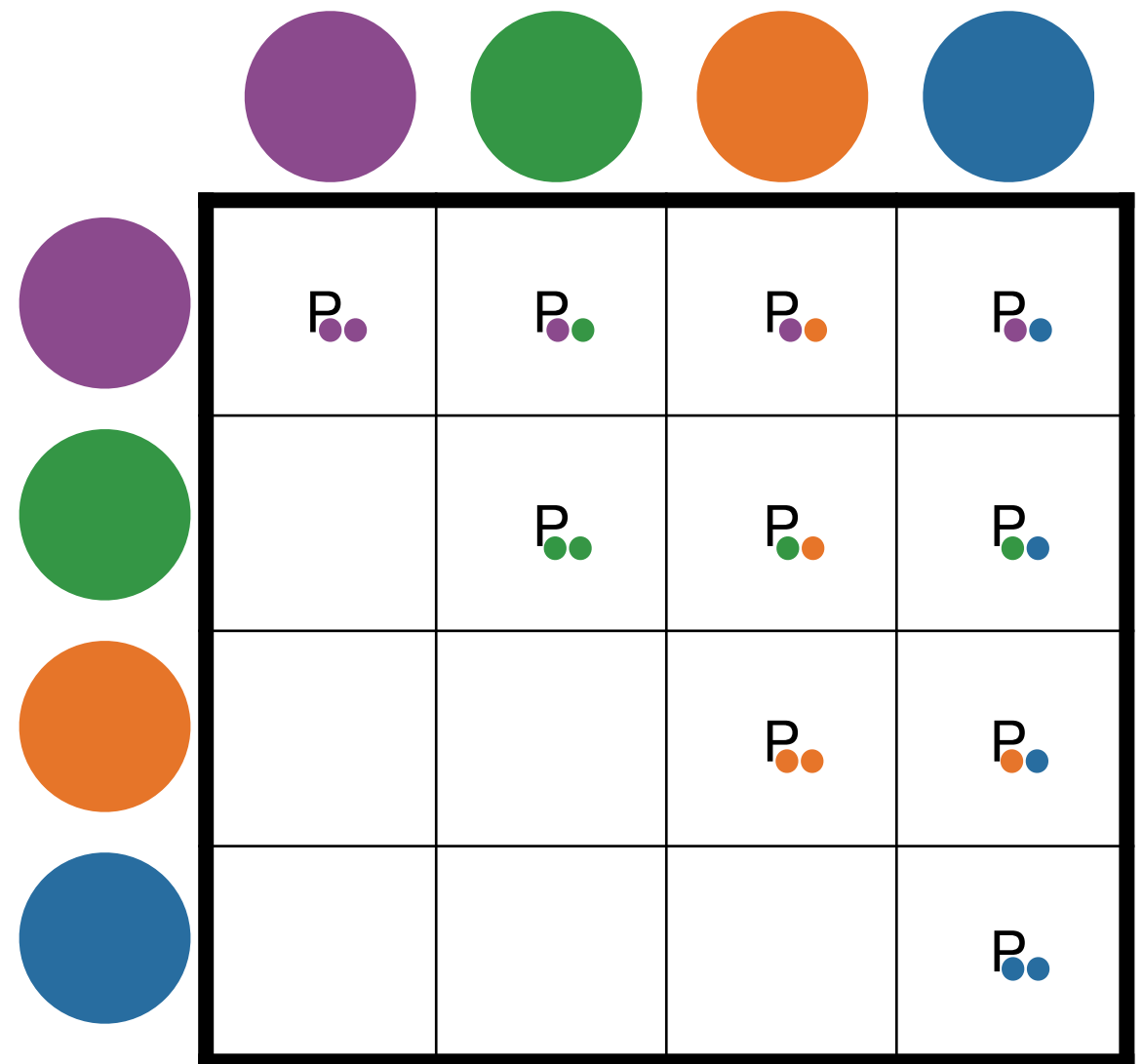
