

Social Networks in Economic Geography

University of Warsaw

February 17-19 2025

Balázs Lengyel

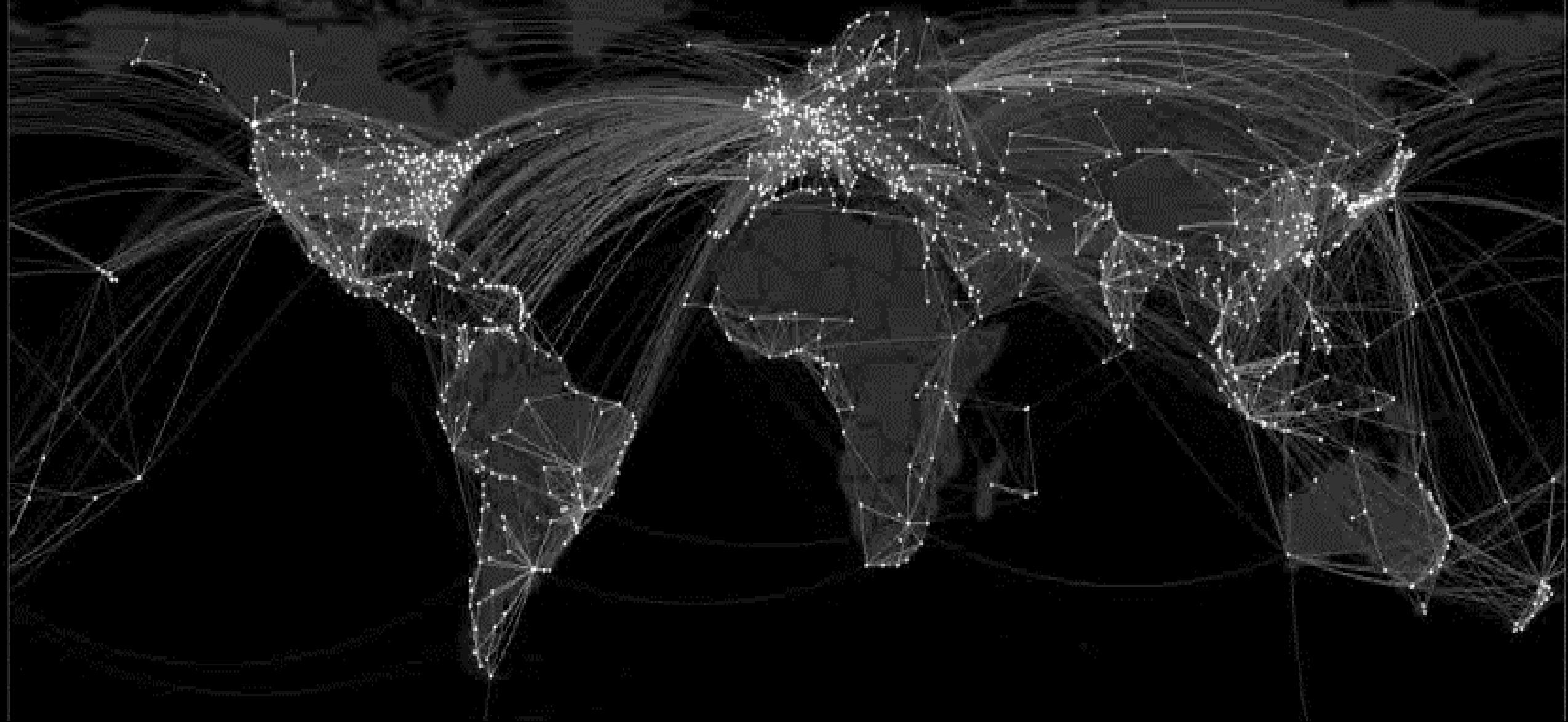
Agglomeration, Networks, and Innovation Research Lab
HUN-REN Centre for Economic- and Regional Studies
Corvinus University of Budapest



We live in small world network.



Diffusion of viruses and information is fast



The global society is segregated

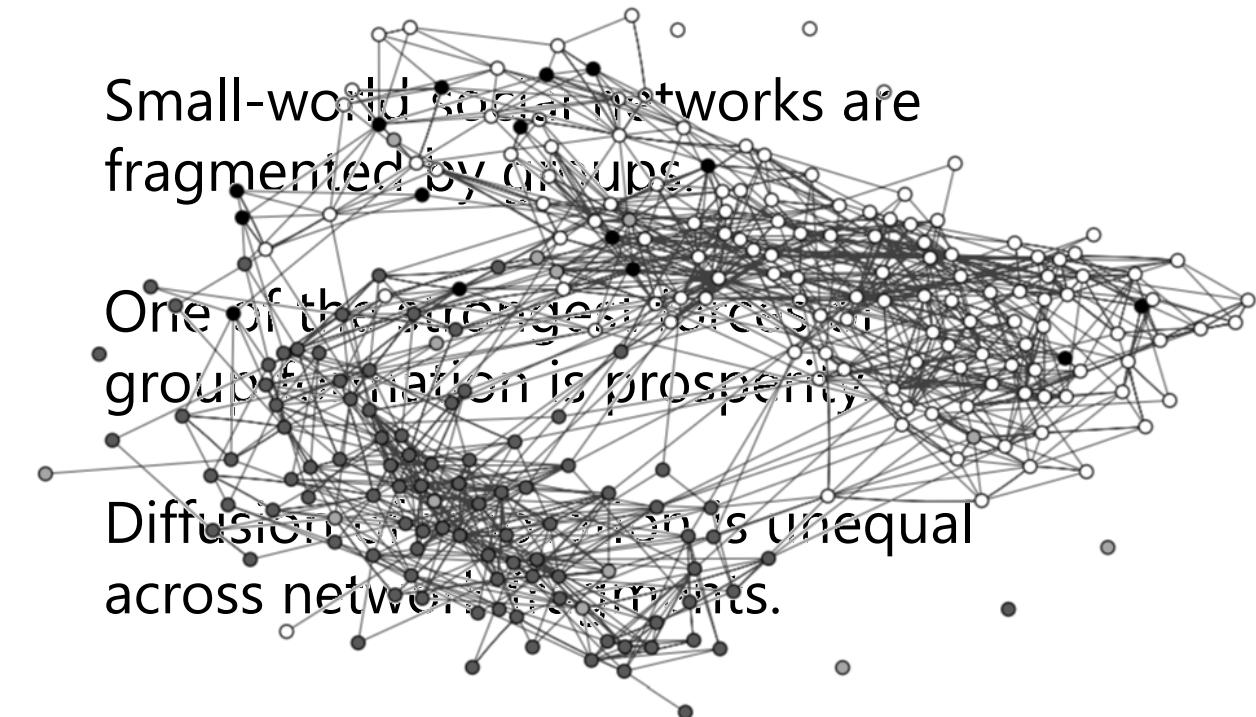
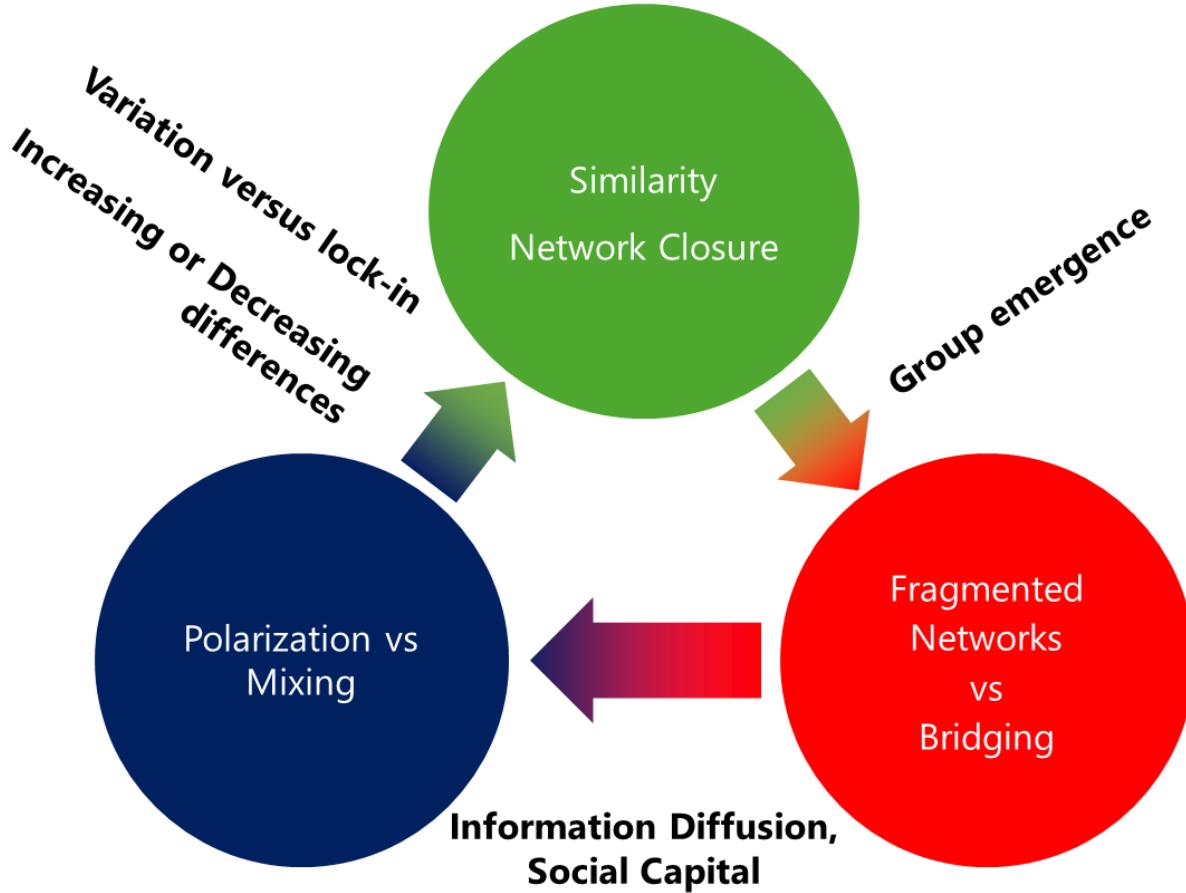
An aerial photograph of Mumbai, India, illustrating the theme of societal segregation. The image shows a vast expanse of urban density, with a clear division between two types of residential areas. In the foreground and middle ground, there is a massive concentration of small, blue-roofed houses forming a dense slum. This is separated by a few narrow roads from a more affluent area in the background, characterized by numerous tall, modern apartment buildings and office towers.

Why do we observe growing inequalities
if diffusion of ideas and knowledge
is very quick
in small-world networks?

Prosperous areas next to a slum in Mumbai.

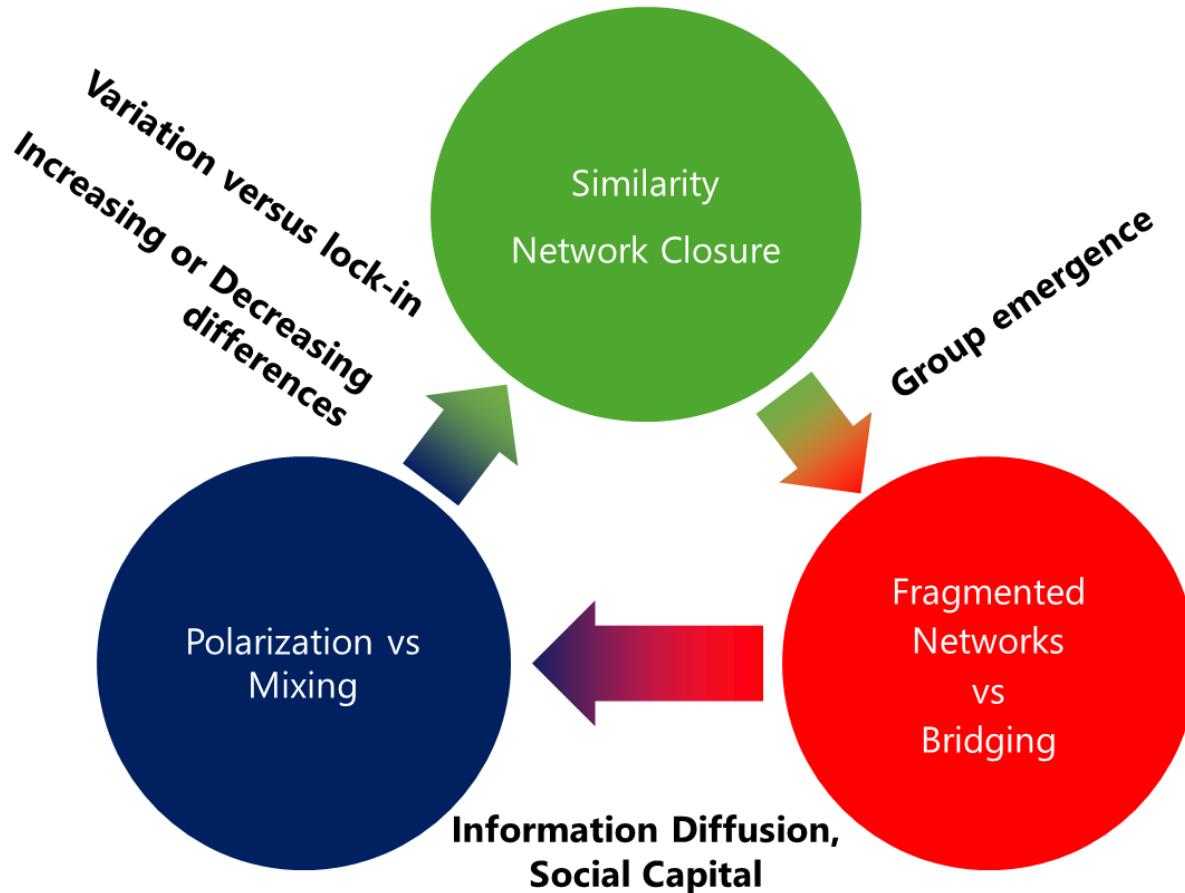
Network fragmentation and inequalities

DiMaggio-Garip (2012) Annual Review of Sociology



Curraini, Pin, Jackson (2009) Econometrica

Network fragmentation and innovation



Similar specialization in knowledge increases link probability.

