# NETWORK COMMUNITY ANALYSIS

Yong-Yeol "YY" Ahn

SCHOOL OF INFORMATICS
AND COMPUTING

INDIANA UNIVERSITY

Bloomington
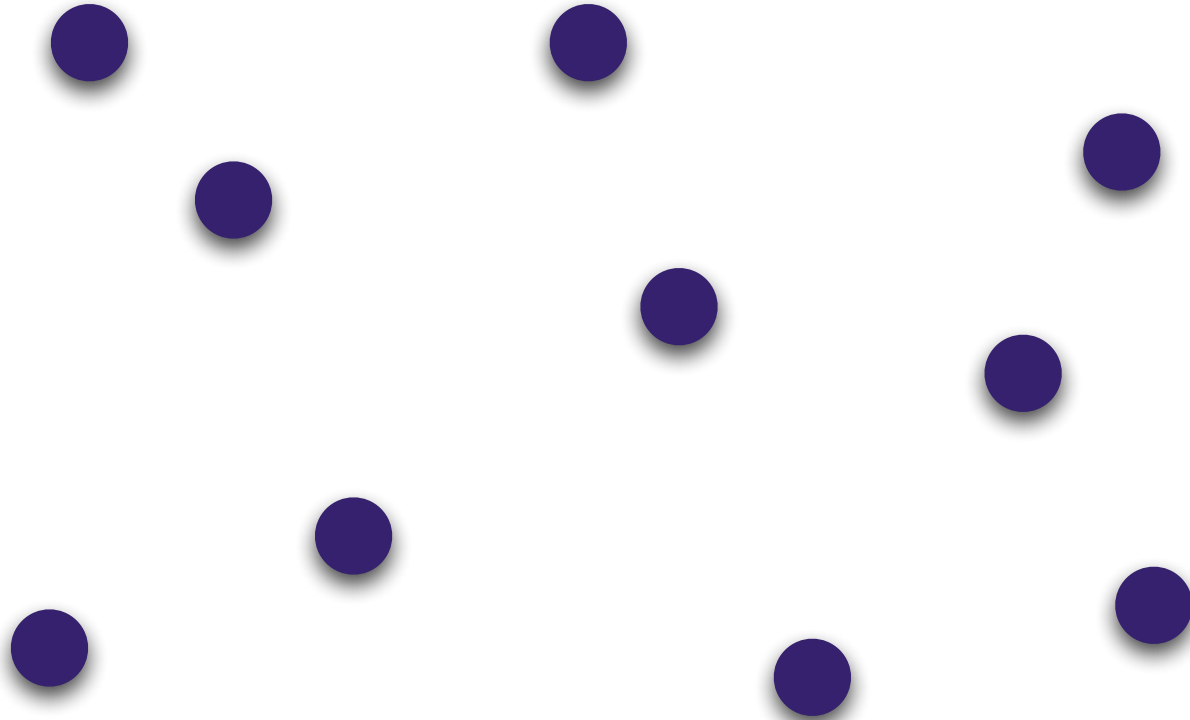
# Disclaimer:

We're just scratching the surface.
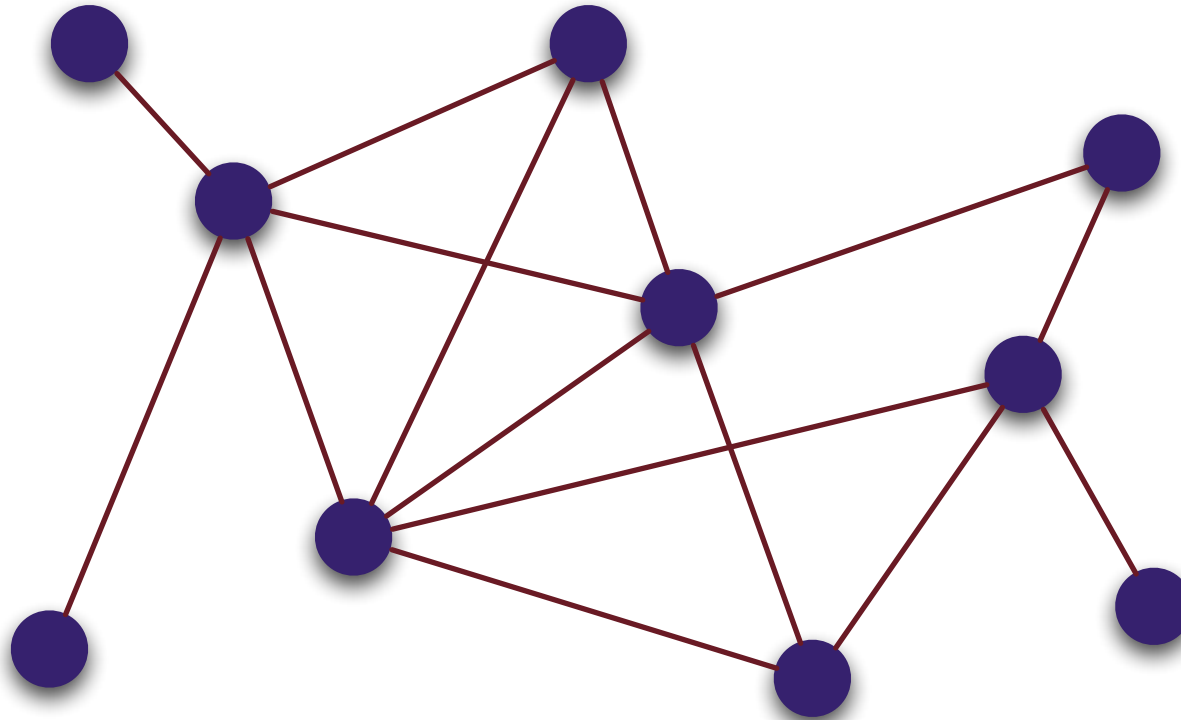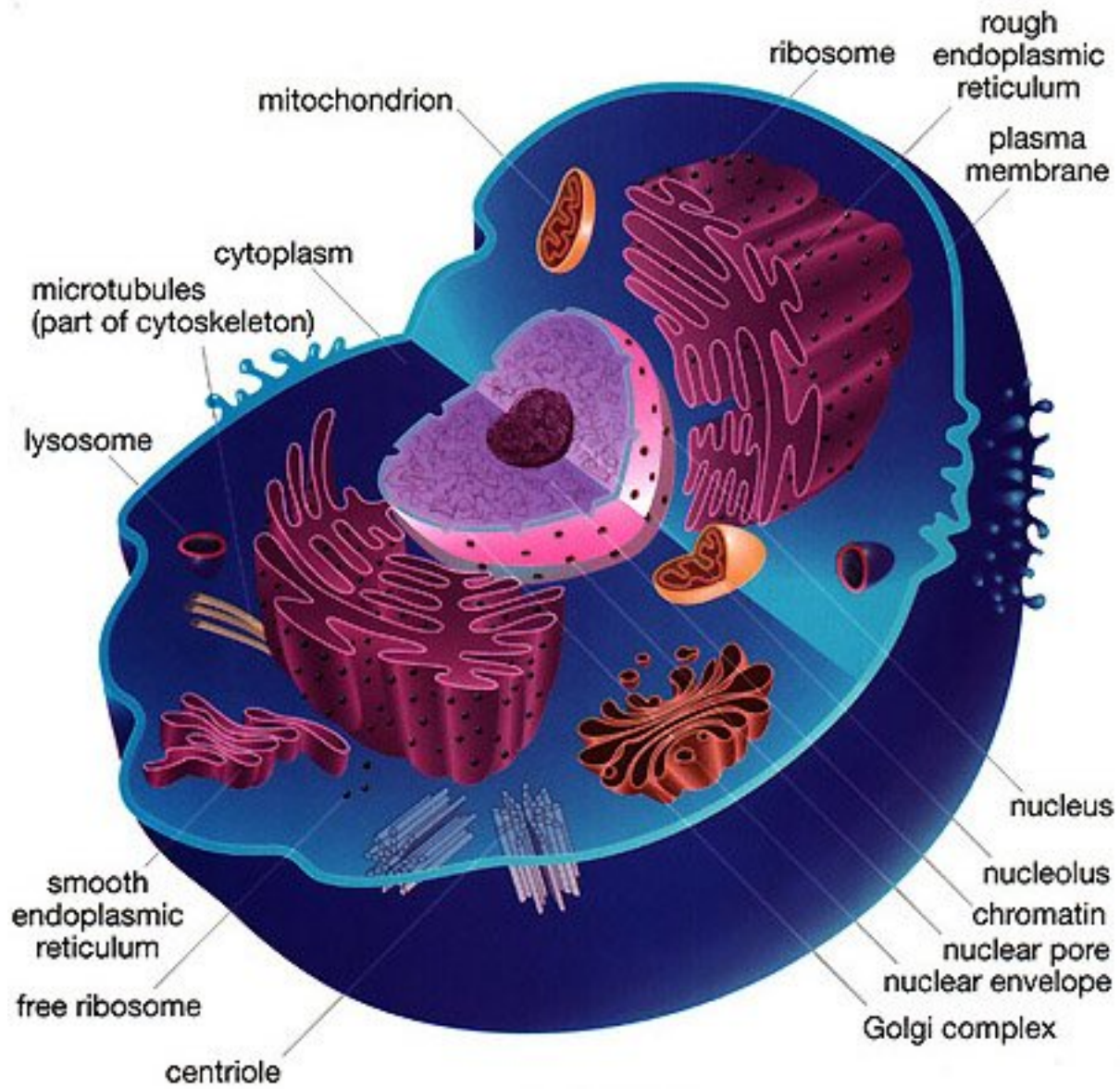There are so many cool studies that
I cannot cover here!