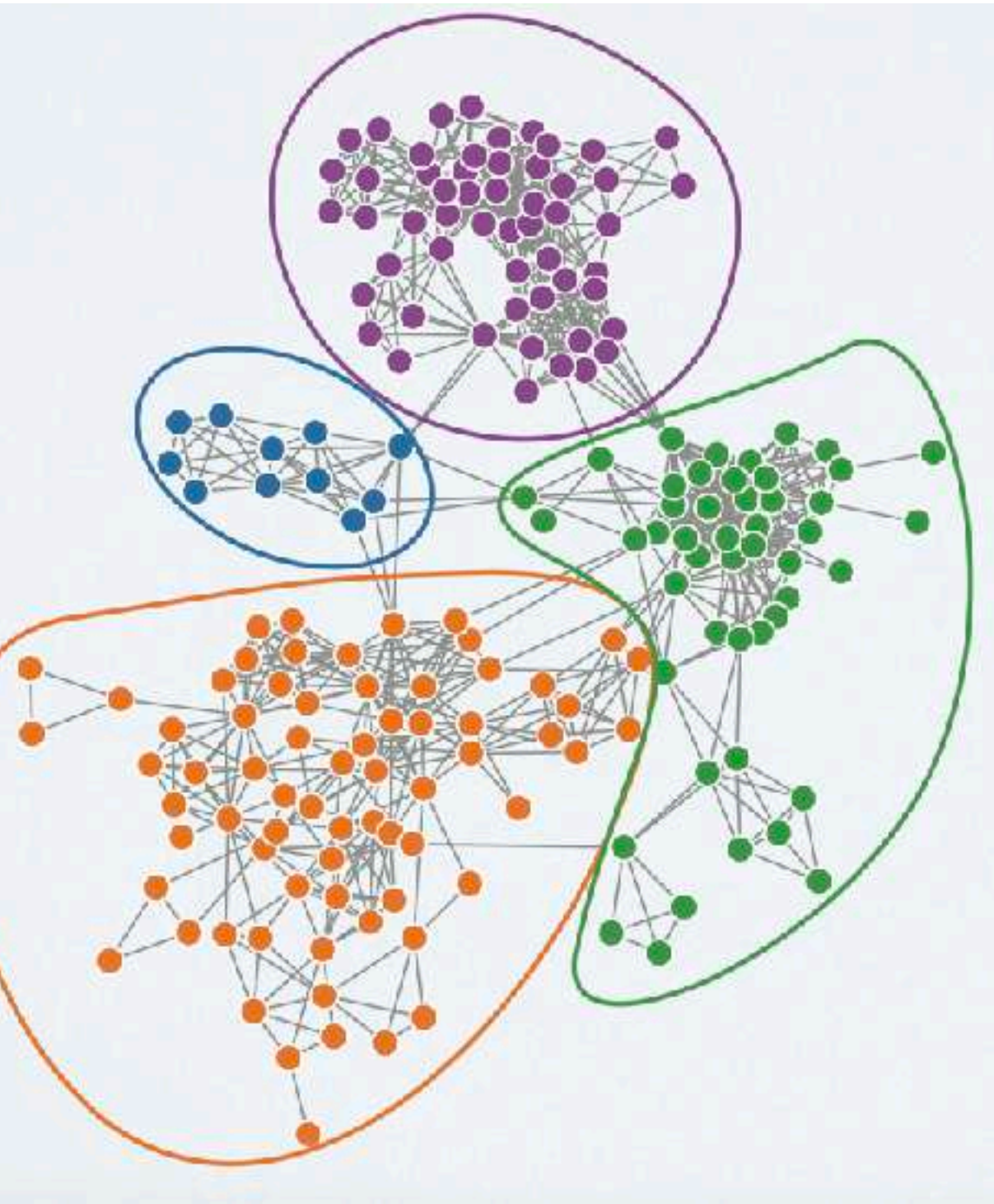
# STOCHASTIC BLOCK MODEL