- Increases understanding and specialization

Local collaboration networks can become too cohesive and locked-in into technologies – Grabher (1993), Boschma & Frenken (2010), Giuliani (2013), Balland et al. (2016)

Bridging is on purpose

- To collect new knowledge.

In this course

Day 1

1. Network science: concepts and tools to study
2. Social science: Economic relevance and the geography of social and collaboration networks
3. Network science: Community detection
4. Data Lab: network data and visualization in Gephi

Day 2

1. Social science: Social capital
2. Data Lab: Coding basics in R
3. Social science: Social and collaboration networks and regional development

Day 3

1. Network science: Diffusion in networks
2. Network and social science: Spatial diffusion
3. Data Lab: R diffusion with ABM on networks

Social Networks and Economic Geography

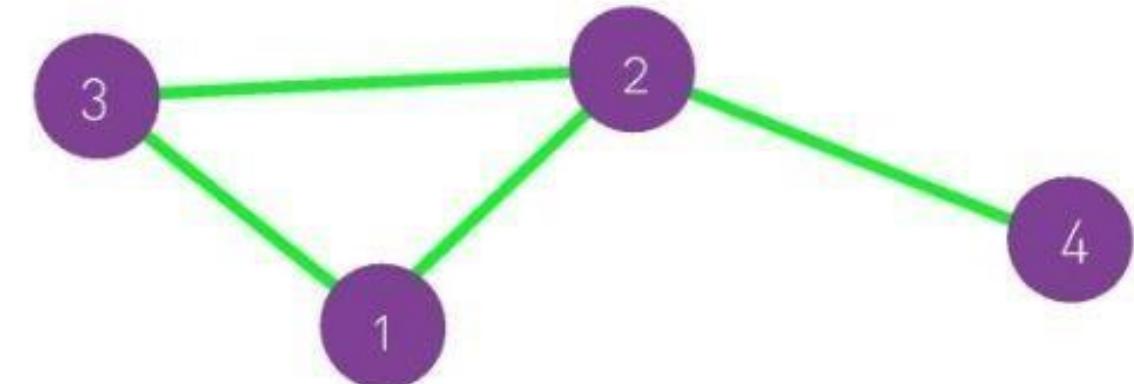
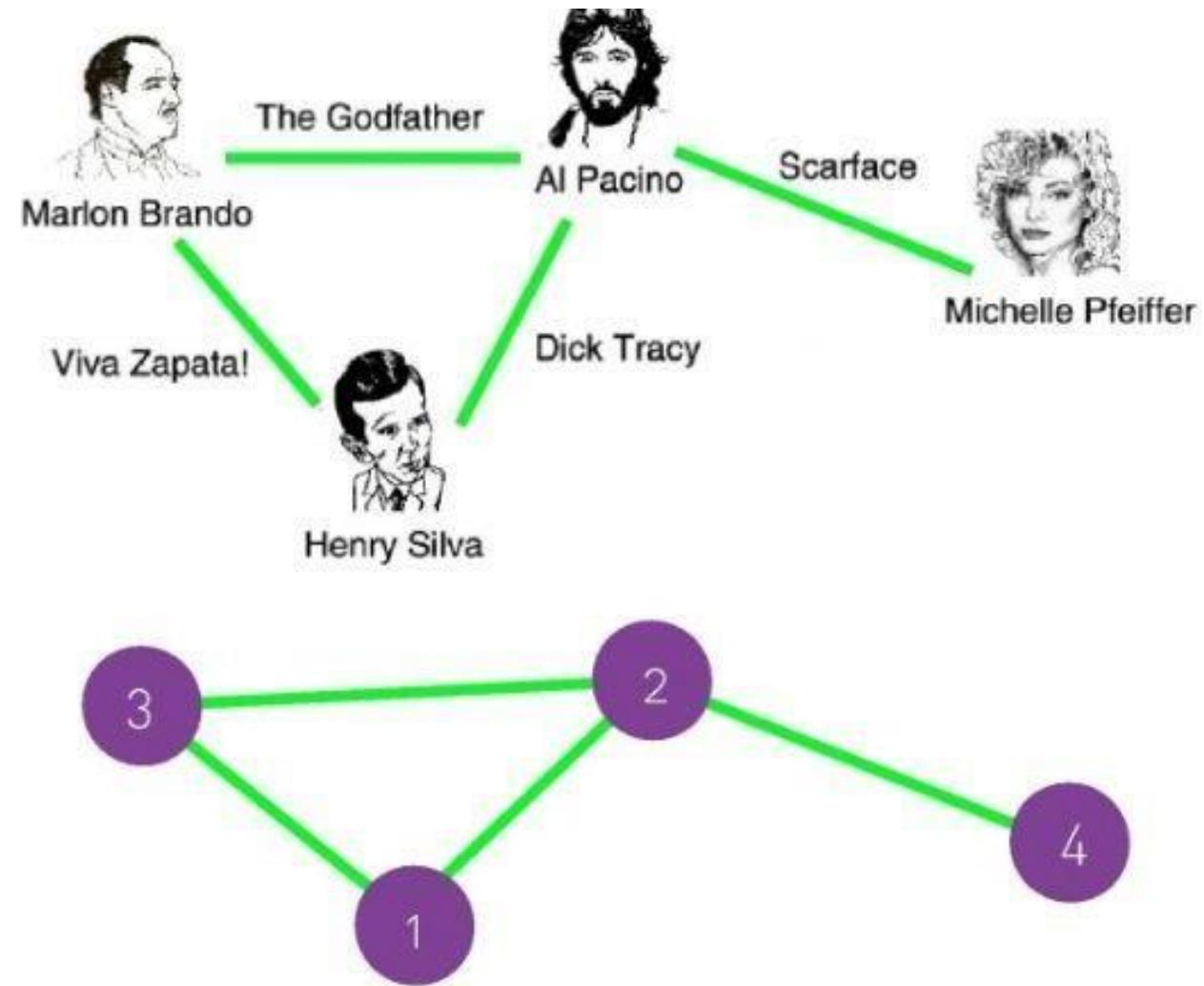
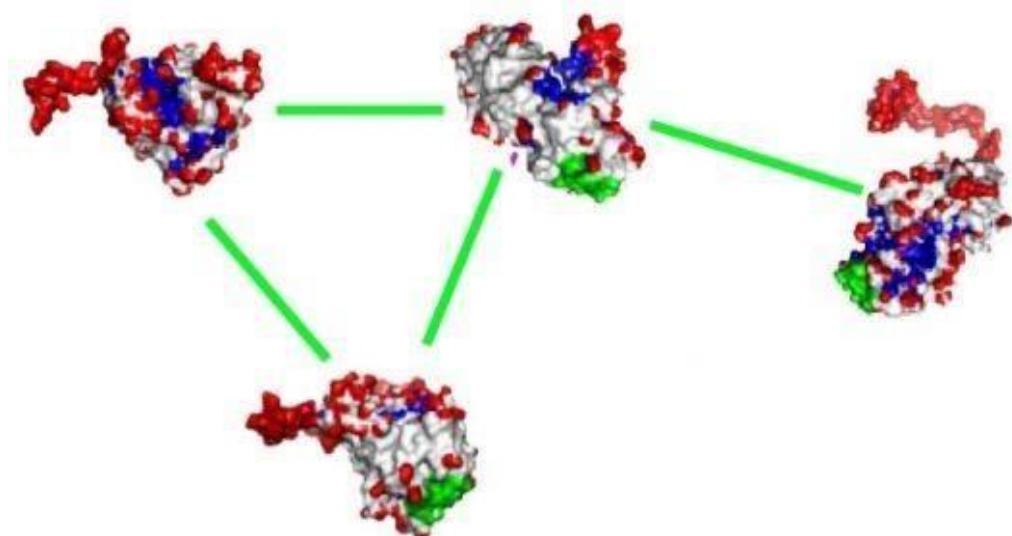
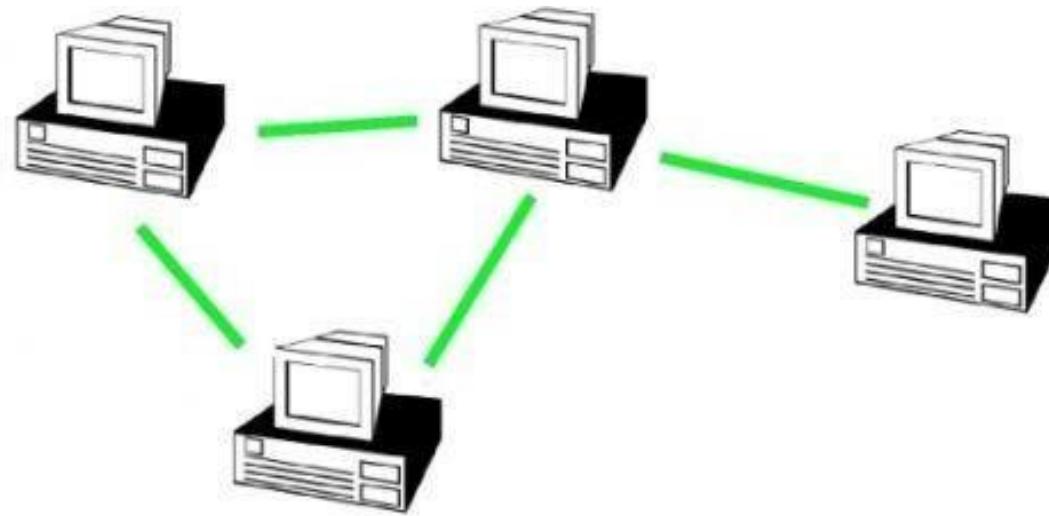
Balázs Lengyel

lengyel.balazs@krtk.hun-ren.hu

Class 1: Network Science

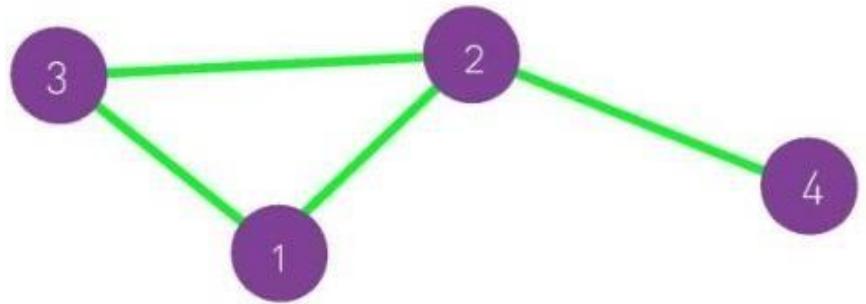
February 17 2024

Networks and Graphs



Network	Nodes	Links	Directed / Undirected	N	L	$\langle k \rangle$
Internet	Routers	Internet connections	Undirected	192,244	609,066	6.34
WWW	Webpages	Links	Directed	325,729	1,497,134	4.60
Power Grid	Power plants, transformers	Cables	Undirected	4,941	6,594	2.67
Mobile-Phone Calls	Subscribers	Calls	Directed	36,595	91,826	2.51
Email	Email addresses	Emails	Directed	57,194	103,731	1.81
Science Collaboration	Scientists	Co-authorships	Undirected	23,133	93,437	8.08
Actor Network	Actors	Co-acting	Undirected	702,388	29,397,908	83.71
Citation Network	Papers	Citations	Directed	449,673	4,689,479	10.43
E. Coli Metabolism	Metabolites	Chemical reactions	Directed	1,039	5,802	5.58
Protein Interactions	Proteins	Binding interactions	Undirected	2,018	2,930	2.90

Networks and Graphs



Average Degree

We can express the total number of links, L , as the sum of the node degrees, $k_1=2$, $k_2=3$, $k_3=2$, $k_4=1$

$$L = \frac{1}{2} \sum_{i=1}^N k_i$$

In undirected networks average degree:

$$\langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i = \frac{2L}{N}$$

In case of directed network the total number of links is a sum of outgoing and ingoing degrees:

$$k_i = k_i^{in} + k_i^{out}$$

$$\langle k^{in} \rangle = \frac{1}{N} \sum_{i=1}^N k_i^{in} = \langle k^{out} \rangle = \frac{1}{N} \sum_{i=1}^N k_i^{out} = \frac{L}{N}$$

Networks and Graphs

Degree Distribution

Degree Distribution p_k gives the probability that a given node has a degree of k , normalized to 1:

$$\sum_{k=1}^{\infty} p_k = 1$$

,where there are N number of nodes, the normalized distribution can be expressed by:

$$p_k = \frac{N_k}{N}$$

,thus the number of degree- k nodes given by the degree distribution $N_k = Np_k$.