# Visit:

# http://goo.gl/opEfTp

# **Network** community analysis

Nodes

Links (edges) between nodes

Biochemical Pathways

Biochemical Pathways




Cell Science
Insulin Signalling
Paul Bevan

facebook

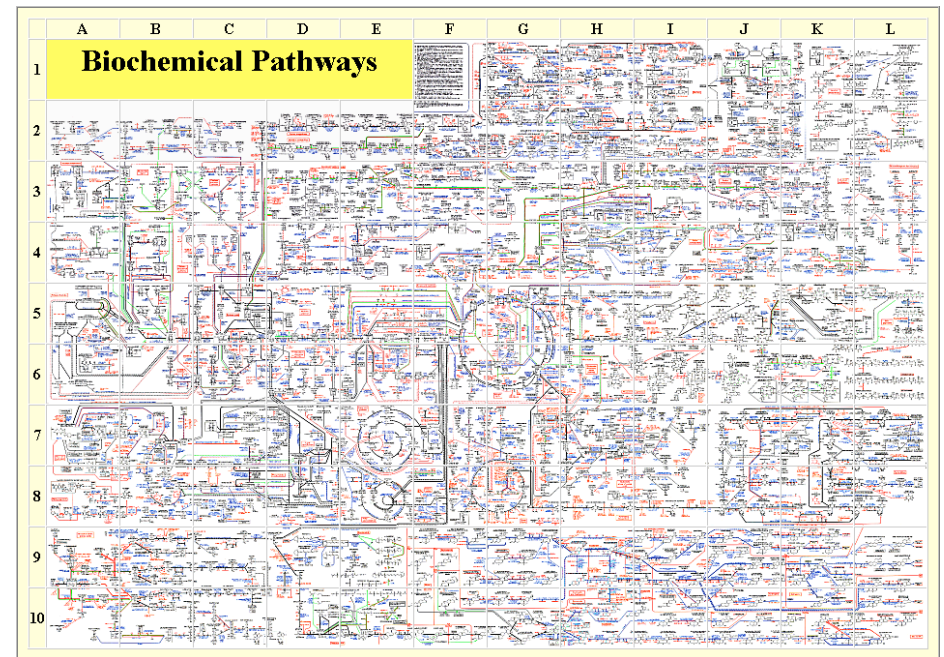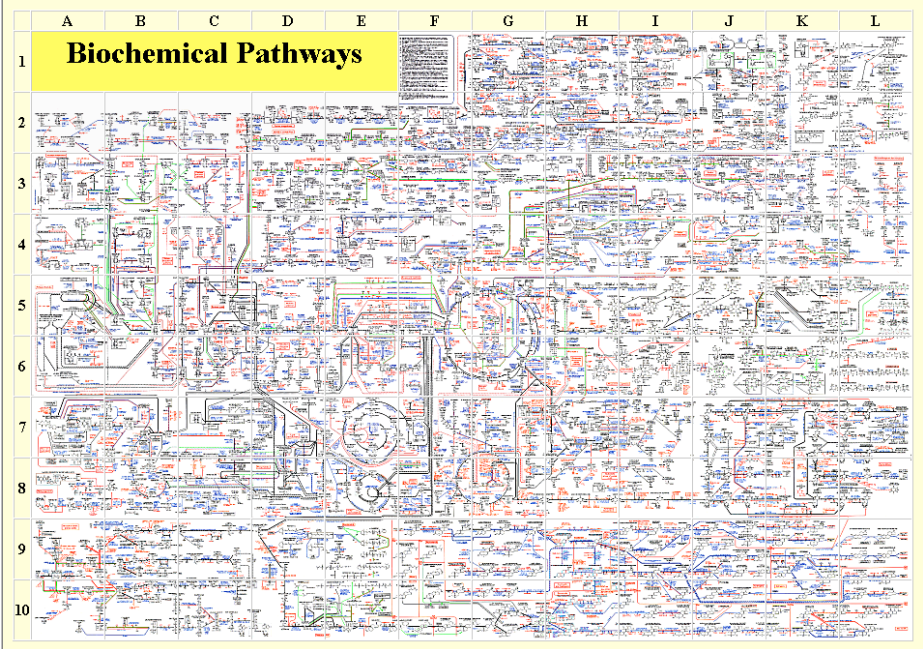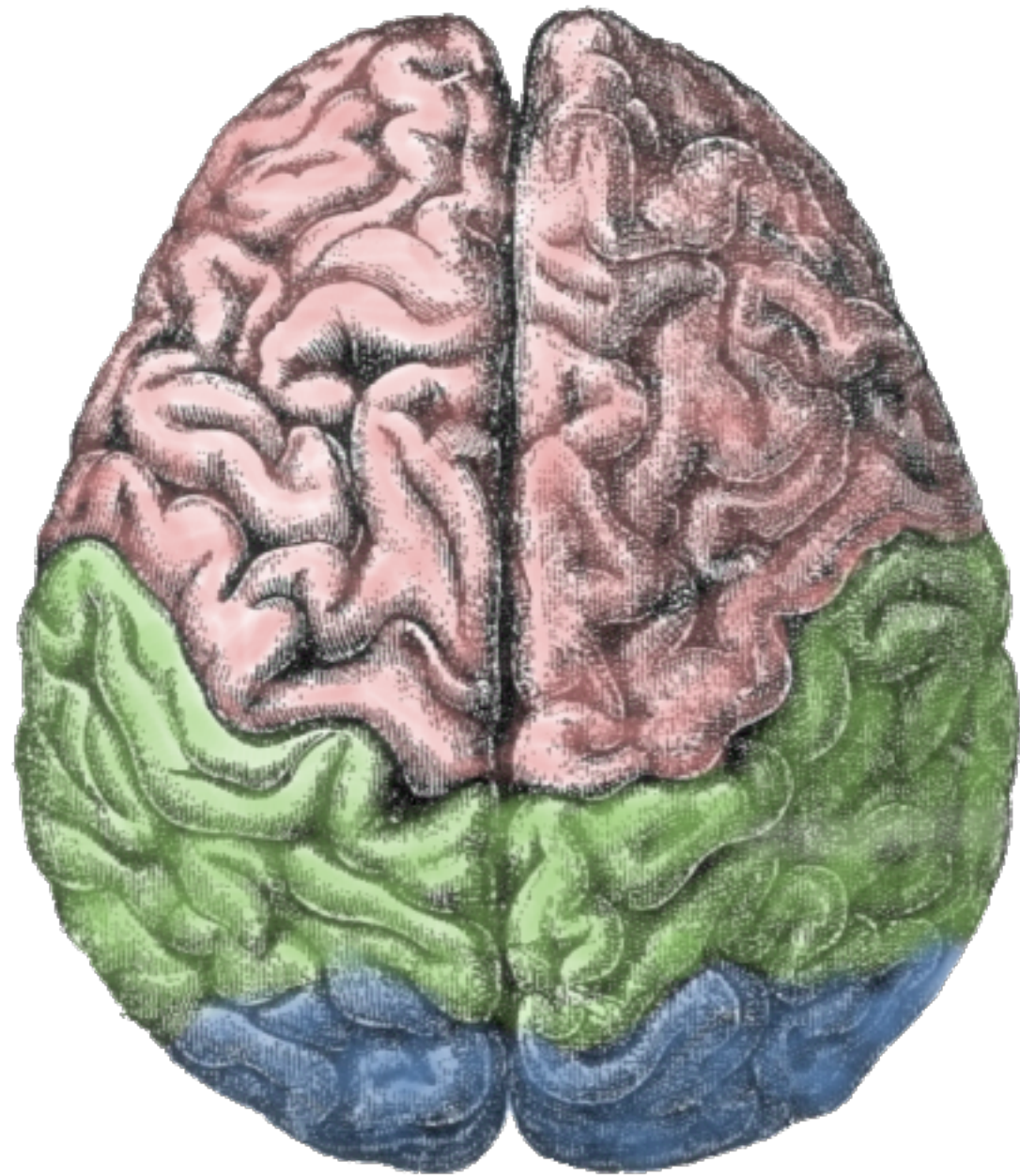December 2010
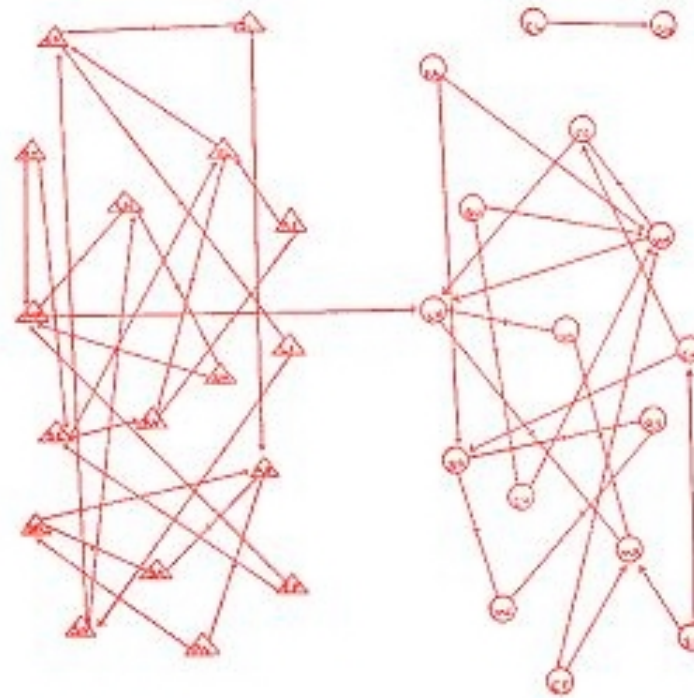
# Networks = The maps of complex systems

# Network **community** analysis

# What is a network community?

EMOTIONS MAPPED
BY NEW GEOGRAPHY

Charts Seek to Portray the
Psychological Currents of
Human Relationships.
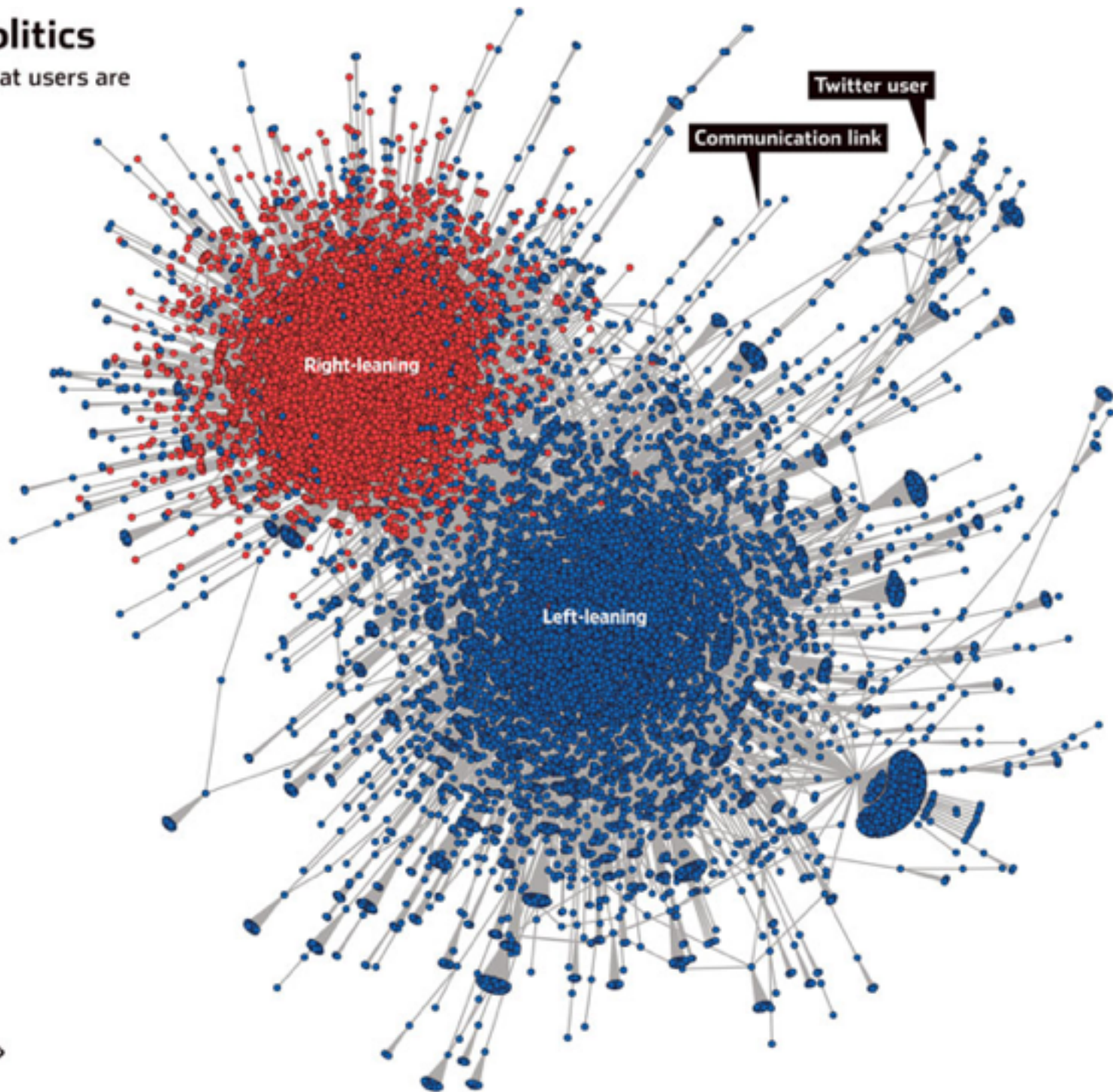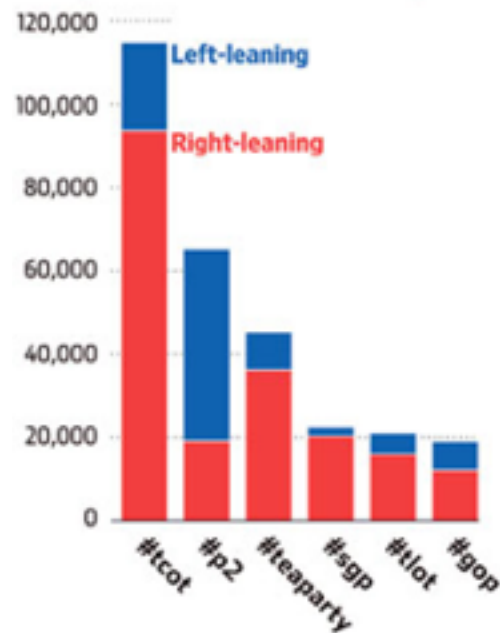
**New York Times**
April 3, 1933

# Moreno's "sociogram"

# Twitter's Divided Politics

## Political Twitter traffic reveals that users are polarized along party lines.*

Researchers at Indiana University analyzed 250,000 Twitter messages on political topics exchanged by 45,000 people during the 2010 mid-term congressional elections. This chart of 'retweets'—in which one user forwards another's message—shows that, though there were more left-leaning users, right-leaning users were more densely connected to one another. (Each dot is a Twitter user, and the lines show retweets.) Even so, as the chart illustrates, lines of communication do sometimes reach across the political divide.
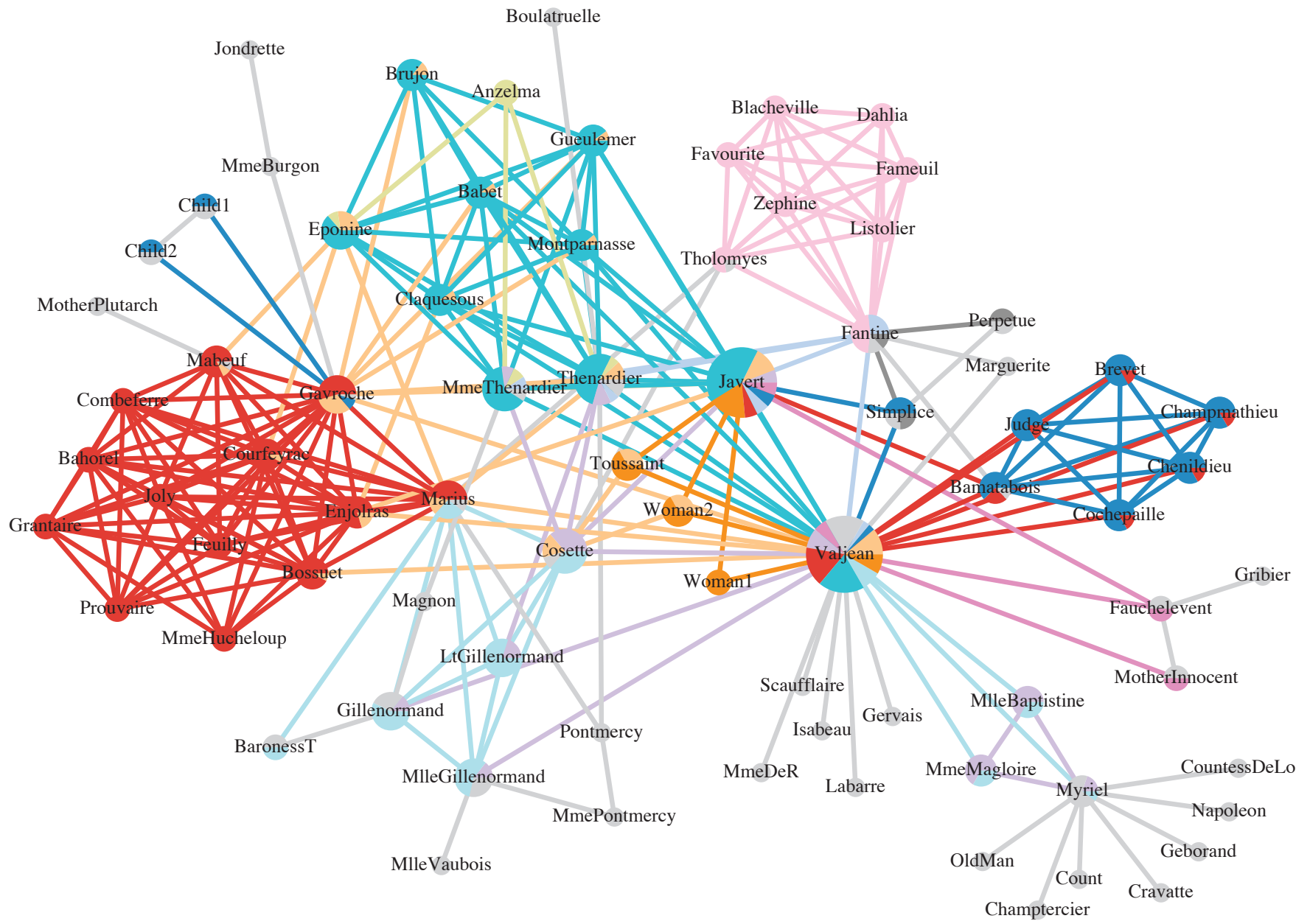
The most popular hashtags (short codes signaling the message's content), shown by number of tweets. Researchers found that users on the left and right use each other's hashtags.