**Generative algorithm**

**generates communities with given probabilities,  
chooses the most likely**



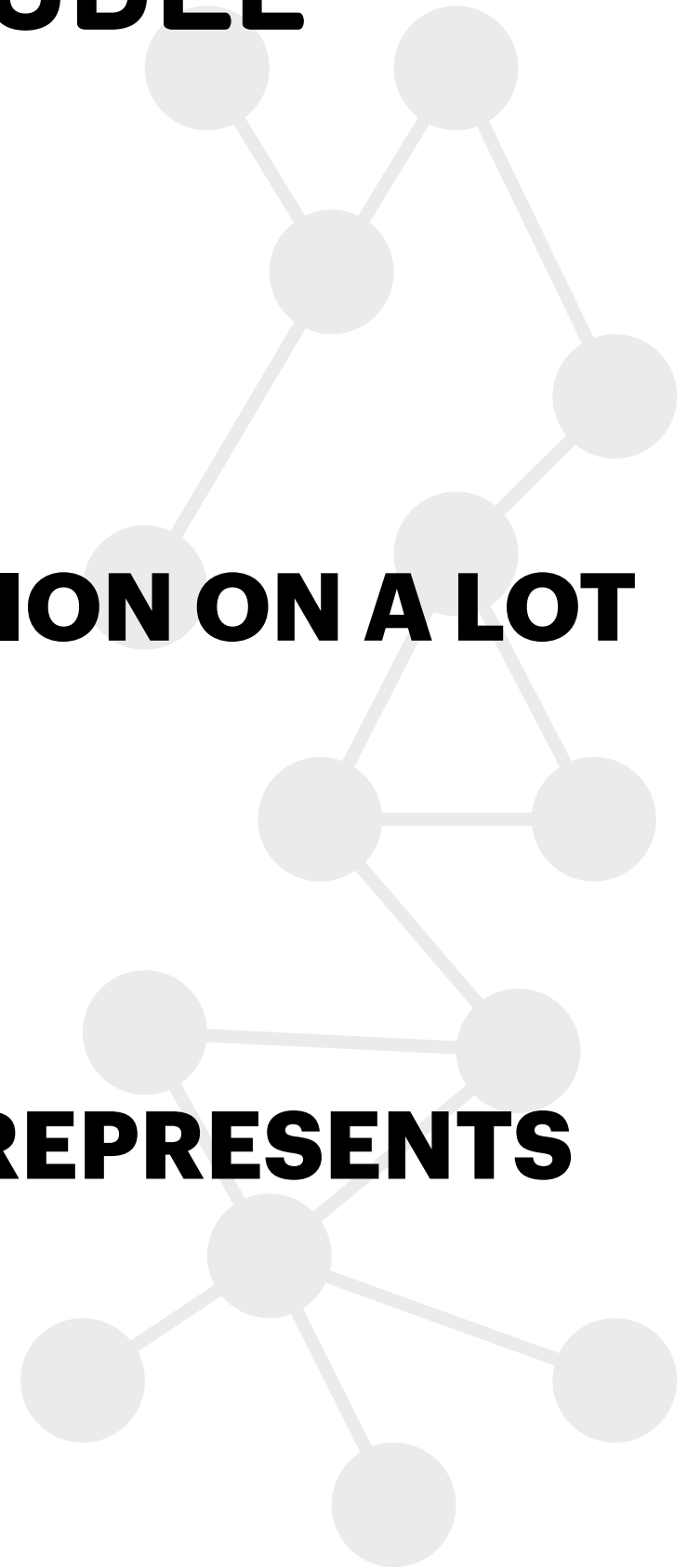
# STOCHASTIC BLOCK MODEL



# STOCHASTIC BLOCK MODEL

**CAN PERFORM COMMUNITY DETECTION ON A LOT OF DIFFERENT NETWORK TYPES**

**FOR EXAMPLE: IF  $\forall