Networks and Graphs

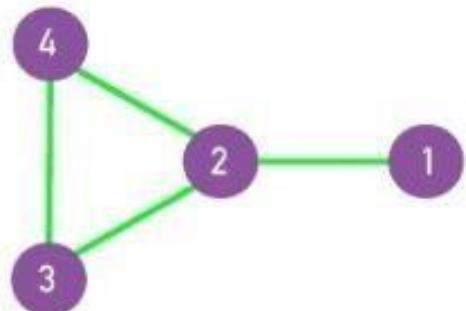
Degree Distribution

Degree Distribution plays a key role in Graph and Network Theory. Most of the structural network properties require the calculation of p_k .

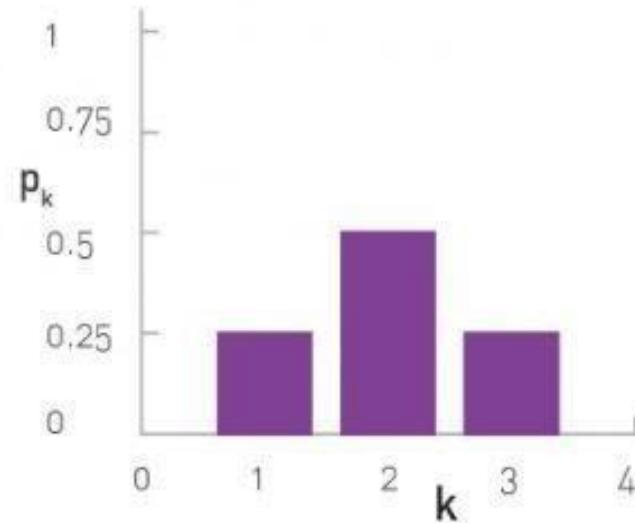
e.g Average Degree can be expressed as a function of Degree Distribution:

$$\langle k \rangle = \sum_{k=0}^{\infty} kp_k$$

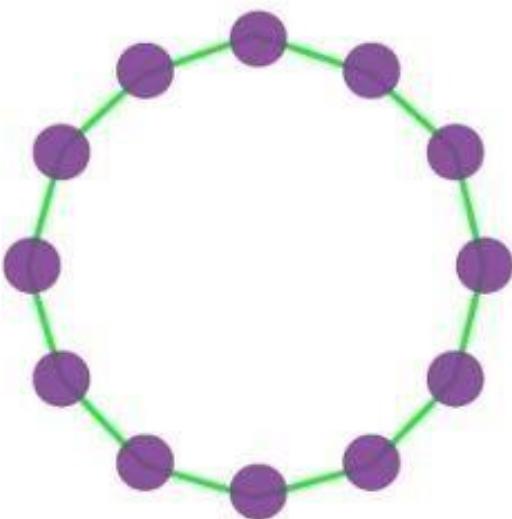
a.



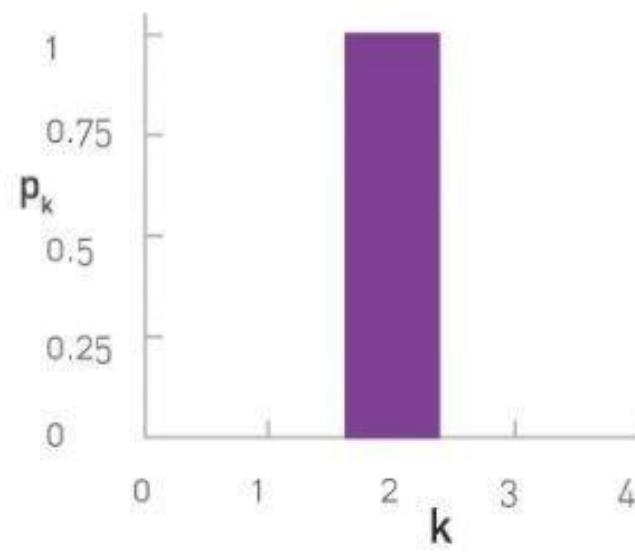
b.

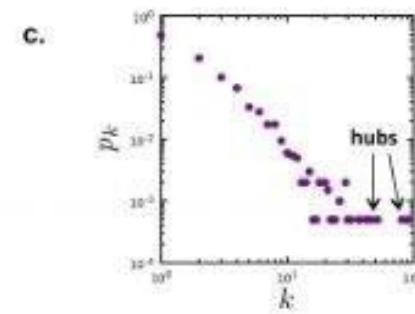
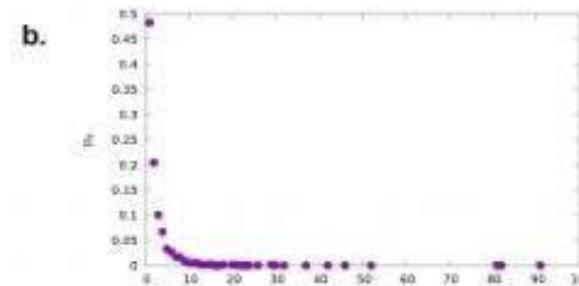
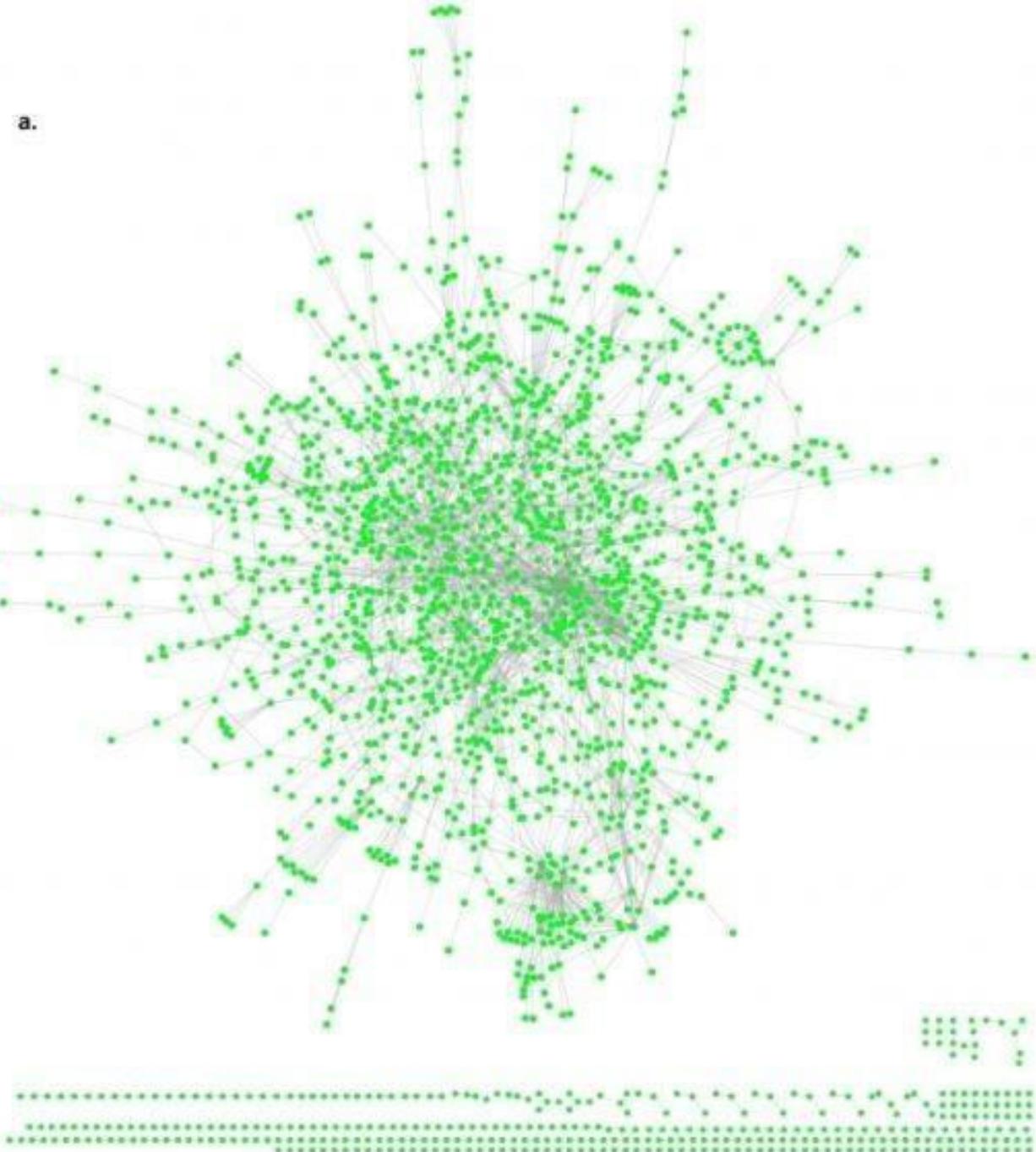


c.

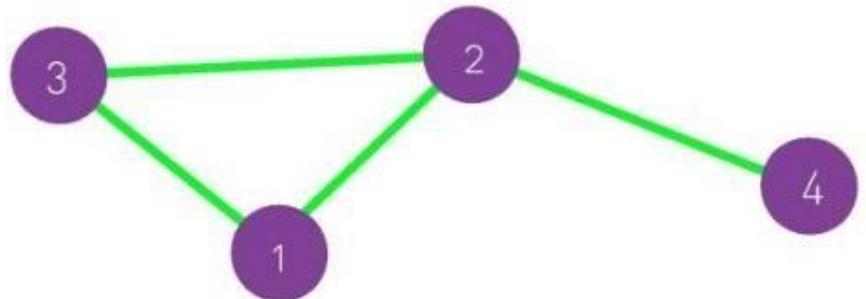


d.





Adjacency Matrix



To keep track of the links we can provide a complete edge list:

$$\{(1,2);(1,3);(2,3);(2,4)\}$$

We can translate it into an Adjacency Matrix, where network of N nodes has N rows and N columns:

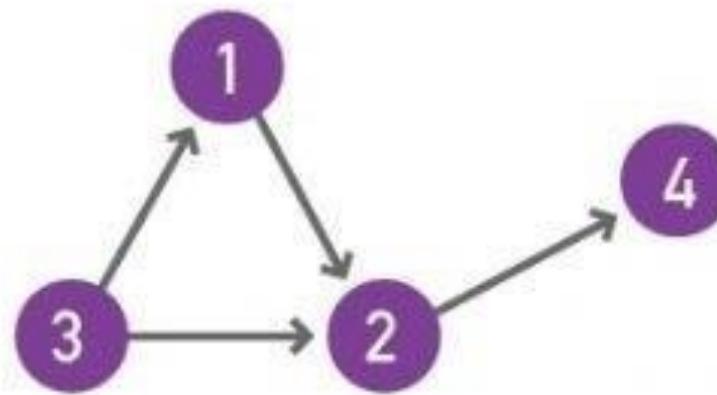
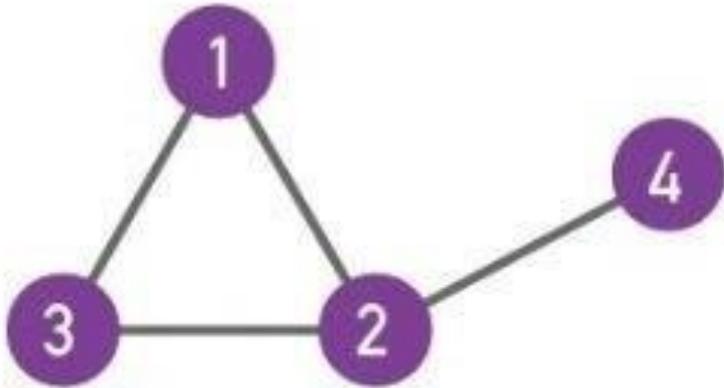
- $A_{ij} = 1$ if there is a connection between i and j ,
- $A_{ij} = 0$ no connection between i and j ,
- in an undirected network $A_{ij} = A_{ji}$.