Hashtags: tcot, top conservatives on Twitter; p2, progressives 2.0; sgp, smart girl politics; tlot, top libertarians on Twitter.
*Data show 'retweets' of other users' messages. Political leaning designations are based on algorithmically-determined communities of users which correlate with political affiliation.
Source: Center for Complex Networks and Systems Research, Indiana University
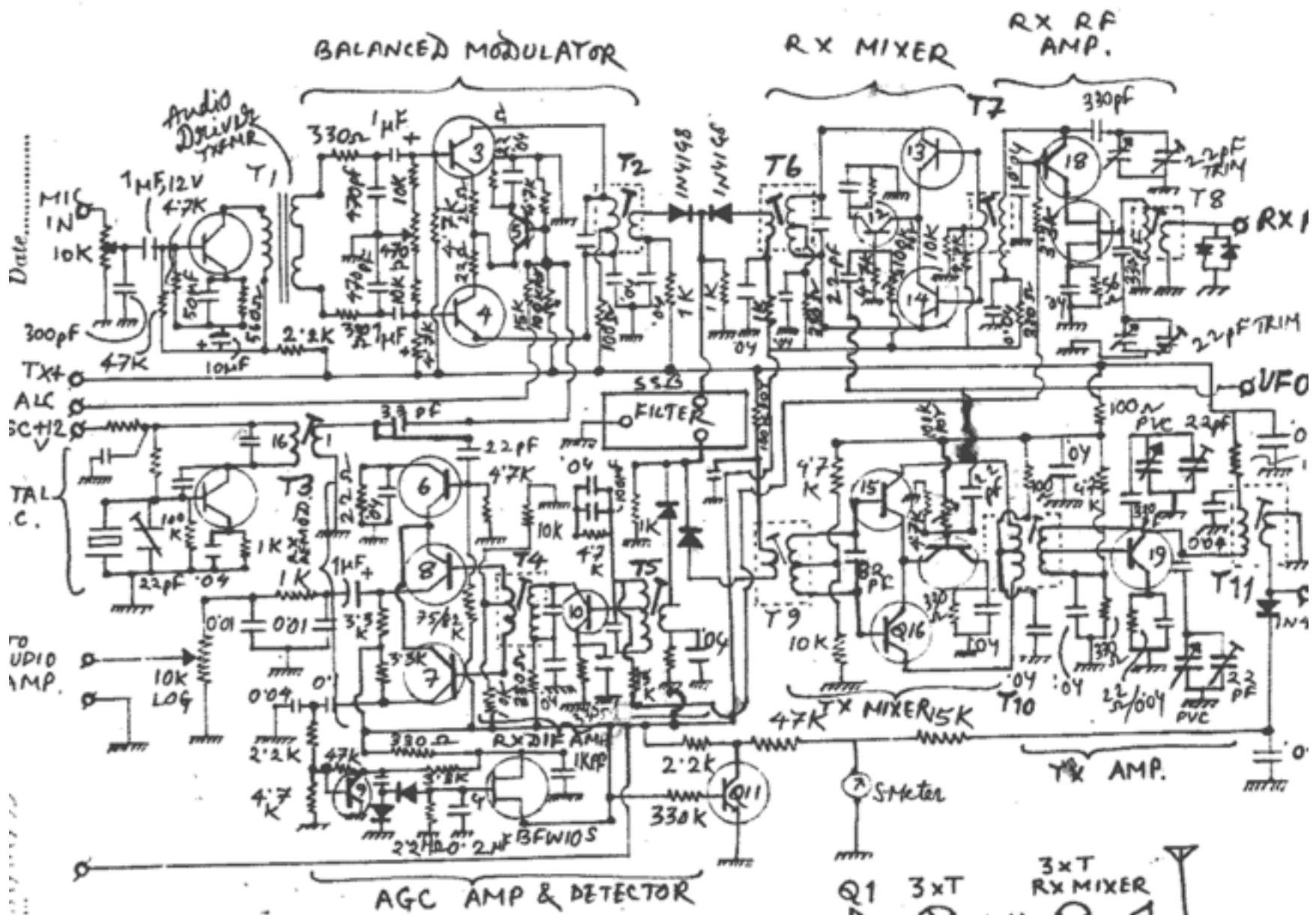
# Cohesiveness


# Separation

# Group **cohesiveness**
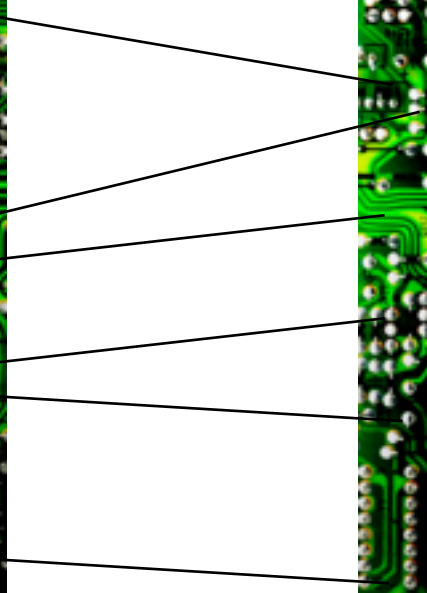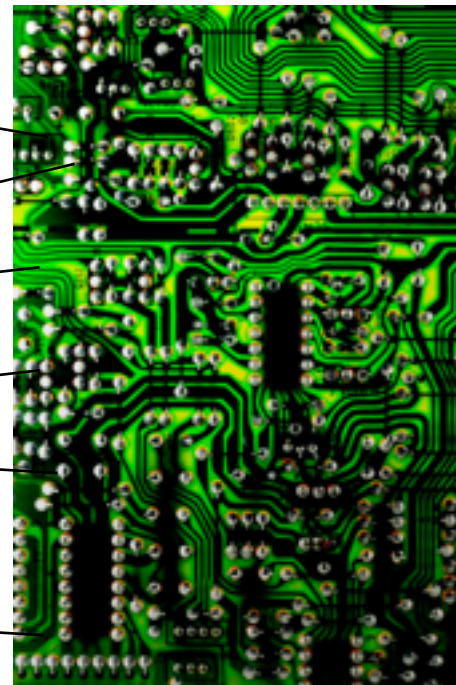
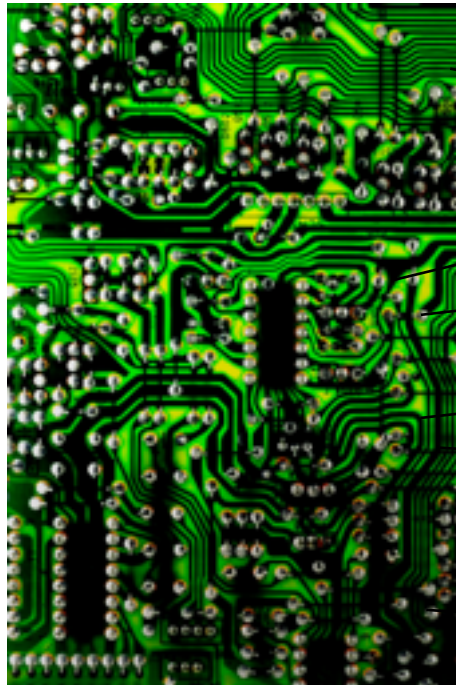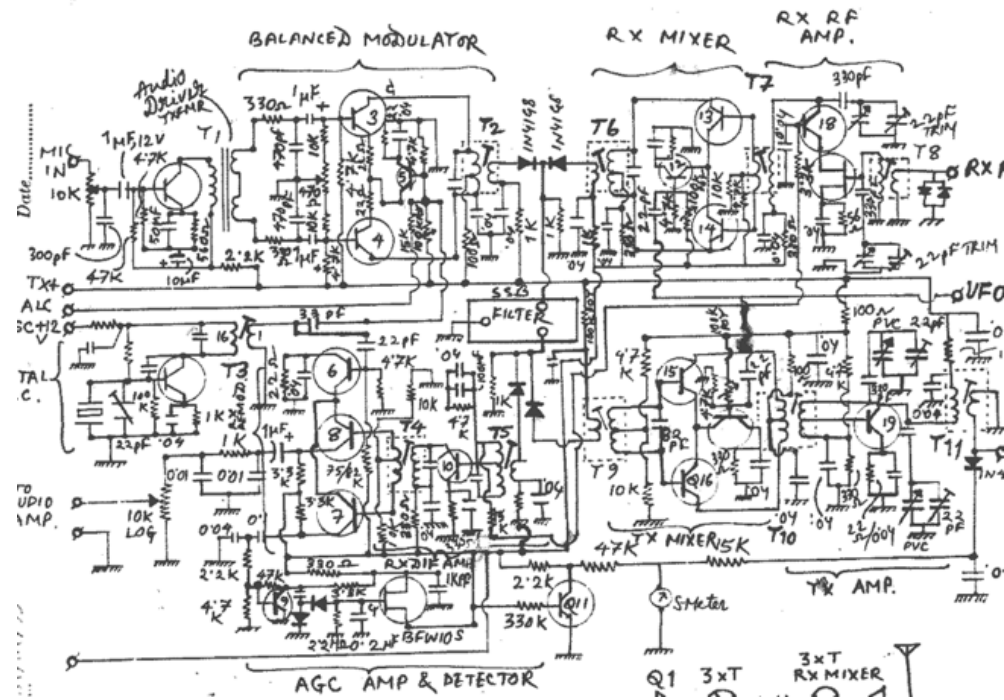(Moreno & Jennigs 1938, Festinger 1950, Gross & Martin 1952)

# Graph **partitioning**

(Kernighan & Lin 1970)

# Why do we care?

# Original motivation:
## Computation

BALANCED MODULATOR     RX MIXER     RX RF AMP.