r, p_{rr} = 0$  THIS REPRESENTS MULTIPARTITE NETWORKS**





# STOCHASTIC BLOCK MODEL

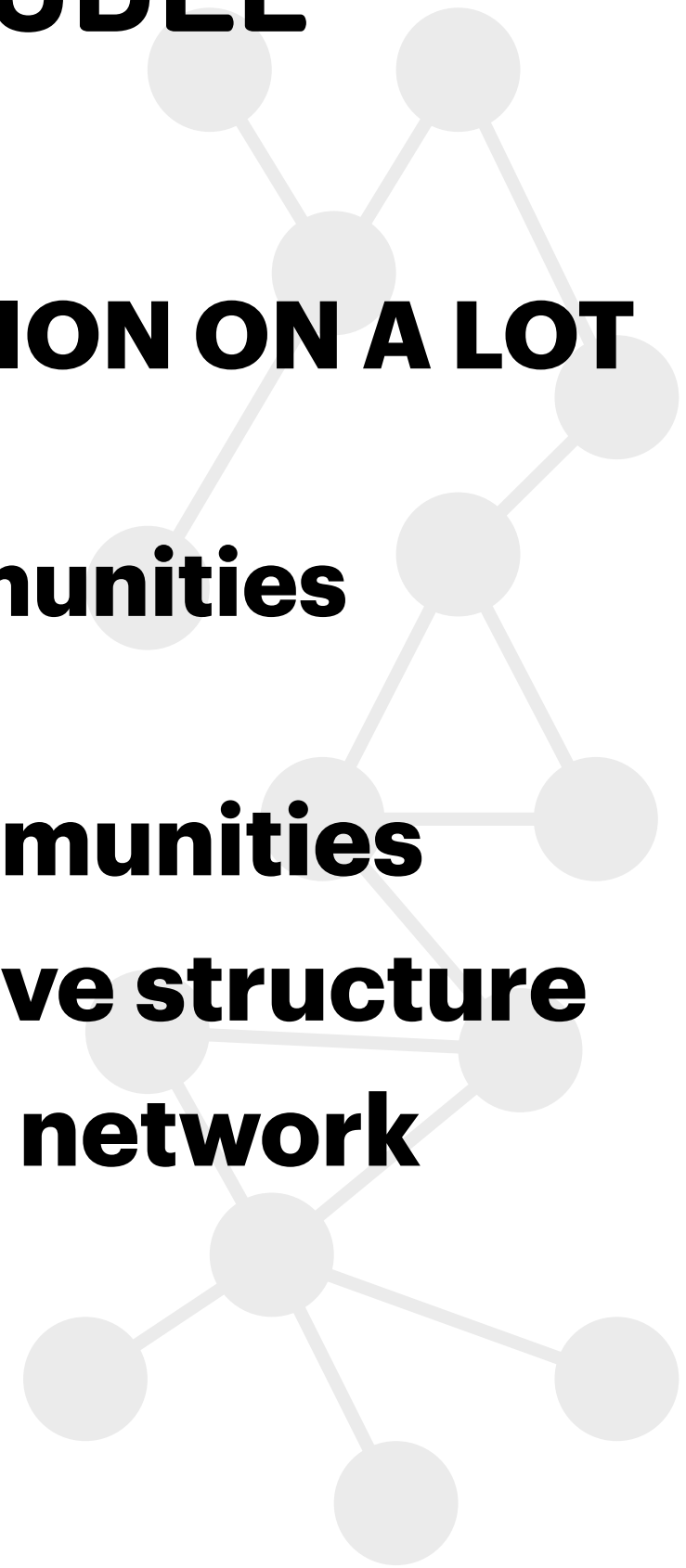
**CAN PERFORM COMMUNITY DETECTION ON A LOT OF DIFFERENT NETWORK TYPES**

**And can discover more than just communities**

$\forall r, s \quad p_{rr} > p_{rs}$  **Classic communities**