The degree k_i can be expressed by the sum of either the rows or the columns of the matrix:

$$k_i = \sum_{i=1}^N A_{ij} = \sum_{j=1}^N A_{ij}$$

$$2L = \sum_{i=1}^N k_i^{in} = \sum_{i=1}^N k_i^{out} = \sum_{ij} A_{ij}$$

Adjacency Matrix



$$A_{ij} = \begin{matrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{matrix}$$

$$A_{ij} = \begin{matrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{matrix}$$

Adjacency Matrix

$$A_{ij} = \begin{matrix} & \begin{matrix} 0 & 1 & 1 & 0 \end{matrix} \\ \begin{matrix} 0 & 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{matrix} & \begin{matrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{matrix} \end{matrix}$$

$$k_2 = \sum_{j=1}^4 A_{2j} = \sum_{i=1}^4 A_{i2} = 3$$

$$k_2^{\text{in}} = \sum_{j=1}^4 A_{2j} = 2, k_2^{\text{out}} = \sum_{i=1}^4 A_{i2} = 1$$

$$A_{ij} = A_{ji} \quad A_{ii} = 0 \quad A_{ij} \neq A_{ji} \quad A_{ii} = 0$$

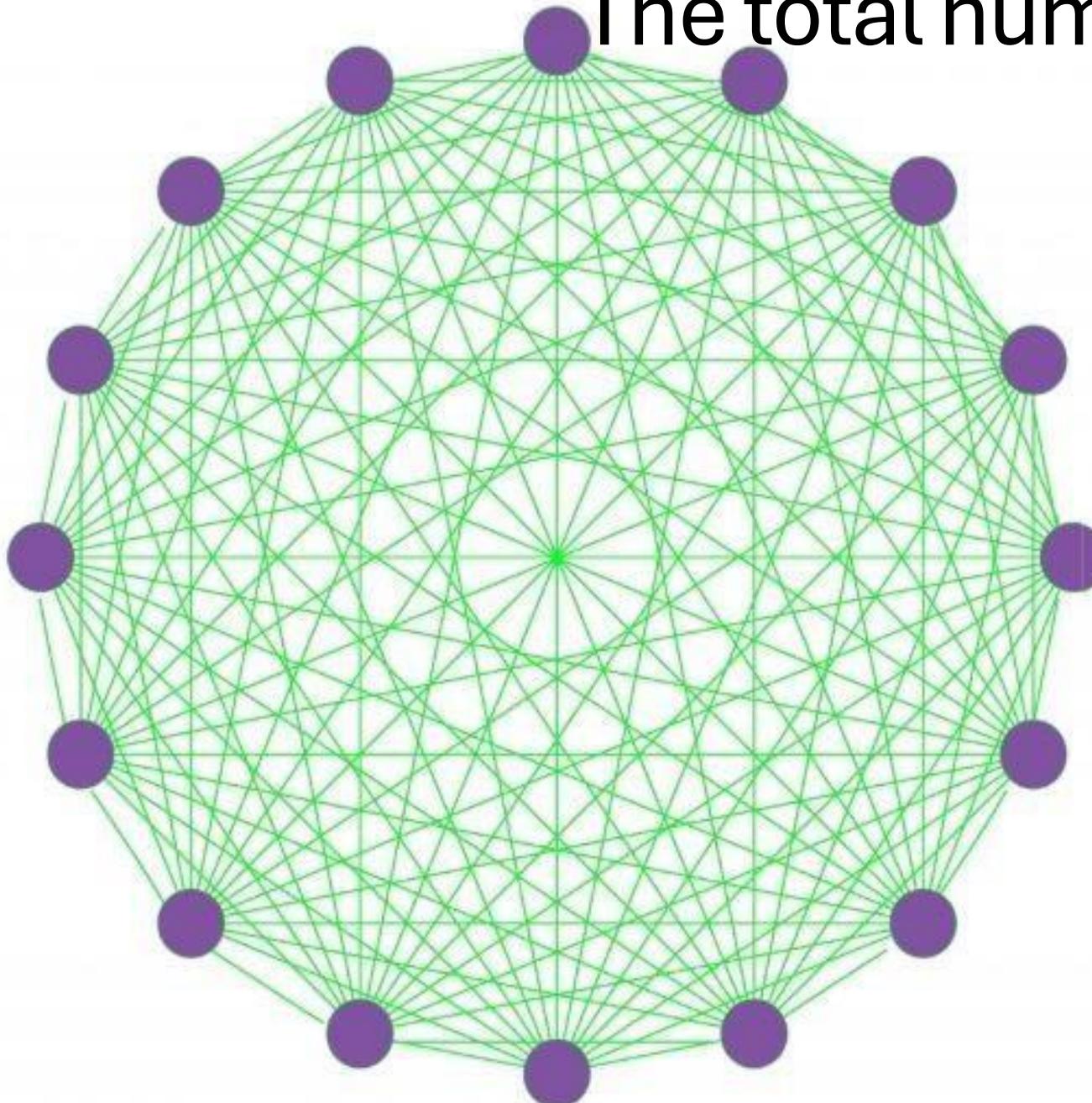
$$L = \frac{1}{2} \sum_{i,j=1}^N A_{ij}$$

$$L = \sum_{i,j=1}^N A_{ij}$$

$$\langle k \rangle = \frac{2L}{N}$$

$$\langle k^{\text{in}} \rangle = \langle k^{\text{out}} \rangle = \frac{L}{N}$$

The total number of links varies widely.

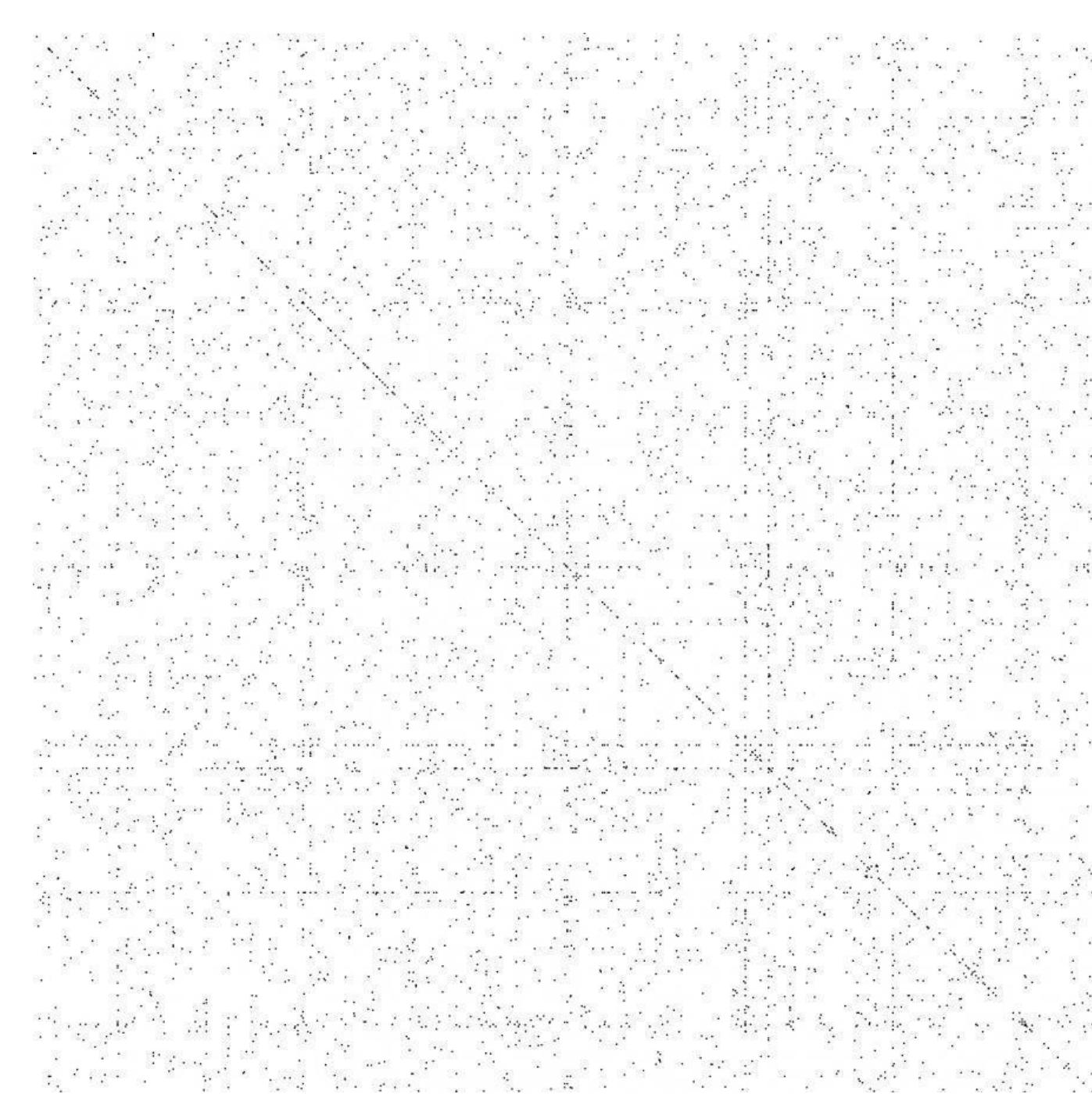


$$L_{min} = 0$$

$$L_{max} = \frac{N(N-1)}{2}$$

Real networks are sparse!

$L \ll L_{max}$.



$$L_{min} = 0$$

$$L_{max} = \frac{N(N-1)}{2}$$

Real networks are sparse!

$$L \ll L_{\max}$$

Overwhelming fraction of elements
are zero!

Weighted Networks

Unweighted Networks

$$A_{ij} = 1$$

$A_{ij} = 1$ if there connection between i and j

- is there a phone call?
- any export/import between i and j ?

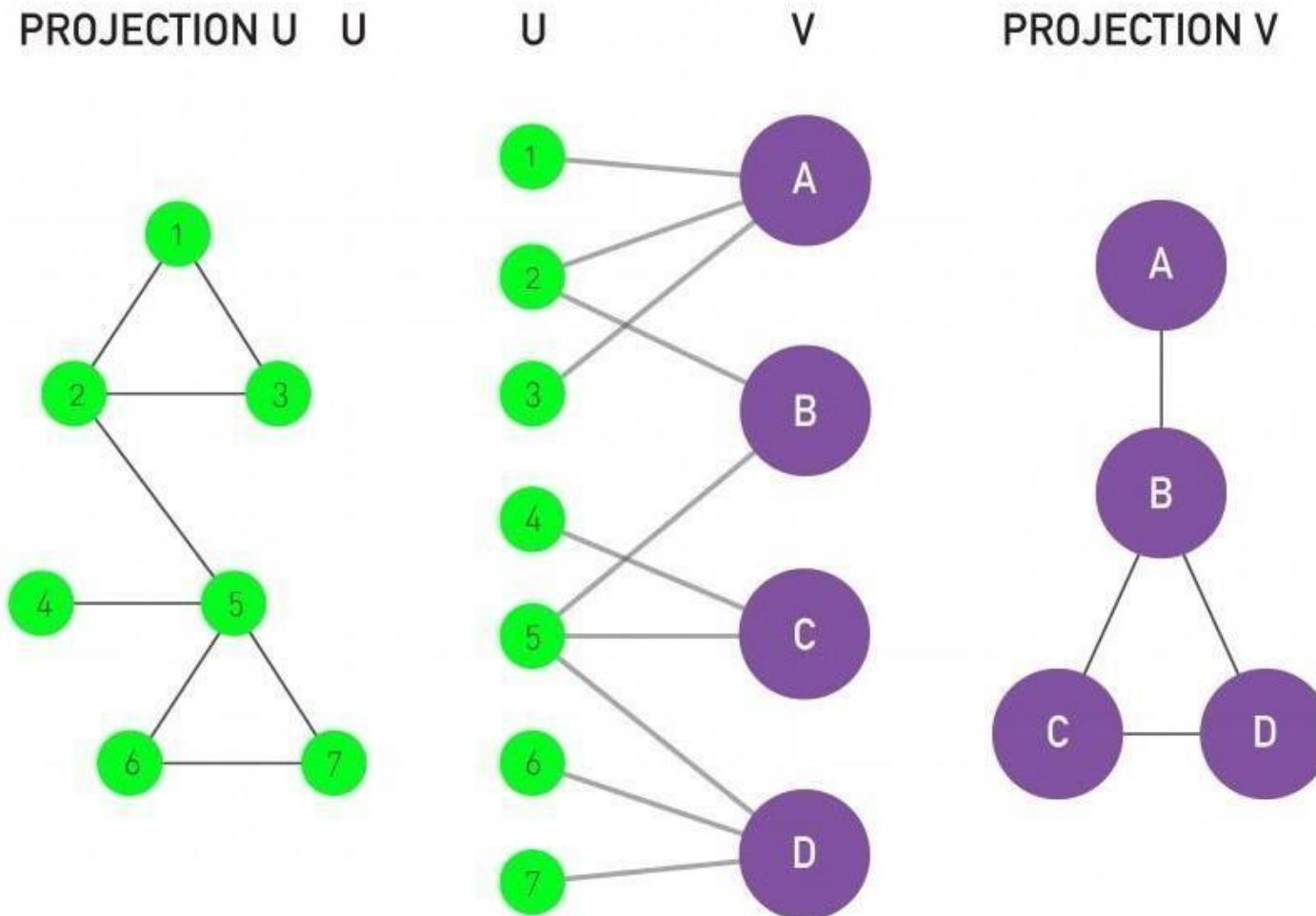
Weighted Networks

$$A_{ij} = w_{ij}$$

$A_{ij} = w_{ij}$ the extent of connection between i and j

- how many phone calls have happened?
- the extent of export/import between countries?