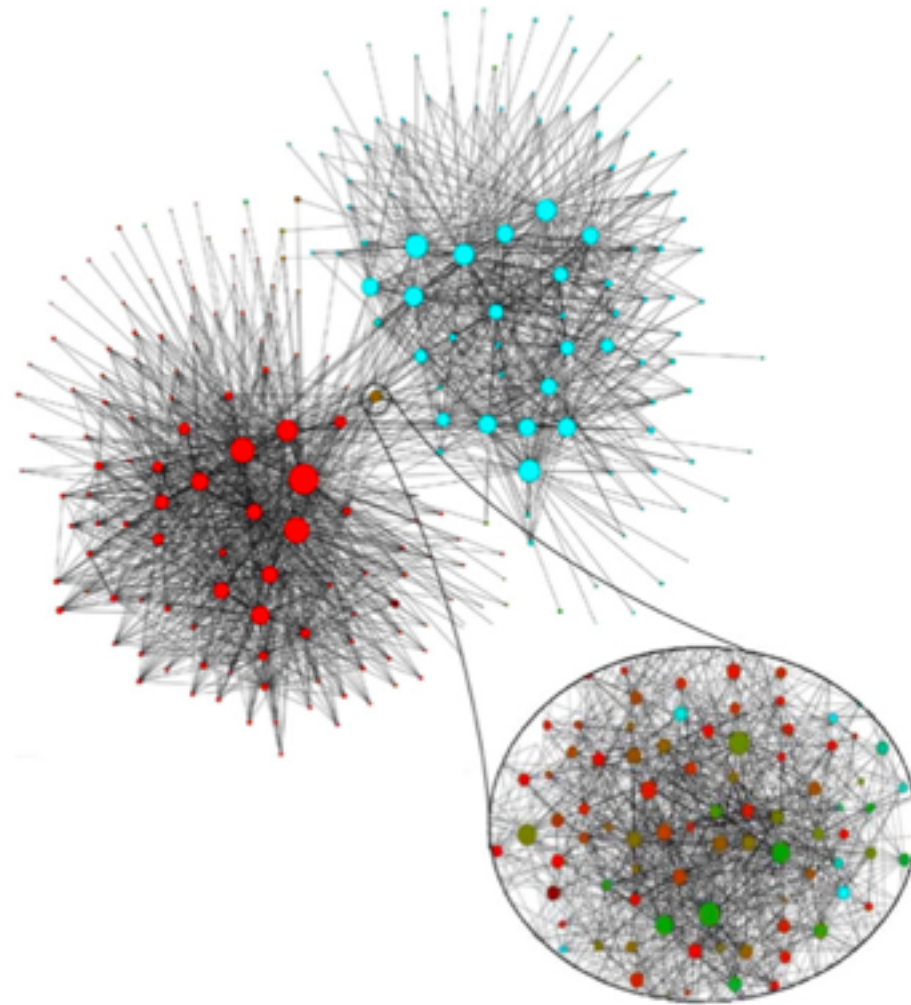AGC AMP & DETECTOR

# How to minimize the number of wires?

# How to minimize the communication between computers?

# Circuits, Communication between softwares ~ **Networks**
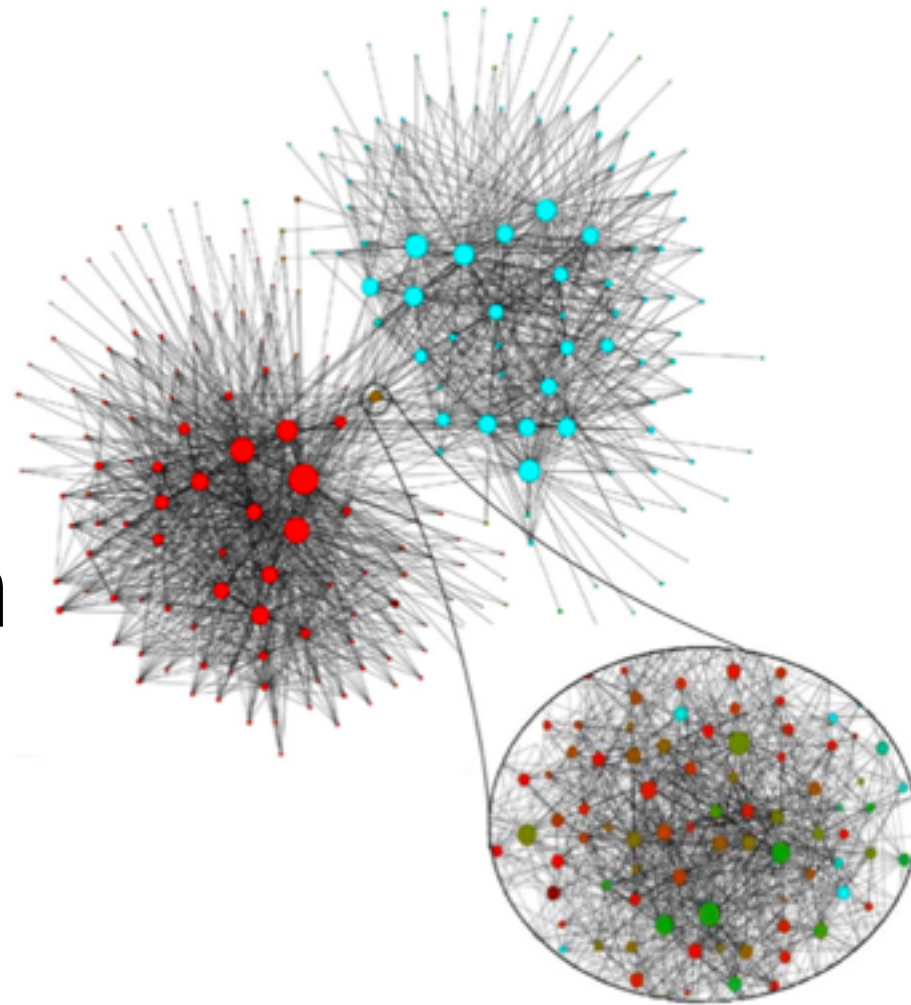

# Functional modules ~ **Communities**

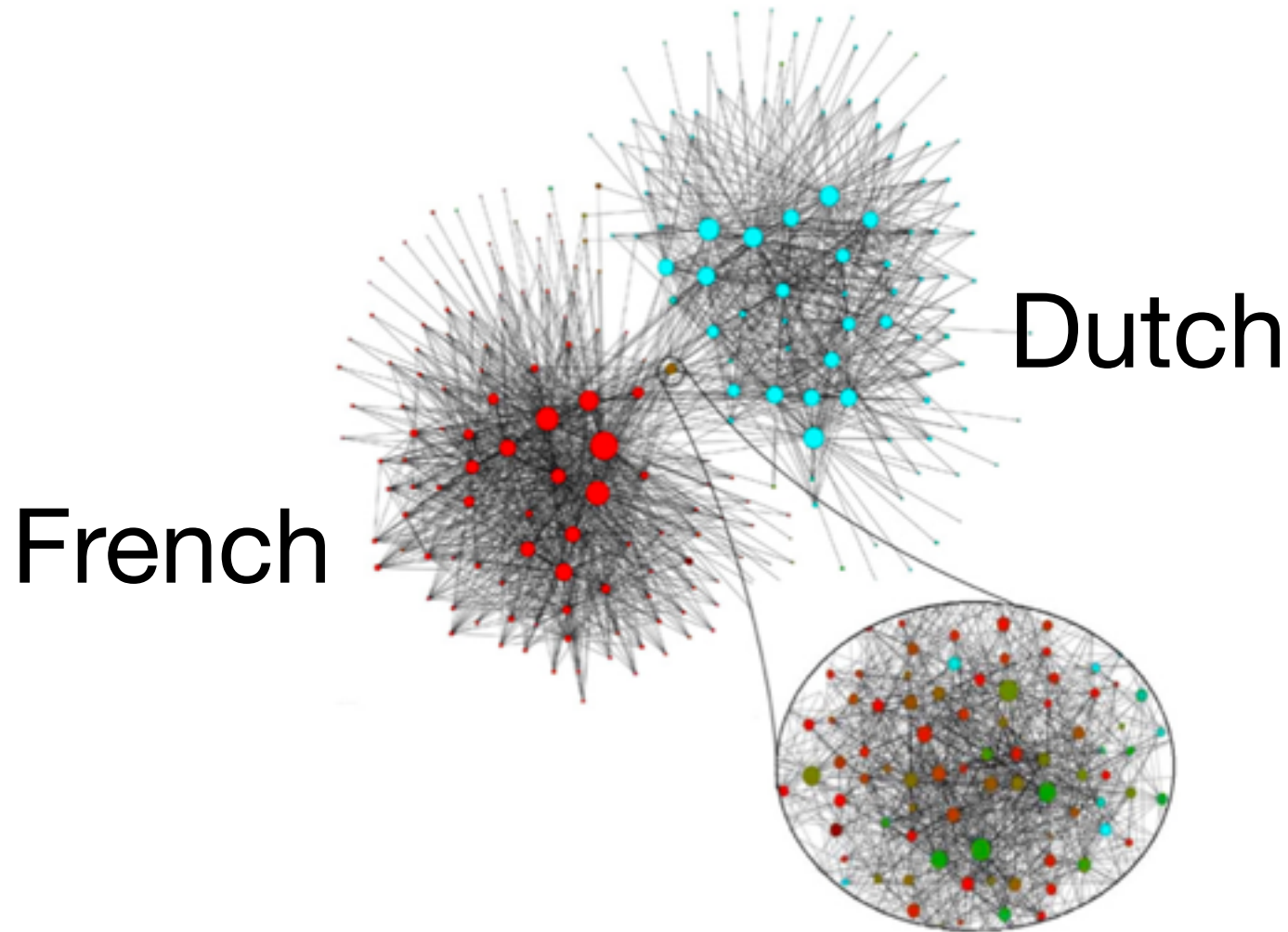# Correspondence to **functional, structural units**

# Belgian communication network



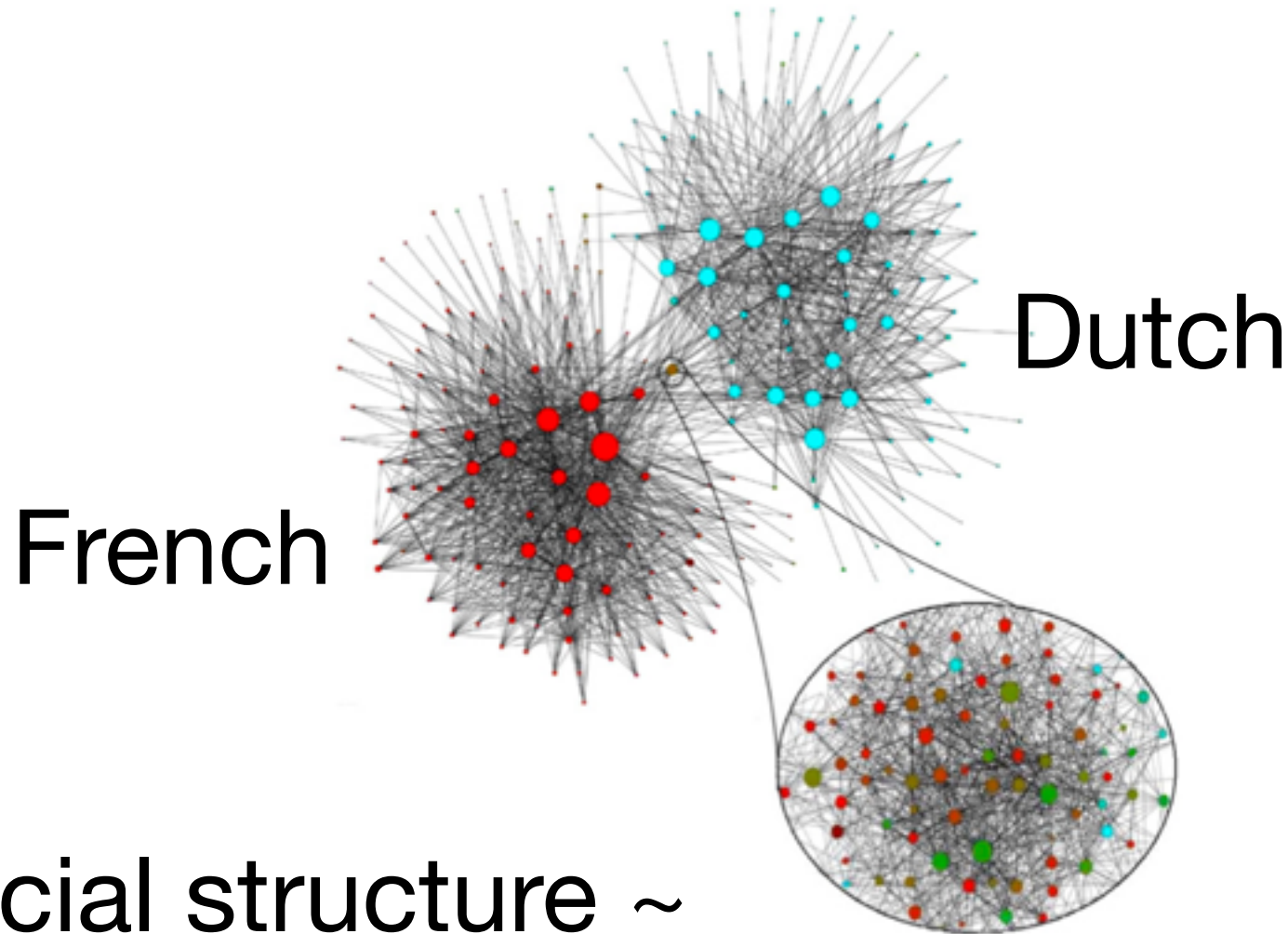V. D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre, *J. Stat. Mech.* (2008)

# Belgian communication network



French

V. D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre, *J. Stat. Mech.* (2008)

# Belgian communication network



Dutch

French

V. D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre, *J. Stat. Mech.* (2008)

# Belgian communication network



Dutch

French

## Social structure ~
**Communities**

V. D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre, *J. Stat. Mech.* (2008)

Y.-Y. Ahn, J. P. Bagrow, S. Lehmann, *Nature* (2010)

Protein complexes

~

**communities**

Y.-Y. Ahn, J. P. Bagrow, S. Lehmann, *Nature* (2010)