$p_{rr} < p_{rs}$  **Disassortative structure**

$\forall r \quad p_{rr} = 0$  **Multipartite network**



# STOCHASTIC BLOCK MODEL

**CAN PERFORM COMMUNITY DETECTION ON A LOT OF DIFFERENT NETWORK TYPES**

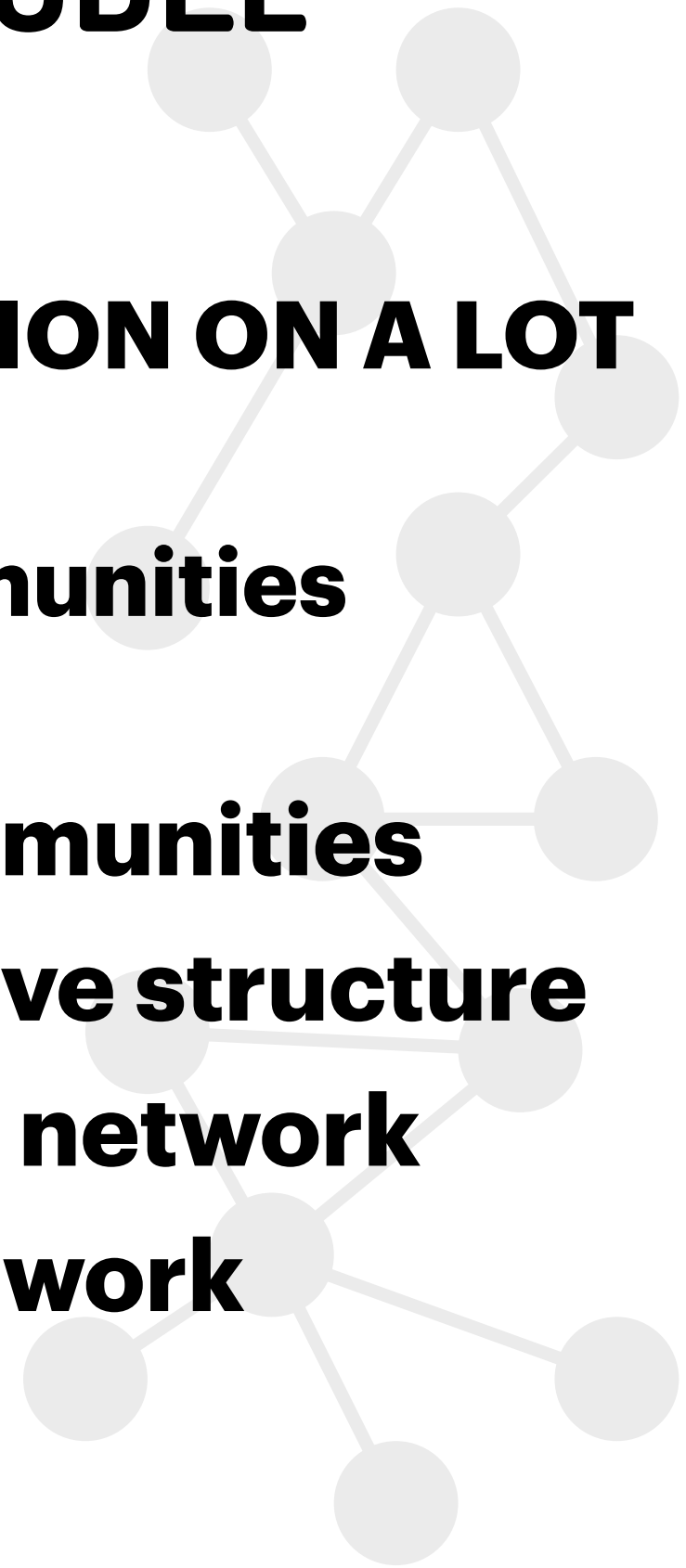
**And can discover more than just communities**