Bipartite Networks



Nodes can be divided into two disjoint sets U and V

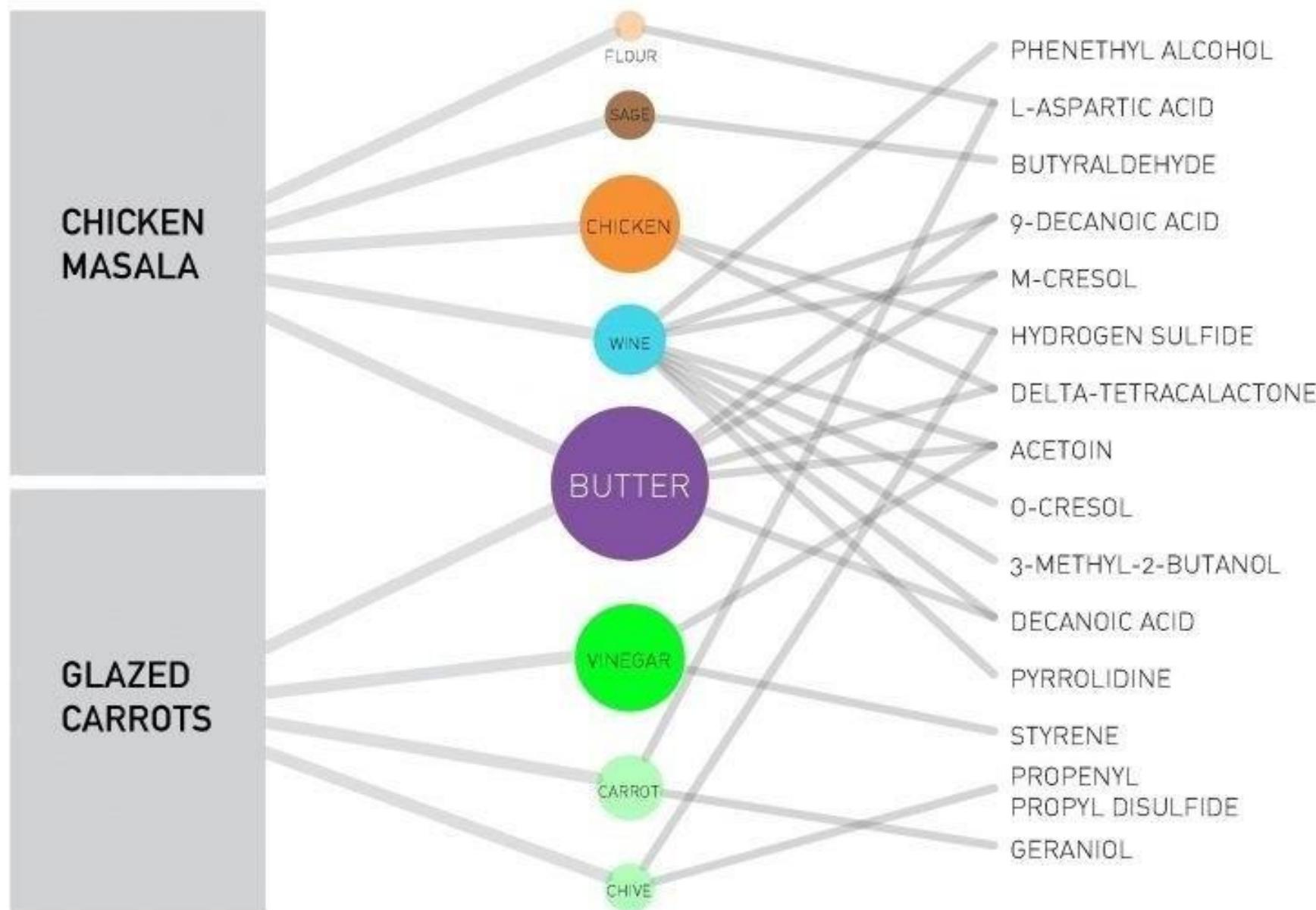
movie network:
 U (actors)
 V (movies)

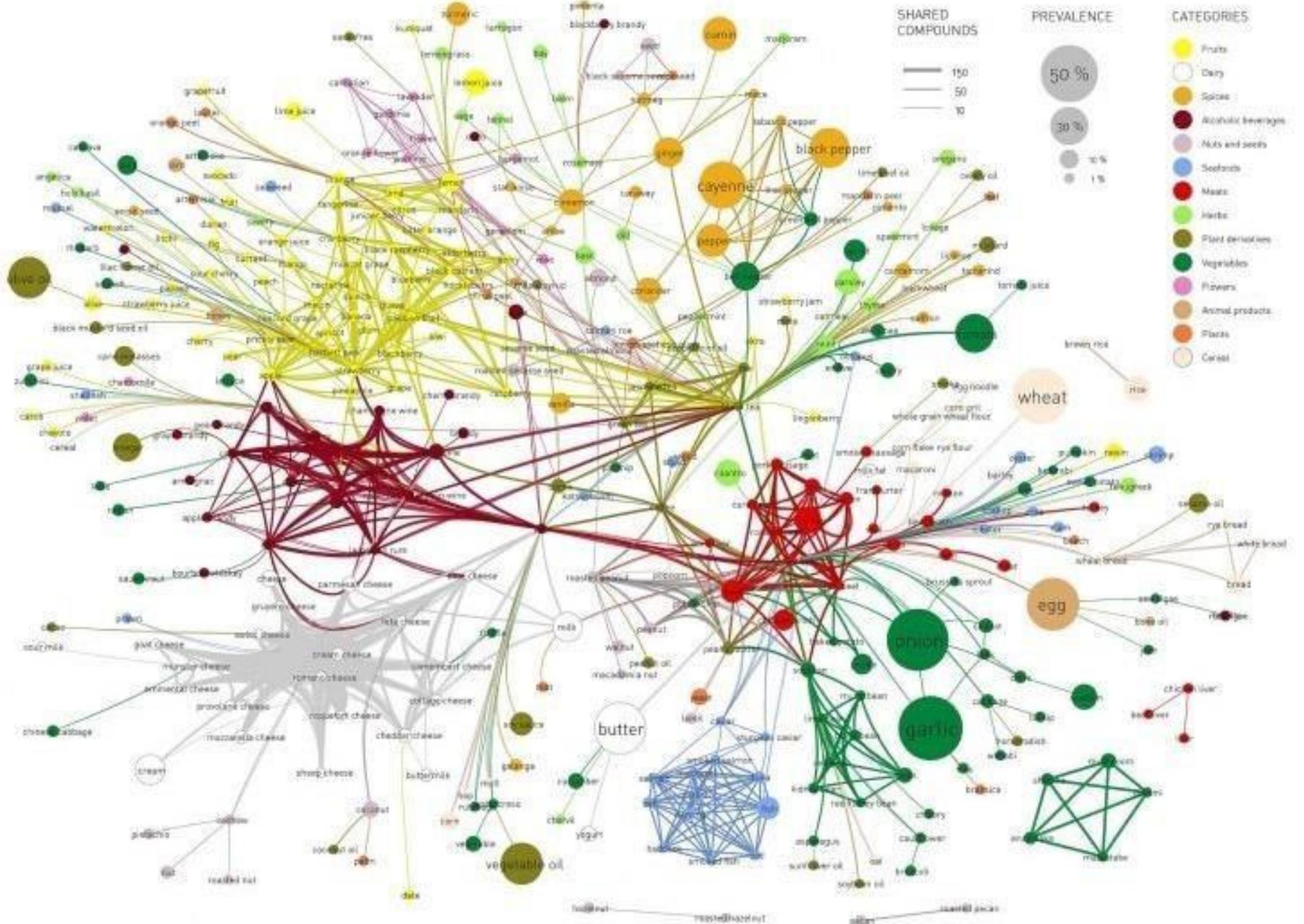
co-publication network:
 U (inventors/researchers)
 V (patents/publications)

RECIPES

INGREDIENTS

COMPOUNDS



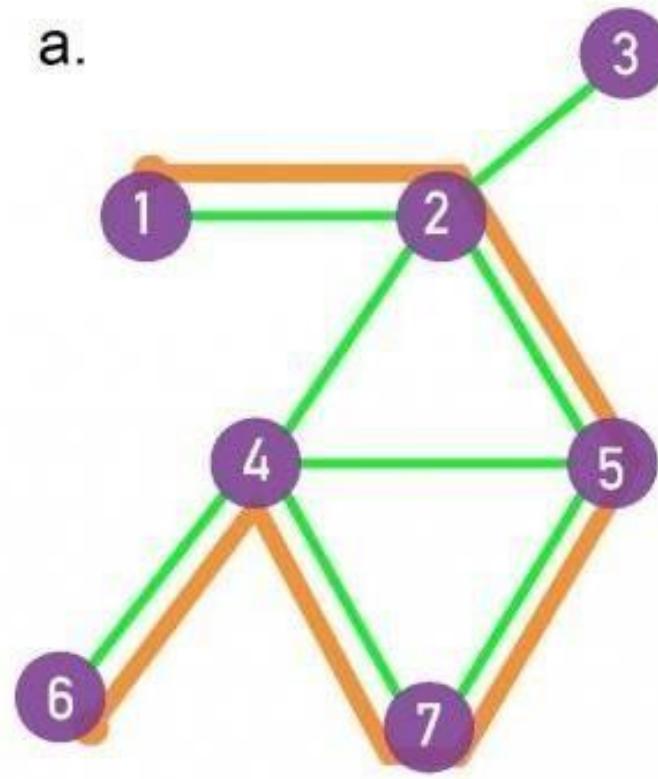


Barabasi (2016) Image 2.11

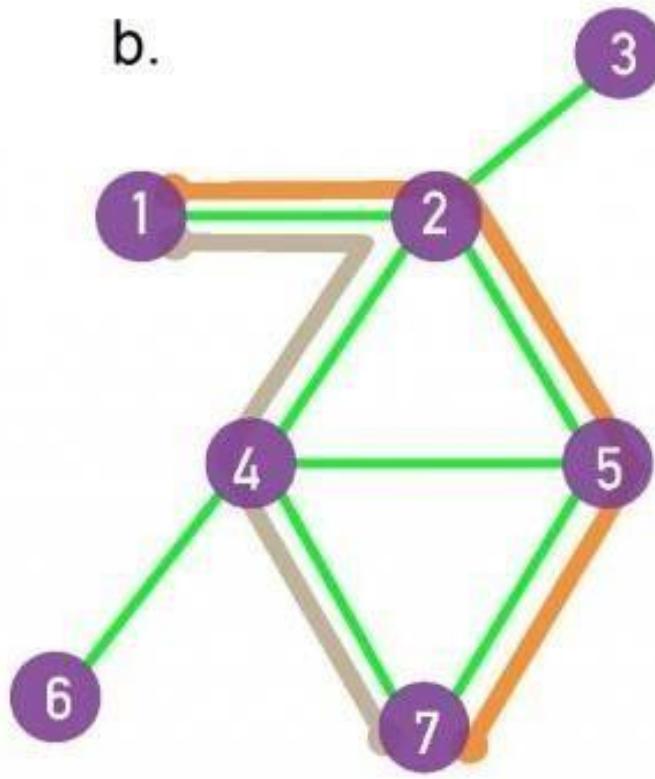
Paths and Distances

To determine the probability of interaction between two components of a system we usually use the physical distance between the agents, e.g. distance between two atoms or distance between two planets.

a.

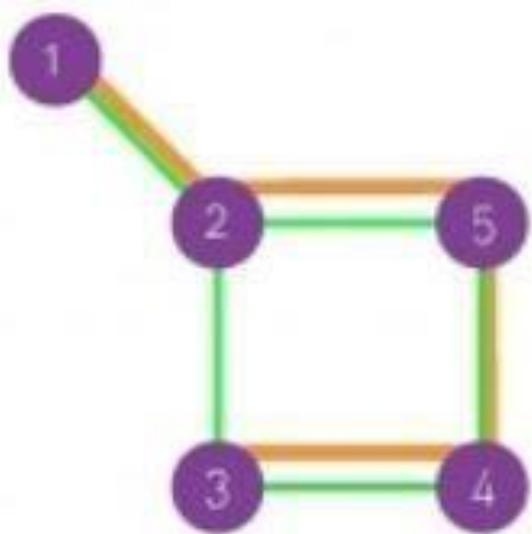


b.