R. Guimerà & L. A. N. Amaral, *Nature* (2005)

Metabolic pathways ~ **communities**

R. Guimerà & L. A. N. Amaral, *Nature* (2005)

Social Sciences

Physical Sciences

Life Sciences

Ecology & Earth Sciences

M. Rosvall, C. T. Bergstrom, PLoS One (2011)

Disciplines ~ **communities**

M. Rosvall, C. T. Bergstrom, PLoS One (2011)

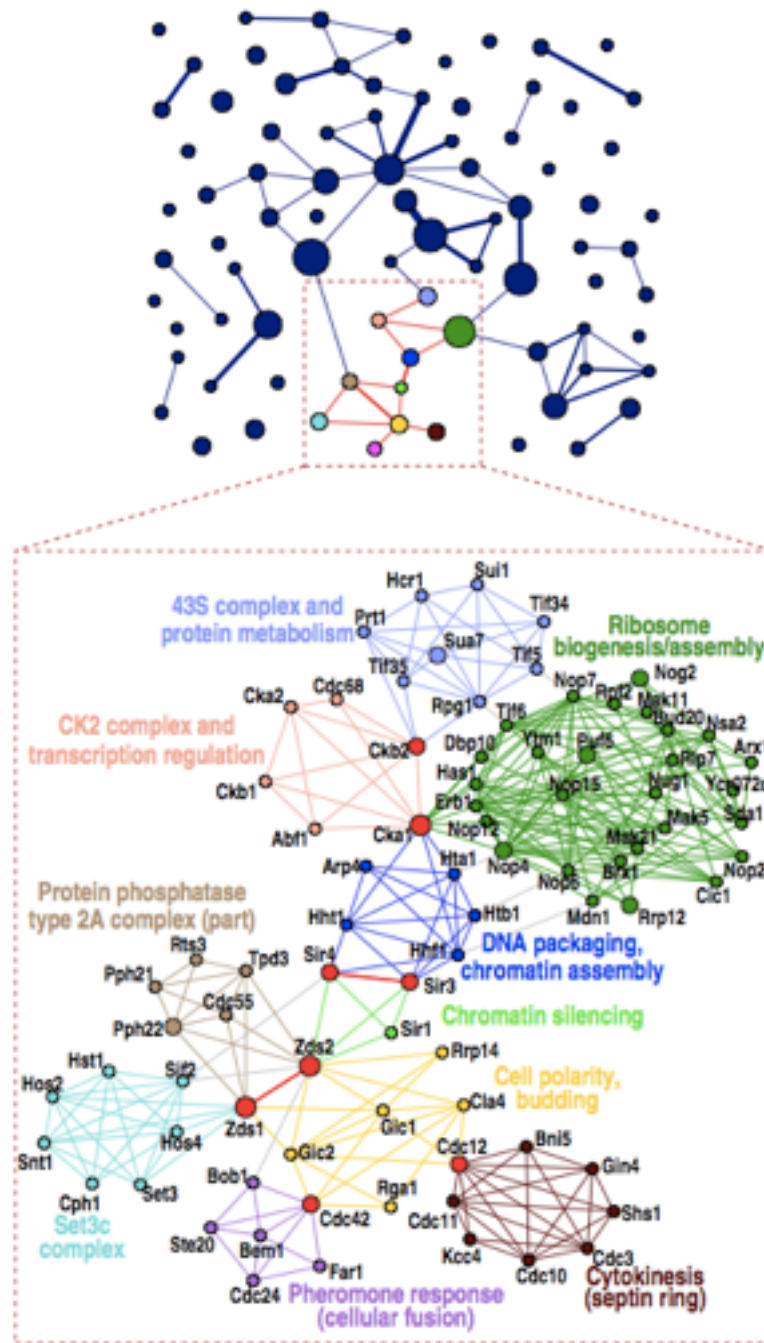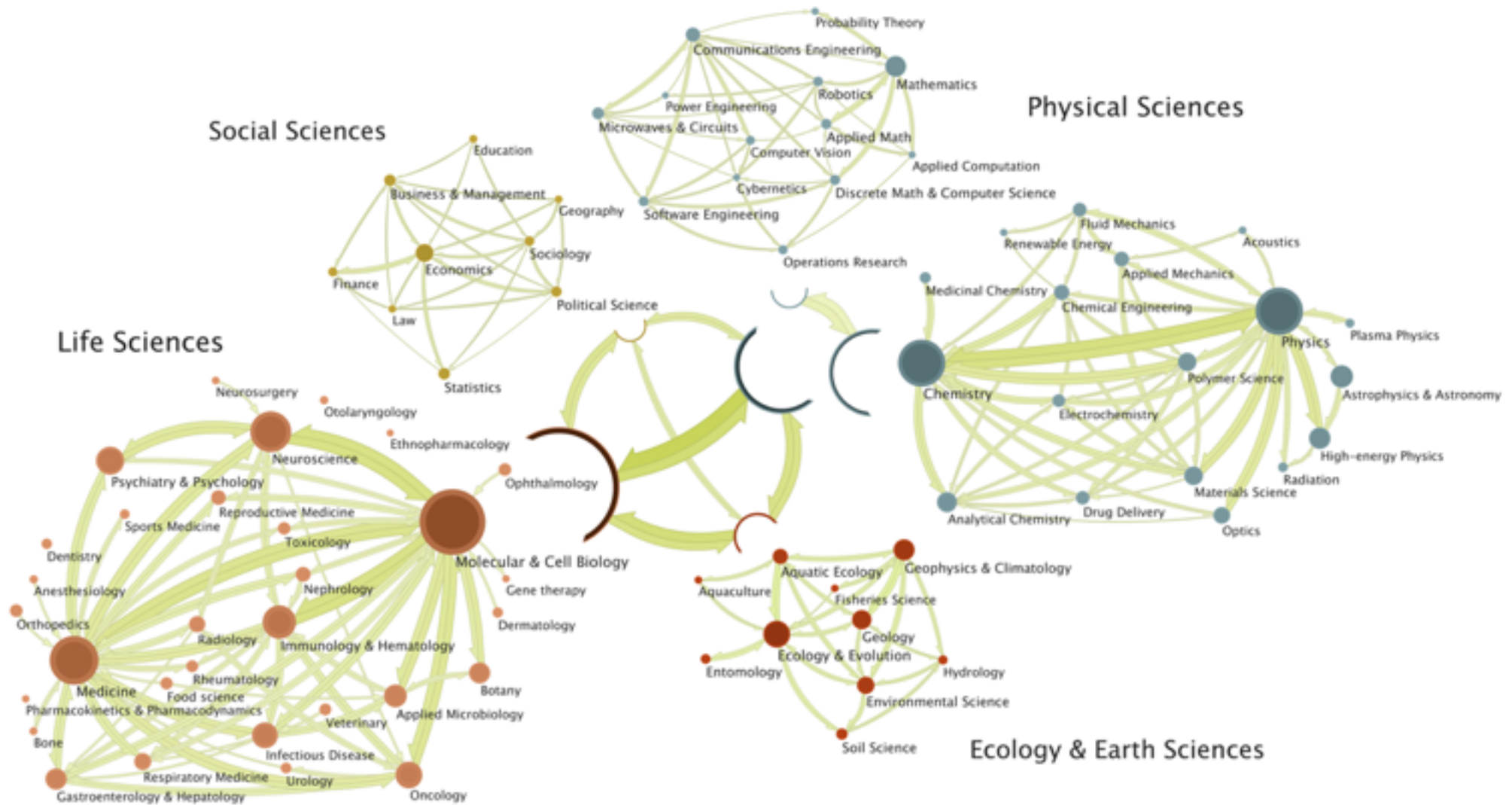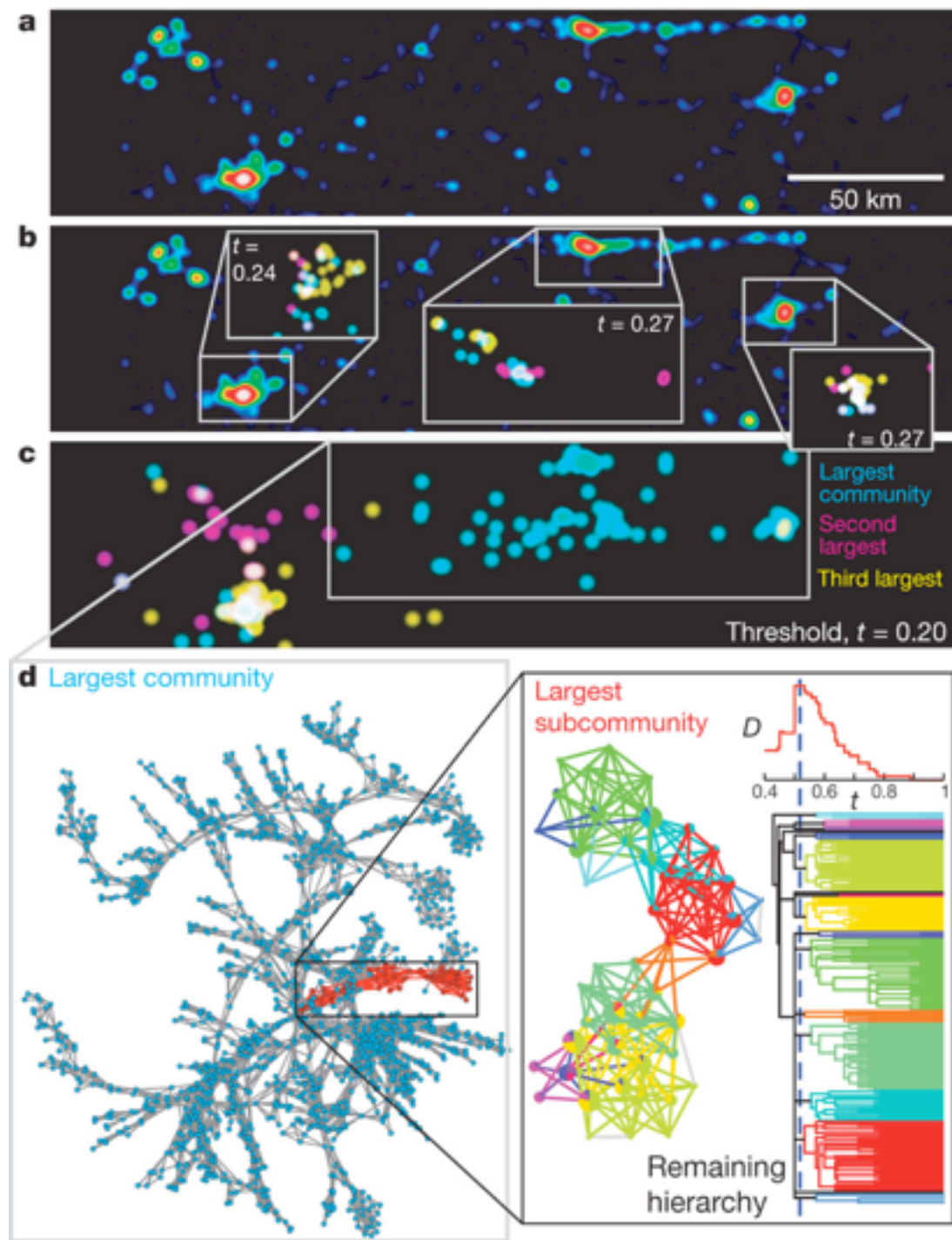| | |
|---|---|
| **Social Networks** | Social circles, communities |
| **Biological networks** | Protein complexes, functional modules |
| **Citation networks** | Disciplines, scientific communities |
| ... | ... |

Finding communities:
A nice way to
**overview**
the whole system

G. Palla, I. Derenyi, I. Farkas, T. Vicsek, *Nature* (2005).

M. Rosvall, C. T. Bergstrom, *PLoS One* (2011)

a

50 km

b

$t = 0.24$

$t = 0.27$

$t = 0.27$

c

Largest community

Second largest

Third largest

Threshold, $t = 0.20$

d Largest community

Largest subcommunity

$D$

0.4   0.6   0.8   1

$t$

Remaining hierarchy

Y.-Y. Ahn, J. P. Bagrow, S. Lehmann, *Nature* (2010)

# Network community **analysis**

# How to define **communities**?

# Cohesiveness, Separation


# or both


A nice review:
J. Yang and J. Leskovec, Defining and Evaluating Network Communities based on Ground-truth, ICDM 2012

# "Cohesiveness"

## Clique percolation
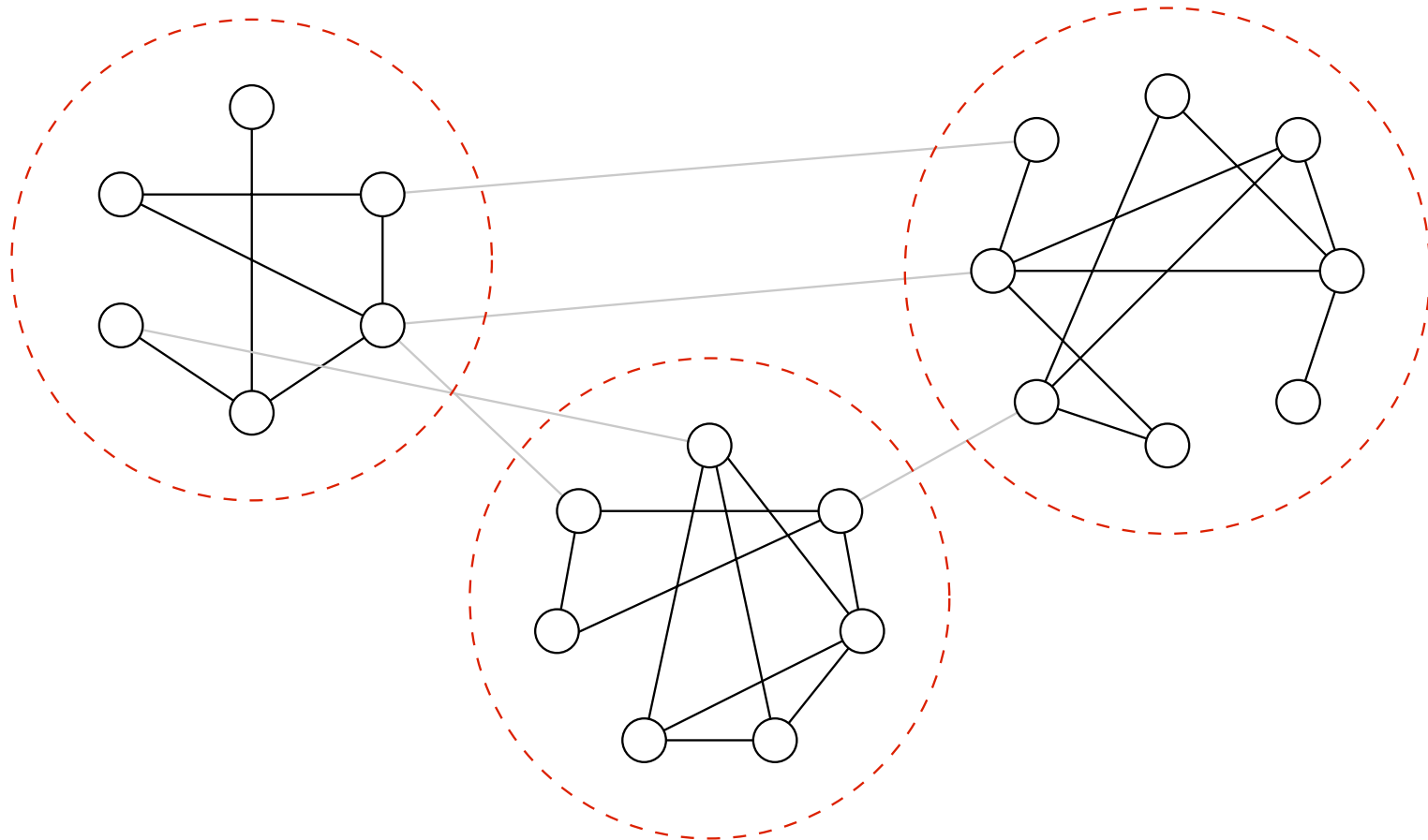## Link communities

# "Separation"

Girvan-Newman algorithm
Graph cuts
Spectral clustering

# Cohesiveness + Separation

# Modularity

$$Q = \frac{1}{2m} \sum_{ij} \left[ A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j)$$