$\forall r, s \quad p_{rr} > p_{rs}$  **Classic communities**

$p_{rr} < p_{rs}$  **Disassortative structure**

$\forall r \quad p_{rr} = 0$  **Multipartite network**

$\forall r, s \quad p_{rr} = p_{rs} = p$  **Random network**



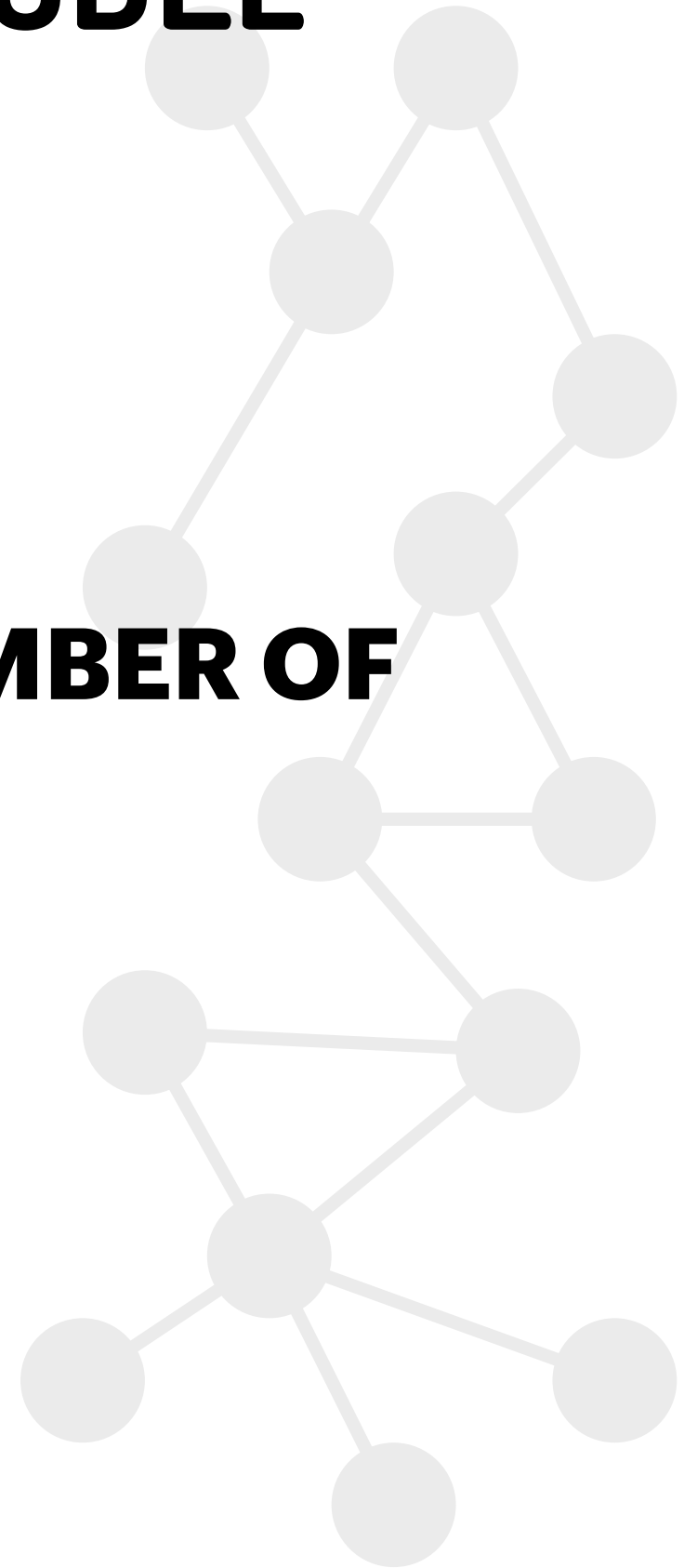
# STOCHASTIC BLOCK MODEL

**LIMITS:**

**NEEDS PRIOR KNOWLEDGE ON NUMBER OF  
COMMUNITIES**

**STRENGTHS:**

**EVERYTHING ELSE**



A circular chord diagram showing complex interactions between nodes on the perimeter. The flows are represented by blue ribbons of varying thickness, indicating the strength of connections. The diagram is dense with many overlapping paths.