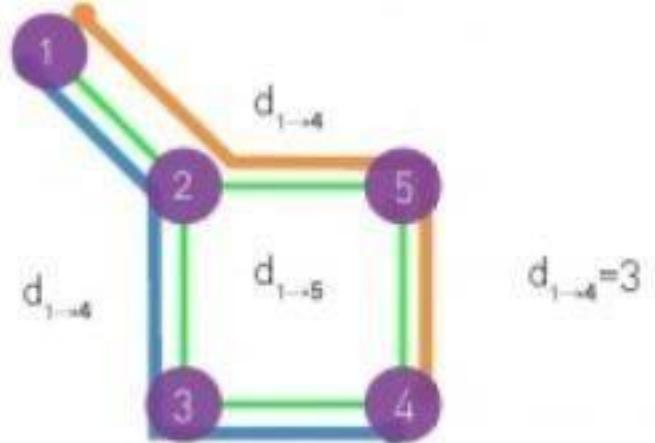
Paths and Distances

- **Path** is a sequence of connected nodes



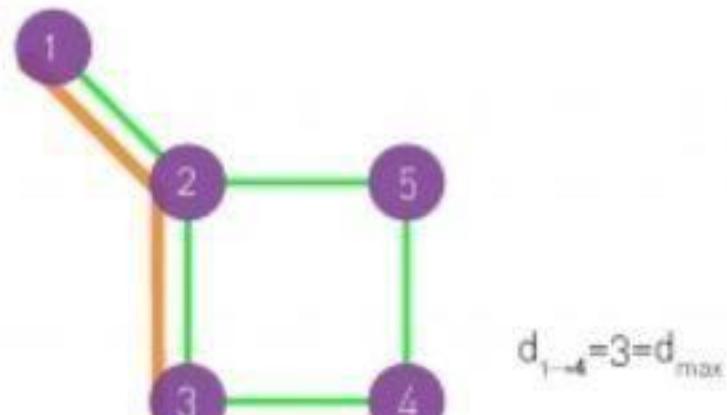
Paths and Distances

- **Path** is a sequence of connected nodes
- **Shortest path (d)** is the shortest distance between two nodes.



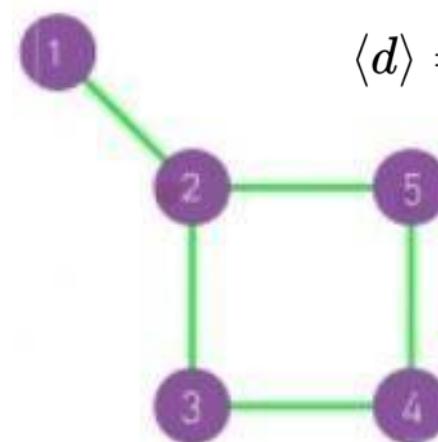
Paths and Distances

- **Path** is a sequence of connected nodes
- **Shortest path (d)** is the shortest distance between two nodes.
- **Diameter (d_{\max})** is the shortest path between the two furthest points of the graph.
 $d_{\max}(1,4)=3$



Paths and Distances

- **Path** is a sequence of connected nodes
- **Shortest path (d)** is the shortest distance between two nodes.
- **Diameter (d_{\max})** is the shortest path between the two furthest points of the graph.
 $d_{\max}(1,4)=3$
- **Average Path Length ($\langle d \rangle$)** is the average of the shortest path between all pair of nodes.
 $\langle d \rangle = 1.6$



$$\langle d \rangle = \frac{1}{N(N-1)} \sum_{i,j=1; i \neq j} d_{ij}$$

$$\begin{aligned}\langle d \rangle &= [d_{1 \rightarrow 2} + d_{1 \rightarrow 3} + d_{1 \rightarrow 4} + d_{1 \rightarrow 5} + \\ &+ d_{2 \rightarrow 3} + d_{2 \rightarrow 4} + d_{2 \rightarrow 5} + \\ &+ d_{3 \rightarrow 4} + d_{3 \rightarrow 5} + \\ &+ d_{4 \rightarrow 5}] / 10 = 1.6\end{aligned}$$

Clustering

Clustering coefficient shows the degree to which the neighbors of a given node connected to each other.

$$C_i = \frac{2L}{k_i(k_i-1)}$$

, where L is the number of links between k_i neighbors of node i .

$C_i = 0$ if there is no connection between the neighbors of i .

$C_i = 1$ if the neighbors of i form a complete graph.

$C_i = 0.5$ implies that there is a 50% chance that two neighbors of i are connected to each other.

Average Clustering Coefficient captures the degree of clustering of the whole network, or in other words it shows the probability that two neighbors of a randomly selected node connected to each other:

$$\langle C \rangle = \frac{1}{N} \sum_{i=1}^N C_i$$

Clustering

Clustering Coefficient

$$C_i = \frac{2L}{k_i(k_i-1)}$$

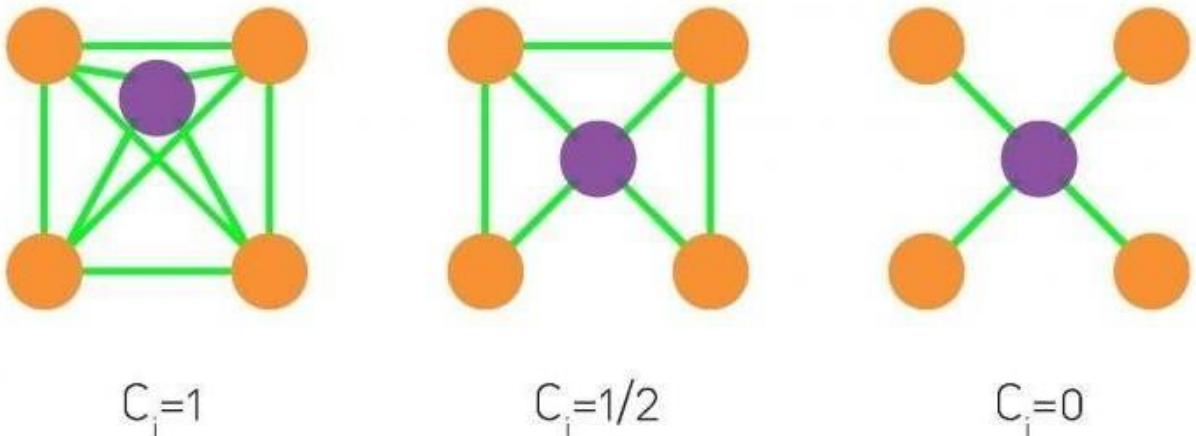
Average Clustering Coefficient

$$\langle C \rangle = \frac{1}{N} \sum_{i=1}^N C_i$$

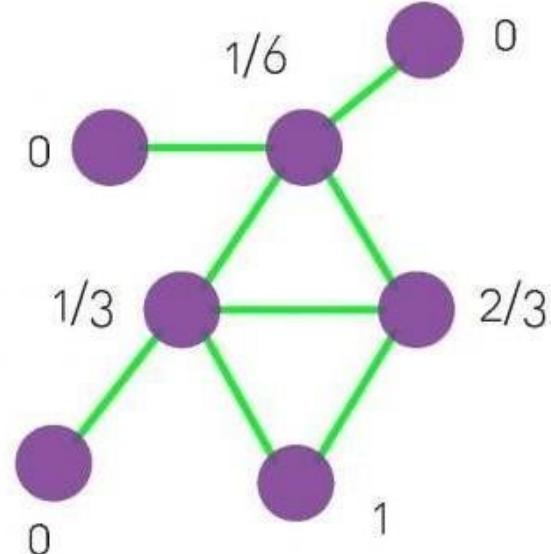
Global Clustering Coefficient

$$C_\Delta = \frac{3 \times \text{Number of } \Delta}{\text{Number of Connected Triples}}$$

a.



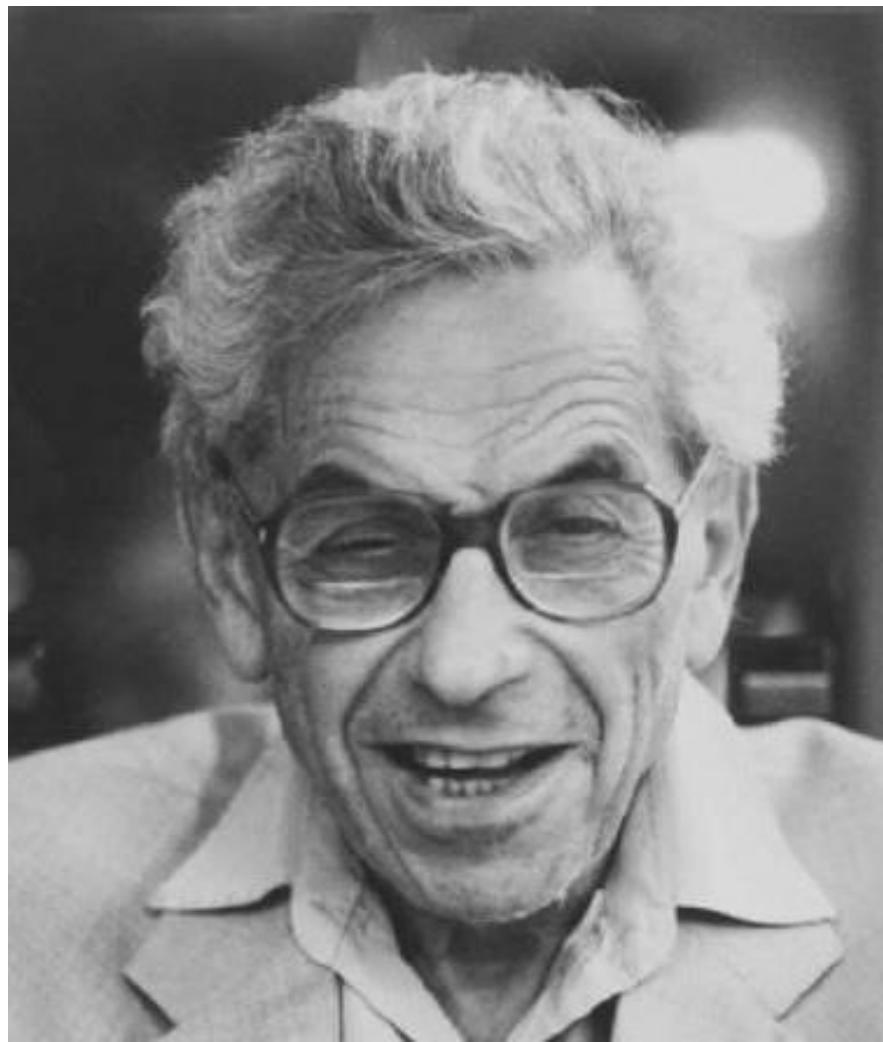
b.



$$\langle C \rangle = \frac{13}{42} \approx 0.310$$

$$C_\Delta = \frac{3}{8} = 0.375$$

Pál Erdős



“A mathematician is a device for turning coffee into theorems”