# of links **within**
- # of **expected** links

M. Girvan and M. E. J. Newman, *PNAS* (2002)
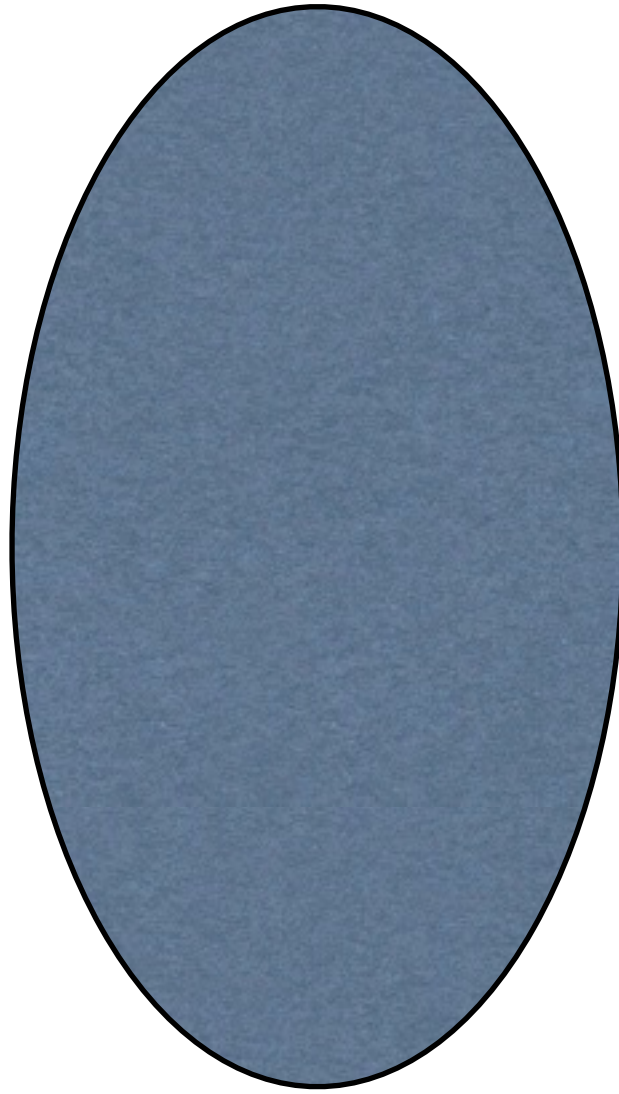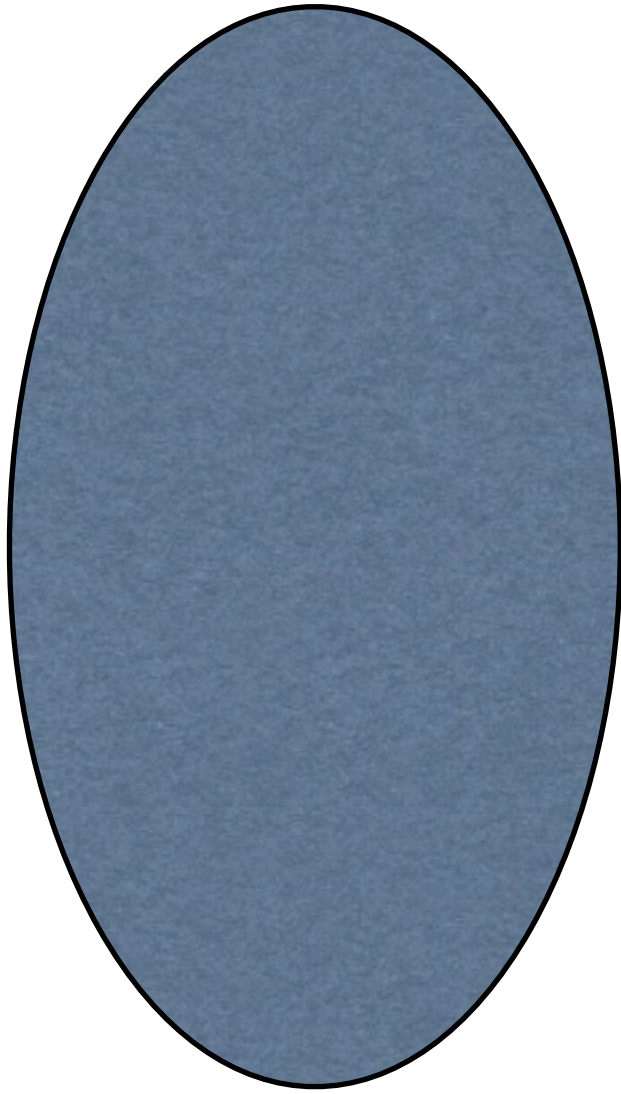
# How to detect communities?

# We should be able to

1. evaluate a community structure
2. explore possible structures effectively

**Wait**, can we just check **every possible** configurations?

# **Bell Number**: # of partitions of a set of size n.

# **Bell Number**: # of partitions of a set of size n.

$$B_3 = 5$$

**Bell Number**: # of partitions of a set of size n.

$$B_3 = 5$$

$$B_{100} = ?$$

16187060274460683058556806
28161135741330684513088812
39989840947008912873079240
70443511081340194490281914
80663320741161870602744606
83058556806281611357413306
84513088812399898409470089
12873079240704435110813401
94490281914806663320741

# Impossible to enumerate

# Fundamental problem of community detection

# 1. evaluate a community structure

- Modularity, cliques, map equation, partition density, …

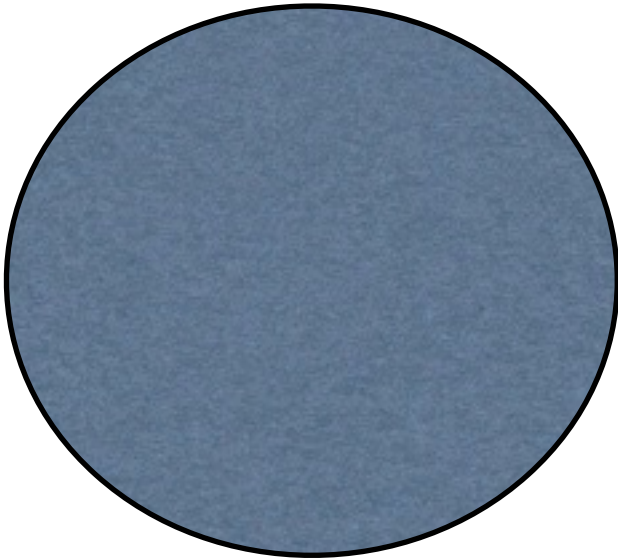# 2. explore possible structures effectively

- Many heuristics, Divisive & agglomerative clustering, Monte-carlo, …

# Modularity-based methods

# Divisive vs. Agglomerative

# Girvan-Newman algorithm
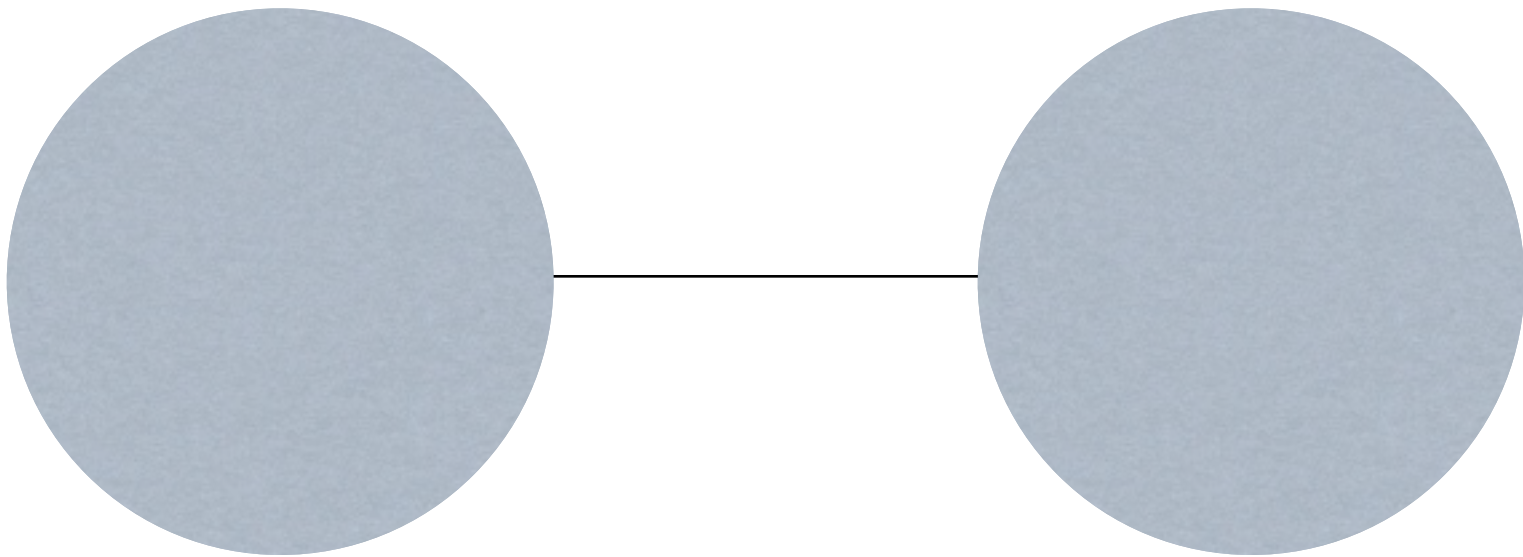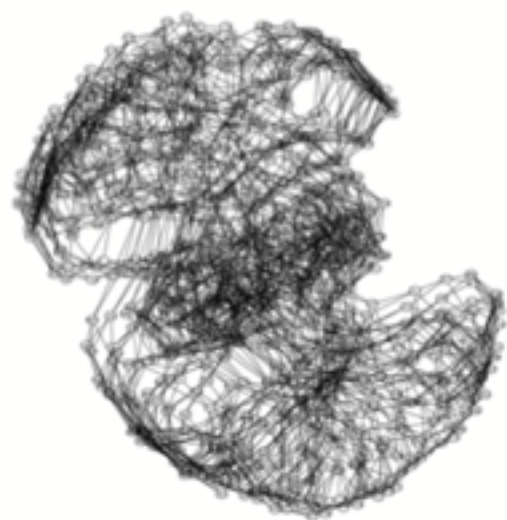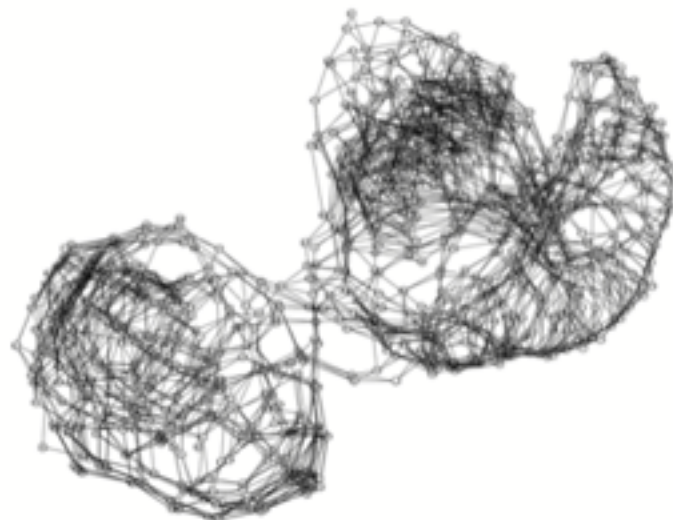
# Idea