Uses Bayesian inference

Does not require prior knowledge

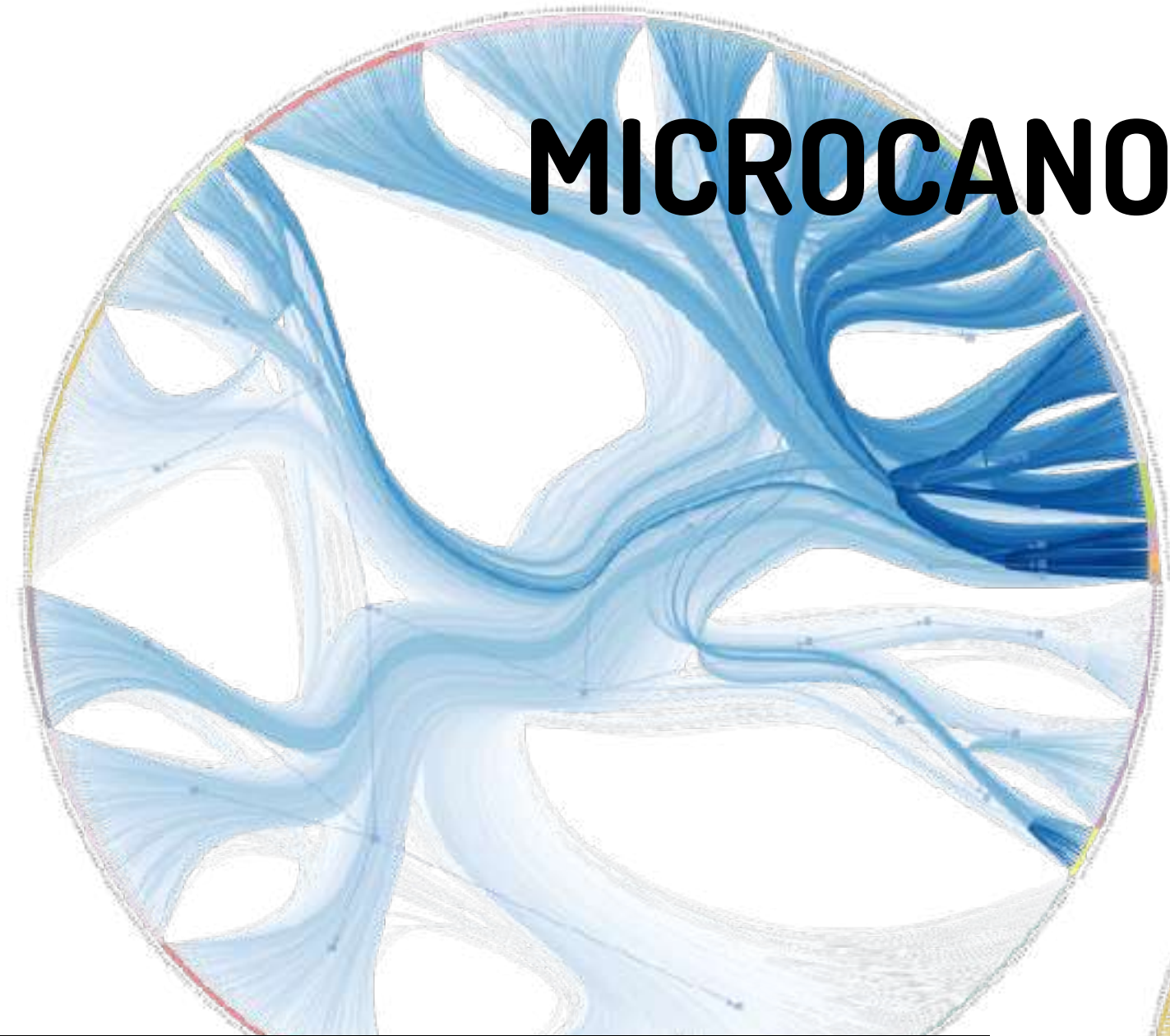
Extremely versatile

A circular chord diagram similar to the one on the left, but with a much wider color palette. The flows are represented by ribbons in various colors including red, yellow, green, purple, and teal, making the structure more visually distinct and complex.

**MICROCANONICAL SBM**



# MICROCANONICAL SBM



Fast and scalable

Explainable

There's a library that does it all and produces beautiful figures