Alfréd Rényi



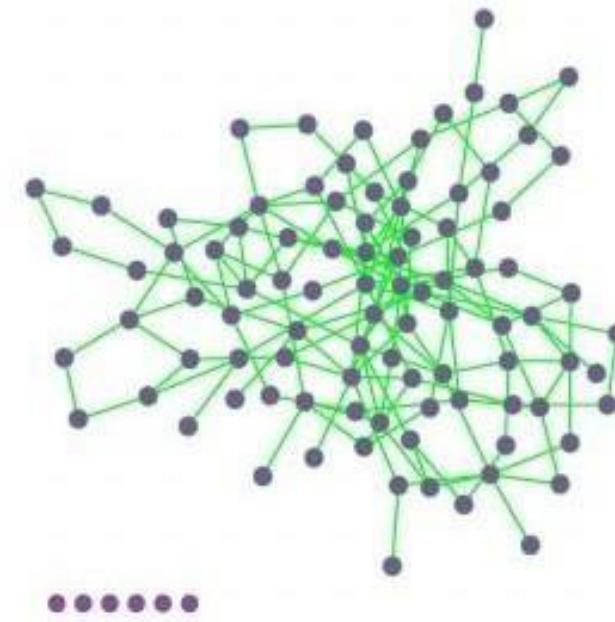
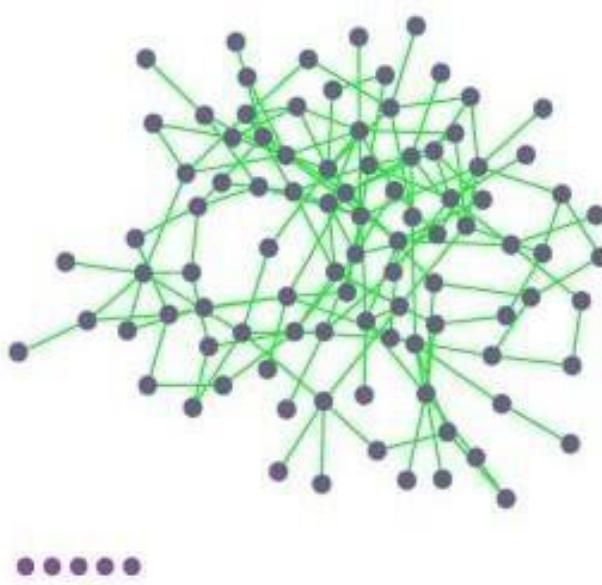
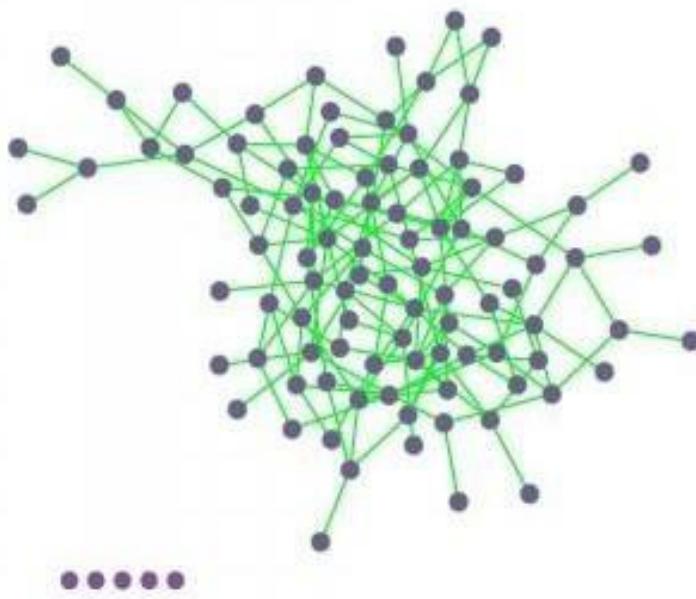
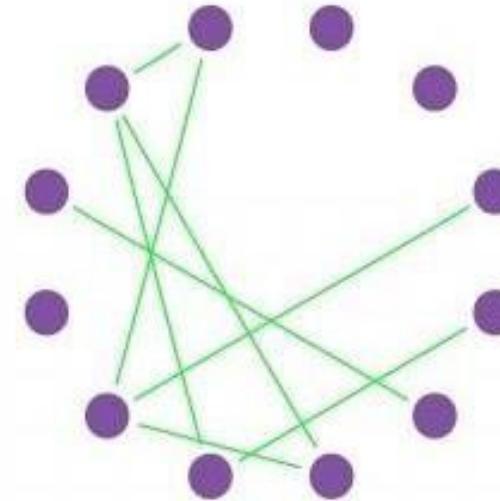
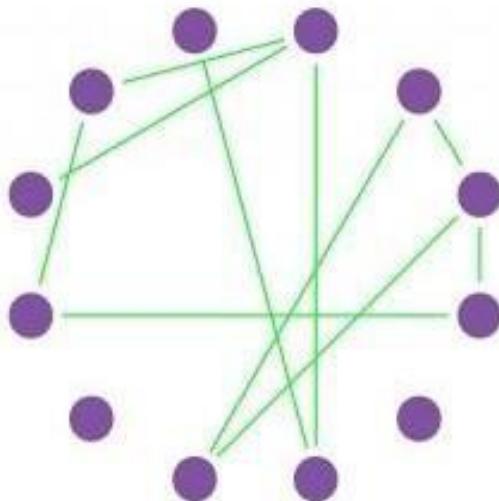
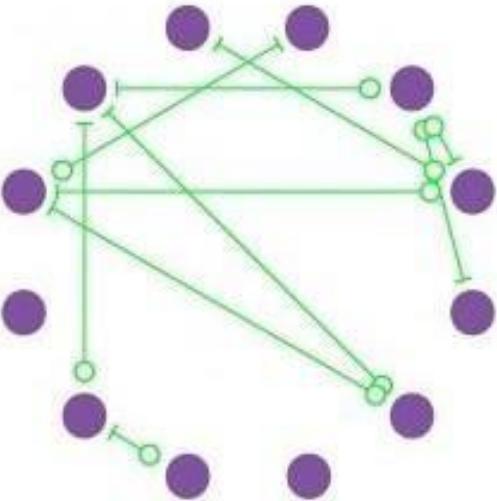
Random Network

$G(N,L)$: have N number of nodes connected by L number of random connections.
Erdős-Rényi network

$G(N,p)$: Between each pair of N_i and N_j there is a probability (p) of connection.
Gilbert

$G(N,p)$ model fixes the probability of connections, while the $G(N,L)$ fixes the total number of connections in the network.

In $G(N,L)$ model the average degree easily calculable: $\langle k \rangle = \frac{2L}{N}$, but the number of connections are rarely fixed, so in other cases we use the $G(N,p)$ model.



Evolution of a Random Network

We have two extreme cases:

- $p = 0$, all the nodes are isolated, then the largest component $N_G = 1$ and N_G/N tends to 0.
- $p = 1$, therefore $\langle k \rangle = N - 1$, then the network is a complete graph and all nodes belong to the giant component $N_G = N$, thus $N_G/N = 1$.

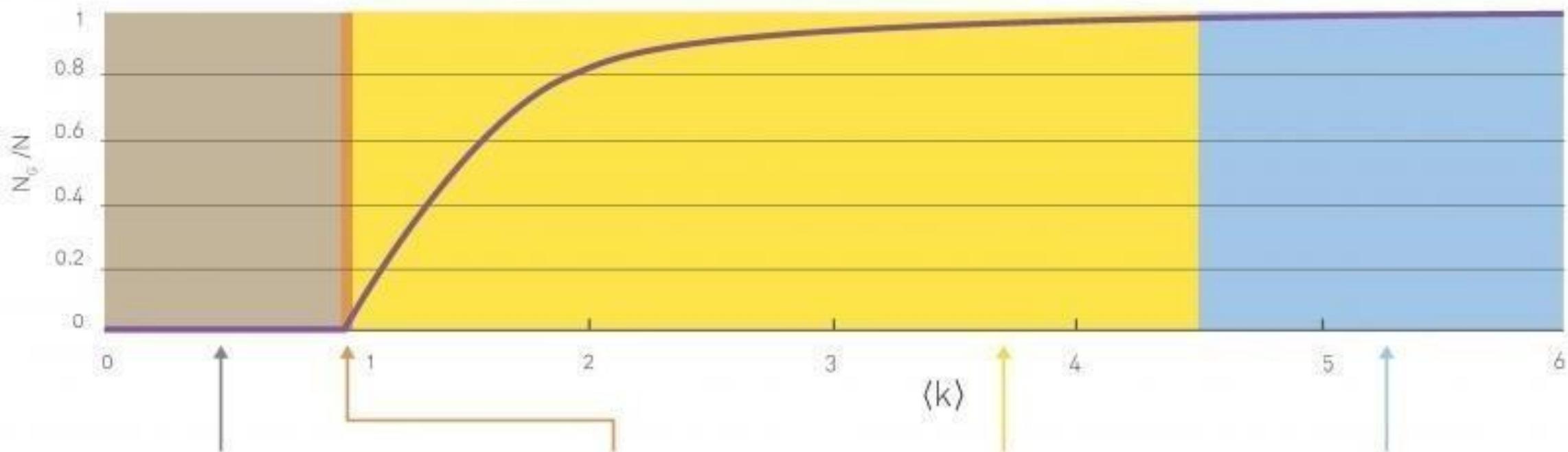
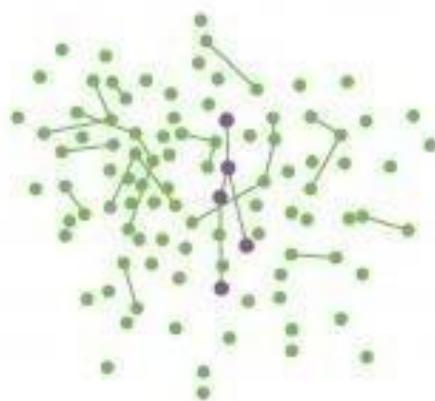
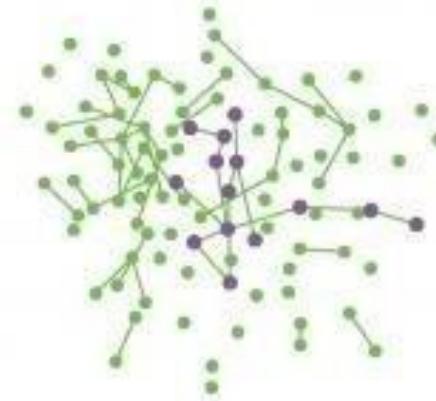
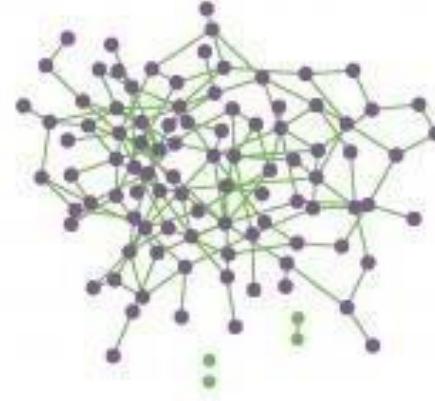
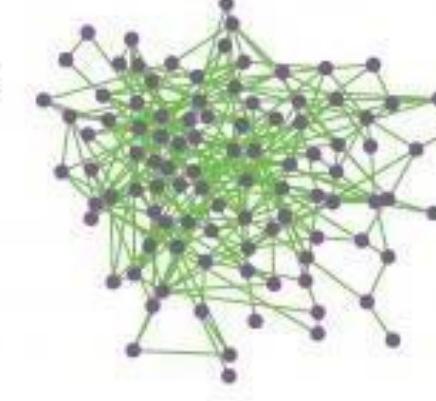
Once $\langle k \rangle$ exceeds a critical value N_G/N increases, which means that a large cluster - what we call as a Giant Component - emerges.

Critical point: $\langle k \rangle = 1$

Subcritical regime : $0 < \langle k \rangle < 1$

Supercritical: $\langle k \rangle > 1$

Connected Regime: $\langle k \rangle > \ln N$

a.**b.****c.****d.****e.**

Real Networks are Supercritical

Critical point: $\langle k \rangle = 1$

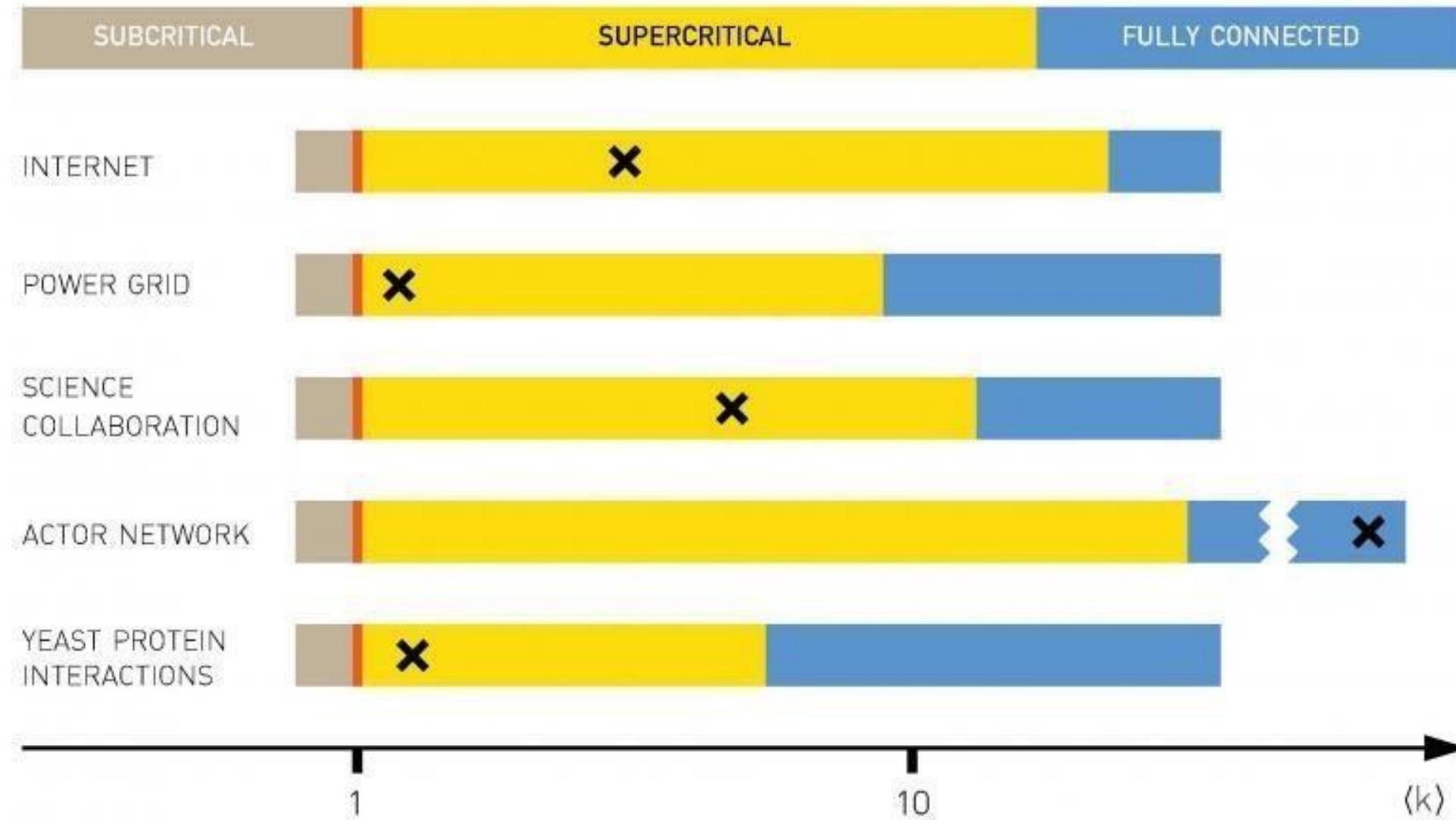
Subcritical regime : $0 < \langle k \rangle < 1$