0 cuts

100 cuts

120 cuts

500 cuts

# Louvain method



V. D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre JSTAT (2008).

# Various optimization techniques

- A. Clauset, M. Newman, C. Moore: Greedy optimization
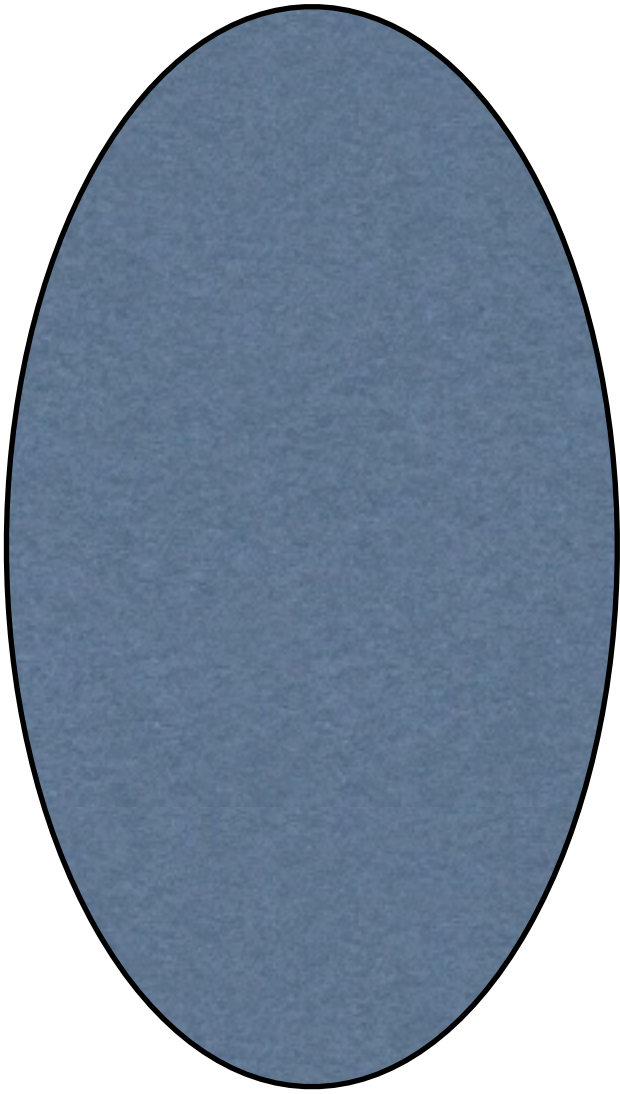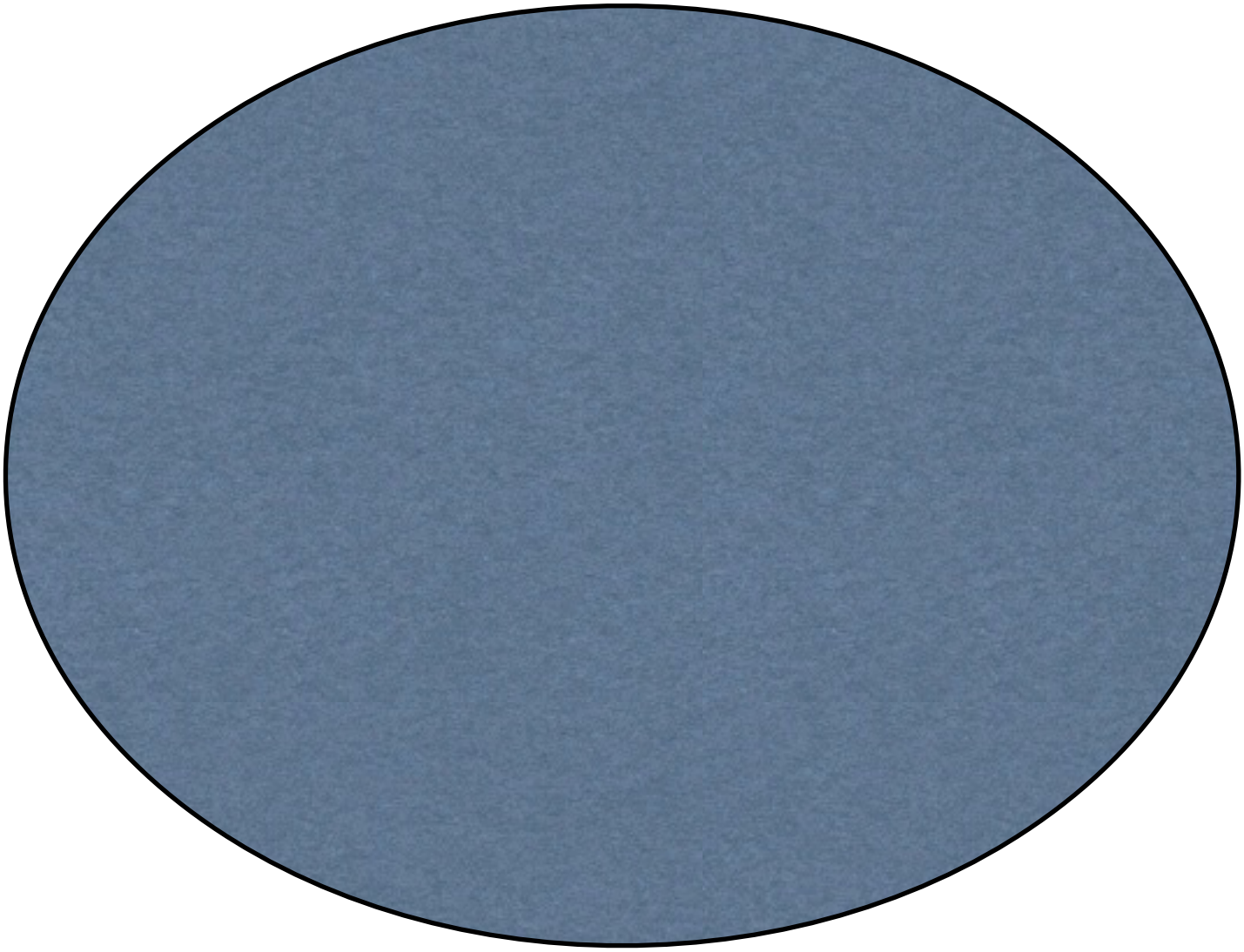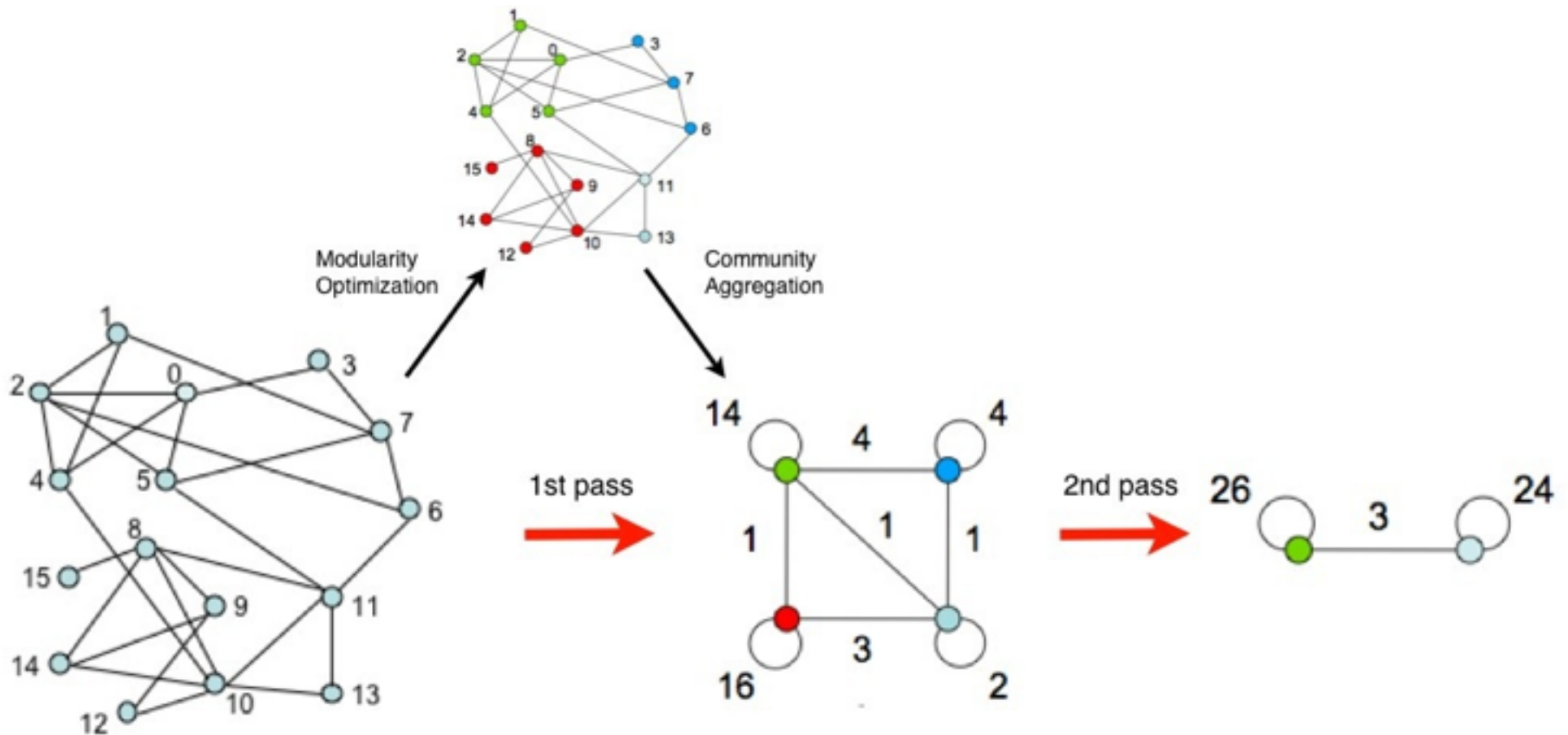
- R. Guimera, L. A. N. Amaral: Extremal optimization

- V. D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre: Hierarchical aggregation

- **Any** optimization technique can be used.

# "Cliques"

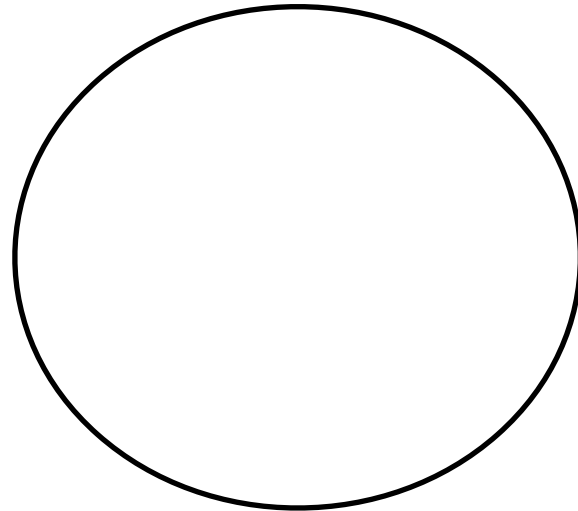# What is a 'perfect community'?

# A clique!

# Then, how about finding **quasi-cliques**?
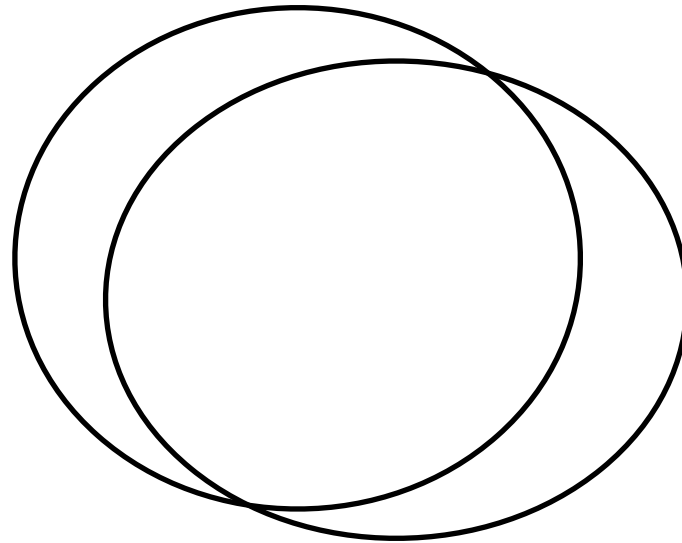
# Clique Percolation Method

- "Rolling" a clique to find a quasi-clique.

- Quasi-cliques are communities.

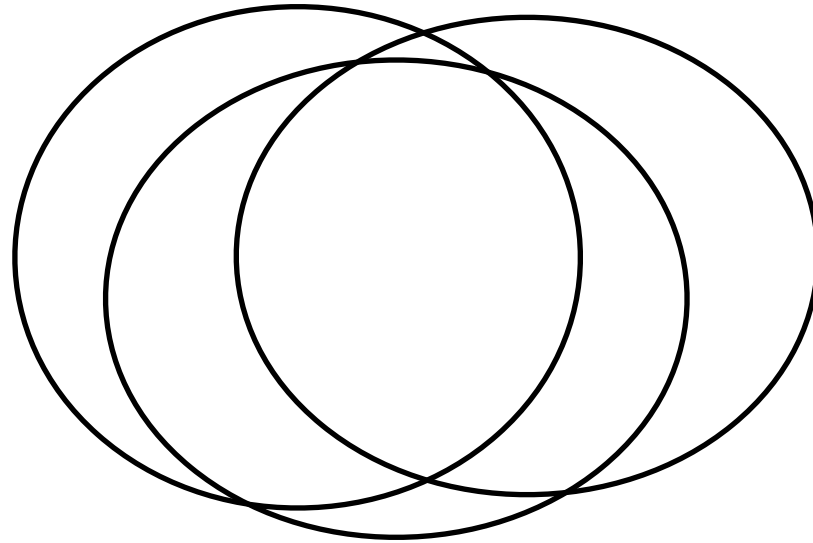G. Palla, I. Derenyi, I. Farkas, T. Vicsek, *Nature* (2005).

# Clique Percolation Method

- "Rolling" a clique to find a quasi-clique.

- Quasi-cliques are communities.

G. Palla, I. Derenyi, I. Farkas, T. Vicsek, *Nature* (2005).

# Clique Percolation Method



- "Rolling" a clique to find a quasi-clique.

- Quasi-cliques are communities.

G. Palla, I. Derenyi, I. Farkas, T. Vicsek, *Nature* (2005).

# Clique Percolation Method



- "Rolling" a clique to find a quasi-clique.

- Quasi-cliques are communities.

G. Palla, I. Derenyi, I. Farkas, T. Vicsek, *Nature* (2005).

# Clique Percolation Method



- "Rolling" a clique to find a quasi-clique.

- Quasi-cliques are communities.

G. Palla, I. Derenyi, I. Farkas, T. Vicsek, *Nature* (2005).
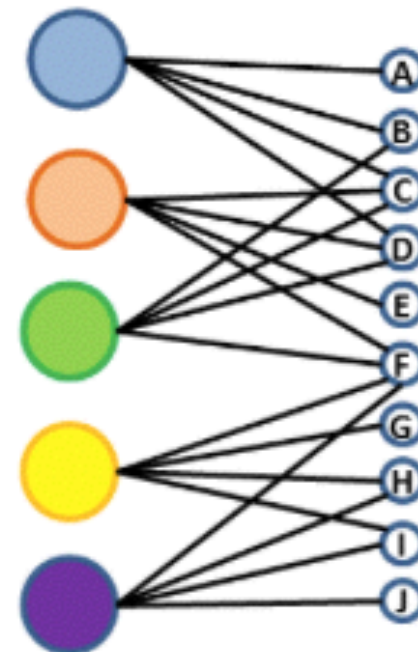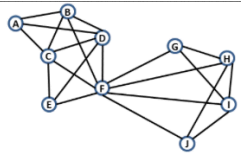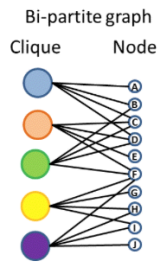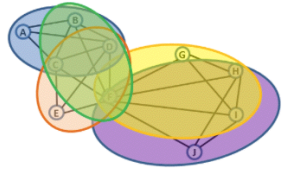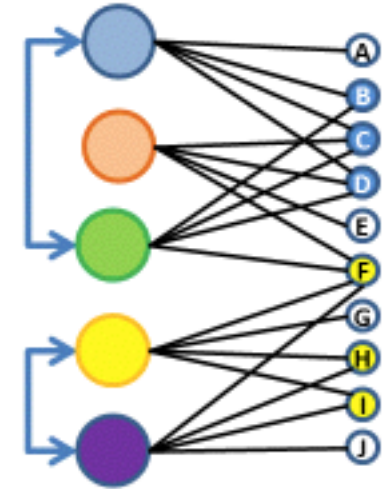
Original Graph

Step 1: Find all K-Cliques (K = 4)

Bi-partite graph

Clique          Node

Original Graph

Step 1: Find all K-Cliques (K = 4)
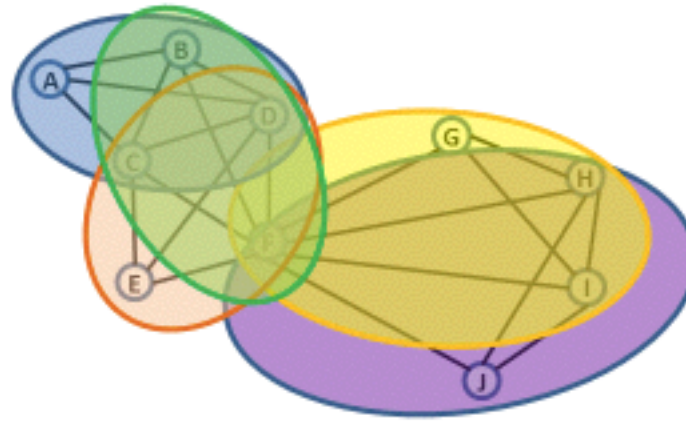
Bi-partite graph

Clique    Node

Step 2: Combine adjacent cliques (with K-1 = 3 shared nodes)

After merging adjacent cliques

After merging adjacent cliques
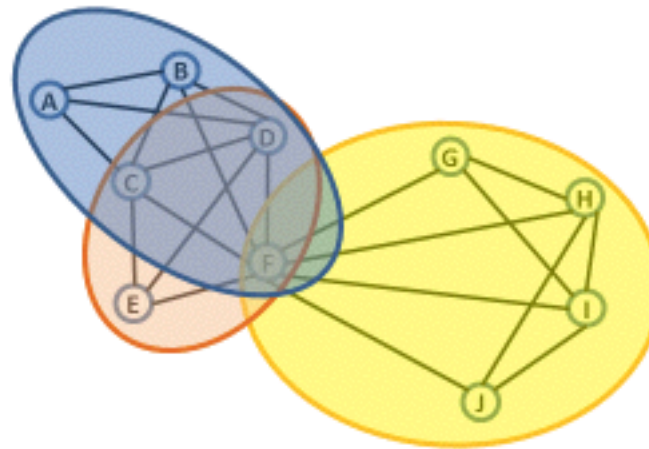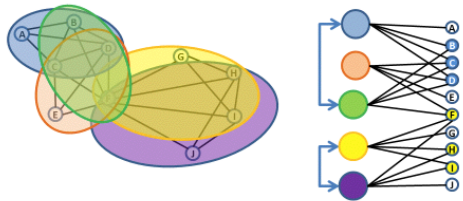
# Step 3: Combine adjacent cliques (with K-1 = 3 shared nodes)



# After merging adjacent cliques, there are 2 overlapping communities

# "Information"

If there is a **random walker** on the network, it will be **trapped** inside each community.

# *Demo*

# "Overlap"

Scientists
Physicists
Department of Biological Physics
Mathematicians
Biologists

'Zoom'

'Zoom'

Scientific community

Hobby

Family

Friends

Schoolmates

G. Palla, I. Derényi, I. Farkas & T. Vicsek, *Nature* (2005)

**Family**

buildings in same
neighborhood

**University**

home and work

joint appointment

**D**

1
2
3
4
7

**F**

3–4

2–4

1–4

It is **impossible** to obtain a single dendrogram.

# Simple local structure

# Complex global structure

# Complex global structure

This is a modular network.

# What is this?

What the xxxx is this?

**Word association network**: Network of  "commonly associated English words"



G. Palla, I. Derényi, I. Farkas & T. Vicsek, *Nature*, 2005

# **Link** communities

Colleagues

Family

Friends

Colleagues

'Family' links

Family

Friends

Colleagues

'Family' links

Friends

Family

'Friends' links

Nodes: multiple membership

Links: unique membership

# Similarity between links

$\downarrow$

# Hierarchical Clustering

$$n_+(i) \equiv \{x \mid d(i,x) \leq 1\}$$

$$S(e_{ik}, e_{jk}) = \frac{|n_+(i) \cap n_+(j)|}{|n_+(i) \cup n_+(j)|}$$

$$S(e_{ac}, e_{bc})$$

$$n_+(i) \equiv \{x \mid d(i, x) \leq 1\}$$

$$S(e_{ik}, e_{jk}) = \frac{|n_+(i) \cap n_+(j)|}{|n_+(i) \cup n_+(j)|} \quad \frac{4}{12}$$

# Partition Density

Community $c$ has $m_c$ edges and $n_c$ induced nodes

# Partition Density

Community $c$ has $m_c$ edges and $n_c$ induced nodes

# Partition Density

Community $c$ has $m_c$ edges and $n_c$ induced nodes



$m_c = 8$

$n_c = 5$

# Partition Density

Community $c$ has $m_c$ edges and $n_c$ induced nodes

 $= m_c$

# Partition Density

Community $c$ has $m_c$ edges and $n_c$ induced nodes

$$\frac{\text{(graph)}}{\text{(complete graph)}} = \frac{m_c}{\frac{n_c(n_c-1)}{2}}$$

# Partition Density

Community $c$ has $m_c$ edges and $n_c$ induced nodes

$$\frac{\ \ }{\ \ } = \frac{m_c}{\frac{n_c(n_c-1)}{2}}$$

A **single** link is maximally dense

# Partition Density

Community $c$ has $m_c$ edges and $n_c$ induced nodes



$$\frac{\text{(graph)} - \text{(graph)}}{\text{(graph)} - \text{(graph)}} = \frac{m_c - (n_c - 1)}{\frac{n_c(n_c - 1)}{2} - (n_c - 1)}$$

# Partition Density



$$\frac{\phantom{xxx} - \phantom{xxx}}{\phantom{xxx} - \phantom{xxx}} = \frac{m_c - (n_c - 1)}{\frac{n_c(n_c - 1)}{2} - (n_c - 1)}$$

$$= 2\frac{m_c - (n_c - 1)}{(n_c - 2)(n_c - 1)}$$

# Partition Density



$$\frac{\,-\,}{\,-\,} = \frac{m_c - (n_c - 1)}{\frac{n_c(n_c-1)}{2} - (n_c - 1)}$$

$$= 2\frac{m_c - (n_c - 1)}{(n_c - 2)(n_c - 1)}$$

$$D \equiv \frac{2}{M}\sum_c m_c \frac{m_c - (n_c - 1)}{(n_c - 2)(n_c - 1)}$$

**Experiment, Science**

FLASK · CHEMICAL · BEAKER · TESTTUBE · LAB · BIOLOGY · BIOLOGIST · RESEARCH · CHEMIST · CHEMISTRY · EXPERIMENT · SCIENTIST

**Smart, Intellect, Scientists**

INVENT · EXCEPTIONAL · BRIGHT · BRILLIANT · INVENTOR · INTELLECT · GENIUS · GIFTED · INTELLIGENT · INTELLIGENCE · SMART · WISDOM · RETARDED

**Science, Scientists**

KINETIC · VELOCITY · SCIENCE · SCIENTIFIC · PHYSICS · HYPOTHESIS · EINSTEIN · THEORY · THEOREM · RELATIVITY · INERTIA · LAW

**Newton, Gravity, Apple**

NEWTON · WEIGHT · GRAVITY · APPLE

**Clever, Wit**

WISE · CUNNING · CLEVER · WIT · SLY · OUTFOX

Mediator

YAP1

NuA4

TRA1, known
subunit of both
NuA4- and SAGA
complex

SAGA

# The first plant (genomic scale) interactome



Arabidopsis Interactome Mapping Consortium, *Science, 2011*

# The first plant (genomic scale) interactome



Arabidopsis Interactome Mapping Consortium, *Science, 2011*

# Statistical inference

Given a graph *G*, and a generative
model with parameters $\theta$

Likelihood

$$P(\theta|G) = \frac{P(G|\theta)P(\theta)}{P(G)}$$

# Stochastic Block Model

# Stochastic Block Model



$$c_i$$

# Stochastic Block Model



$$c_i \qquad p_{c_i c_j}$$

# Stochastic Block Model



$$c_i \qquad p_{c_i c_j}$$

$$\prod_{i<j} p_{c_i c_j}^{A_{ij}} (1 - p_{c_i c_j})^{1 - A_{ij}}$$

# So, what should I use?

1. No silver bullet.

2. Hard to know beforehand.

# Accuracy

# Accuracy

- It is very hard to compare performance of methods on a fair ground because each method is usually very good at finding what it is looking for.

# Accuracy

- It is very hard to compare performance of methods on a fair ground because each method is usually very good at finding what it is looking for.

- Most studies use '**benchmark networks**' to evaluate the performance.

# Accuracy

- It is very hard to compare performance of methods on a fair ground because each method is usually very good at finding what it is looking for.

- Most studies use '**benchmark networks**' to evaluate the performance.

- **Infomap** and **Louvain** method are the best in these benchmarks.

# Accuracy

- It is very hard to compare performance of methods on a fair ground because each method is usually very good at finding what it is looking for.

- Most studies use '**benchmark networks**' to evaluate the performance.

- **Infomap** and **Louvain** method are the best in these benchmarks.

- However, the performance depends on what kinds of community structure the benchmark networks assume.

# Accuracy

- It is very hard to compare performance of methods on a fair ground because each method is usually very good at finding what it is looking for.

- Most studies use '**benchmark networks**' to evaluate the performance.

- **Infomap** and **Louvain** method are the best in these benchmarks.

- However, the performance depends on what kinds of community structure the benchmark networks assume.

- Good performance in the benchmarks *does not guarantee* good performance in real cases.

# Computational complexity

# Computational complexity

- Some methods are much faster than others.

# Computational complexity

- Some methods are much faster than others.
- *O(exp(n)) vs. O(m^2 n) vs. O(n log n)*

# Computational complexity

- Some methods are much faster than others.
- *O(exp(n))* vs. *O(m^2 n)* vs. *O(n log n)*
- Usually a good choice for huge (> 1m ~ 1b) networks: **Louvain method** (*~O(n log n)*)

# Computational complexity

- Some methods are much faster than others.
- *O(exp(n))* vs. *O(m^2 n)* vs. *O(n log n)*
- Usually a good choice for huge (> 1m ~ 1b) networks: **Louvain method** (*~O(n log n)*)
- Current version of **infomap** also uses louvain-type multilevel optimization and very fast.

# Computational complexity

- Some methods are much faster than others.
- *O(exp(n))* vs. *O(m^2 n)* vs. *O(n log n)*
- Usually a good choice for huge (> 1m ~ 1b) networks: **Louvain method** (*~O(n log n)*)
- Current version of **infomap** also uses louvain-type multilevel optimization and very fast.
- **Link clustering** can also handle large graphs (but it becomes slow with large hubs).

# Overlap

# Overlap

- If you expect pervasive overlap of communities, you should use overlapping community detection methods.

# Overlap

- If you expect pervasive overlap of communities, you should use overlapping community detection methods.
- **Link clustering** and **clique percolation** methods are common choices.

# Overlap

- If you expect pervasive overlap of communities, you should use overlapping community detection methods.
- **Link clustering** and **clique percolation** methods are common choices.
- These methods can detect highly overlapping communities. There are many other methods but most methods only deal with '**fuzzy**' overlaps.

# Resolution limit

# Resolution limit



- Modularity has a resolution limit that depends on the system size.

# Resolution limit



- Modularity has a resolution limit that depends on the system size.
- If a community is smaller than this limit, modularity-based optimization cannot find the communities, even though they are cliques.

# My heuristic

# My heuristic

- I don't care too much and I just want to get rough clusters in my network — **Infomap** or **Louvain**

# My heuristic

- I don't care too much and I just want to get rough clusters in my network — **Infomap** or **Louvain**
- I expect multiple community membership for many nodes — **Link clustering** (**Clique Percolation**)

# My heuristic

- I don't care too much and I just want to get rough clusters in my network — **Infomap** or **Louvain**
- I expect multiple community membership for many nodes — **Link clustering** (**Clique Percolation**)
- My network is HUGE and doesn't have super-large hubs — **Louvain** (**Infomap, link clustering**)

# My heuristic

- I don't care too much and I just want to get rough clusters in my network — **Infomap** or **Louvain**
- I expect multiple community membership for many nodes — **Link clustering** (**Clique Percolation**)
- My network is HUGE and doesn't have super-large hubs — **Louvain** (**Infomap, link clustering**)
- My network is HUGE and has lots of super-hubs — **Louvain** (**Infomap**)

# My heuristic

- I don't care too much and I just want to get rough clusters in my network — **Infomap** or **Louvain**
- I expect multiple community membership for many nodes — **Link clustering** (**Clique Percolation**)
- My network is HUGE and doesn't have super-large hubs — **Louvain** (**Infomap, link clustering**)
- My network is HUGE and has lots of super-hubs — **Louvain** (**Infomap**)
- I'd like to see the detailed hierarchical structure — **Link clustering**

# THANK YOU!

@yy

yyahn@indiana.edu