Supercritical: $\langle k \rangle > 1$

Connected Regime: $\langle k \rangle > \ln N$

Network	N	L	$\langle k \rangle$	$\ln N$
Internet	192,244	609,066	6.34	12.17
Power Grid	4,941	6,594	2.67	8.51
Science Collaboration	23,133	94,437	8.08	10.05
Actor Network	702,388	29,397,908	83.71	13.46
Protein Interactions	2,018	2,930	2.90	7.61

Real Networks are Supercritical



Small World Networks

Six degrees of separation.

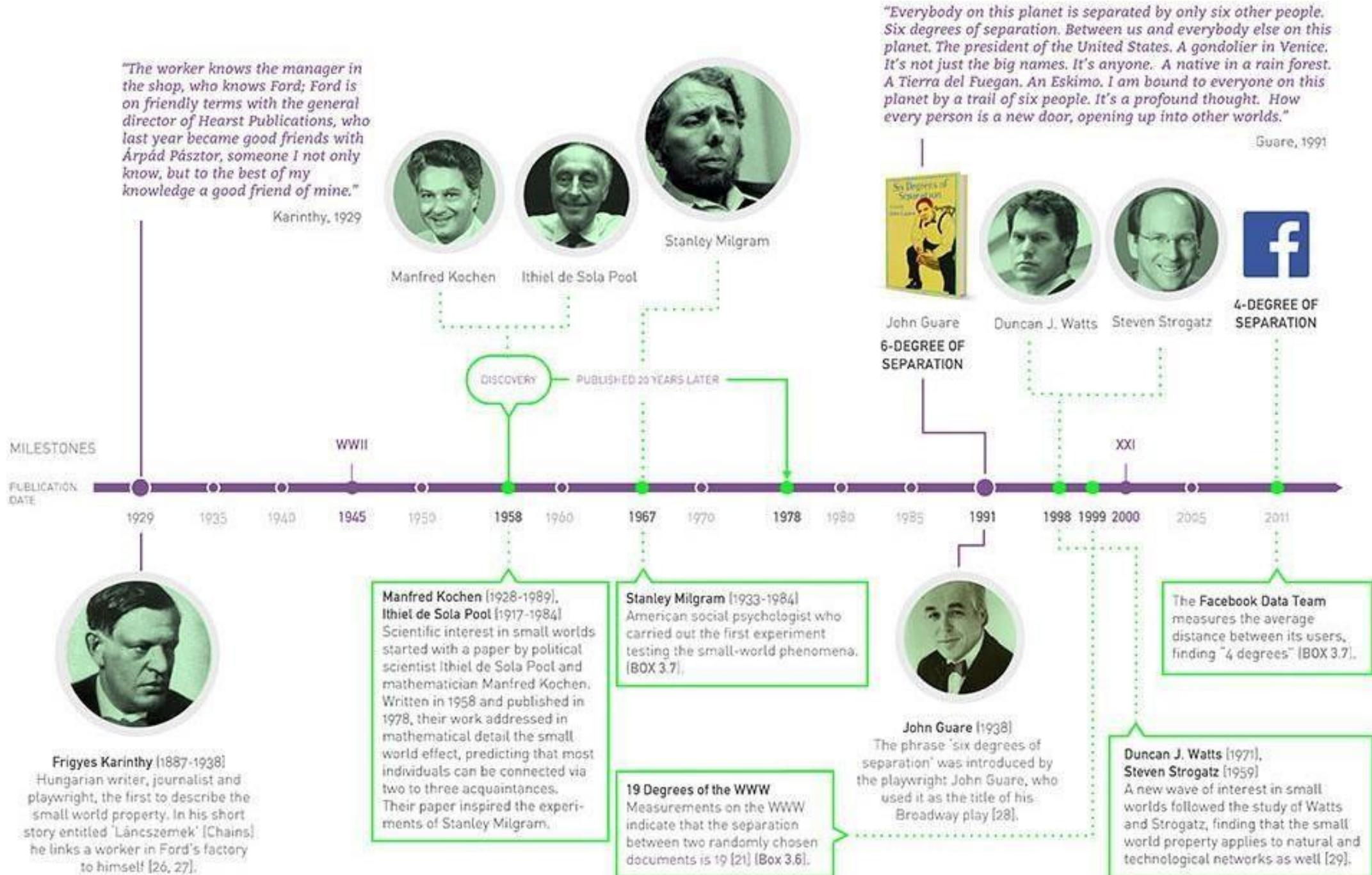
The distance between two randomly chosen node is short.

- What does short mean?
- Short compare to what?

Consider a random network with average degree of $\langle k \rangle$, then a node has on average:

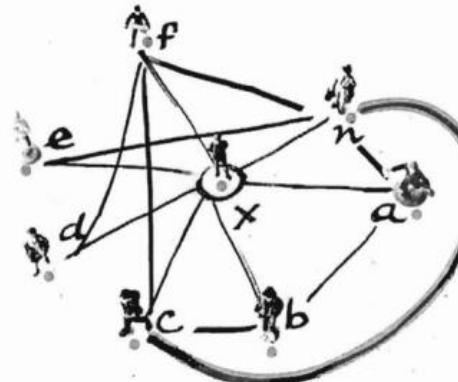
- $\langle k \rangle$ nodes at distance $d=1$,
- $\langle k \rangle^2$ nodes at distance $d=2$
- ...
- $\langle k \rangle^d$ nodes at distance d .
-

So if we assume that a social network has a $\langle k \rangle \sim 1000$, then we know 10^6 people at $d=2$, and the whole earth population is at $d=3$.

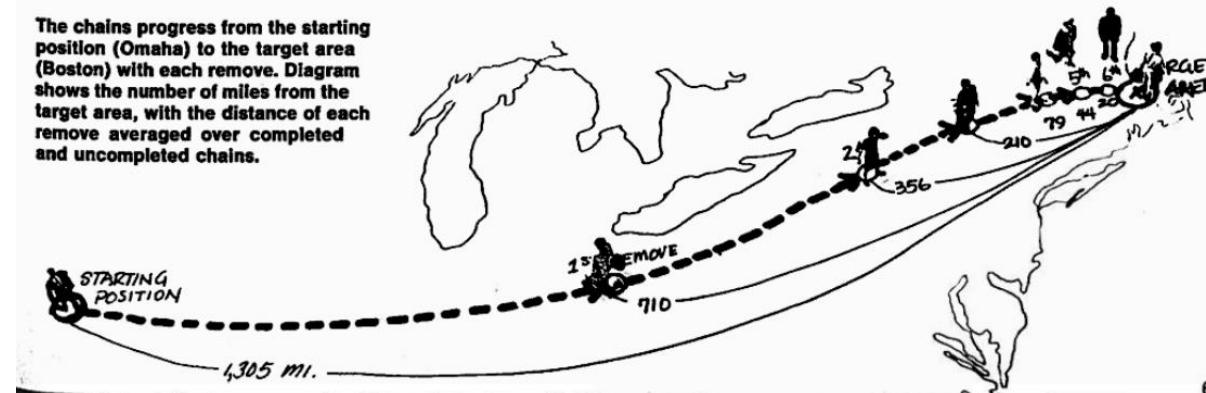




Small-world networks



The chains progress from the starting position (Omaha) to the target area (Boston) with each remove. Diagram shows the number of miles from the target area, with the distance of each remove averaged over completed and uncompleted chains.



Small World Networks

Clustering Coefficient

Degree of a node does not include information about the neighbors of the node. Local clustering measures the density of connections between the immediate neighbors of the ego. In a random network to calculate C_i first we need to know the expected number of links of the ego network of i :

$$\langle L_i \rangle = p \frac{k_i(k_i - 1)}{2}$$

Thus, the local clustering coefficient in a random random network can be expressed by:

$$C_i = \frac{2\langle L_i \rangle}{k_i(k_i - 1)} = p = \frac{\langle k \rangle}{N}$$

- If we fix $\langle k \rangle$, the larger the network, the smaller is the node's local clustering.
- Local clustering coefficient of a node is independent from the number of connections of the node.

Small World Networks

Watts-Strogatz Model

What gives the word *small* in Small World Networks?

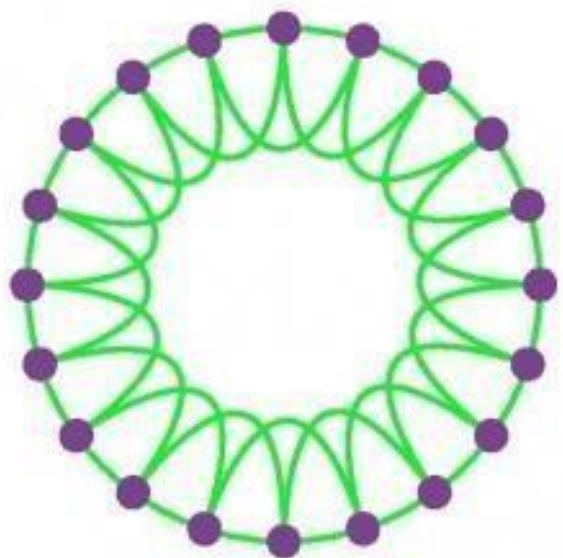
In real networks the average distance between two points is a logarithmic function of N.

What makes it a real lifenetwork?

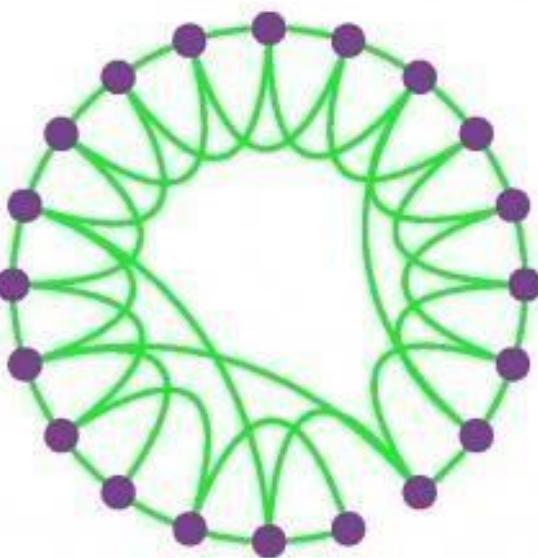
The average clustering coefficient of a real network is much higher than we would expect from a random network with the same properties of N and L.

a.

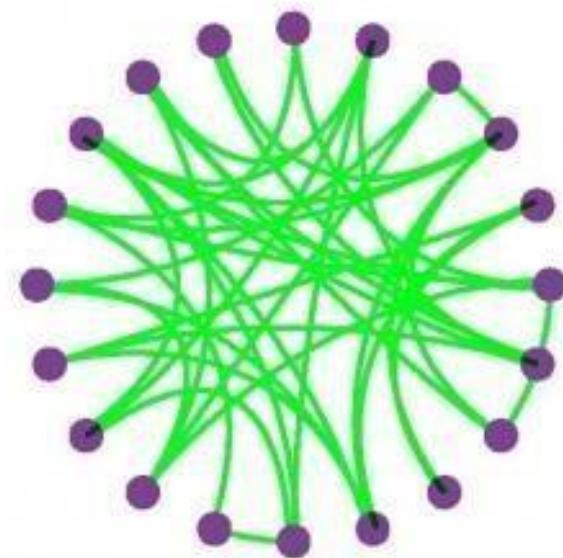
REGULAR

**b.**

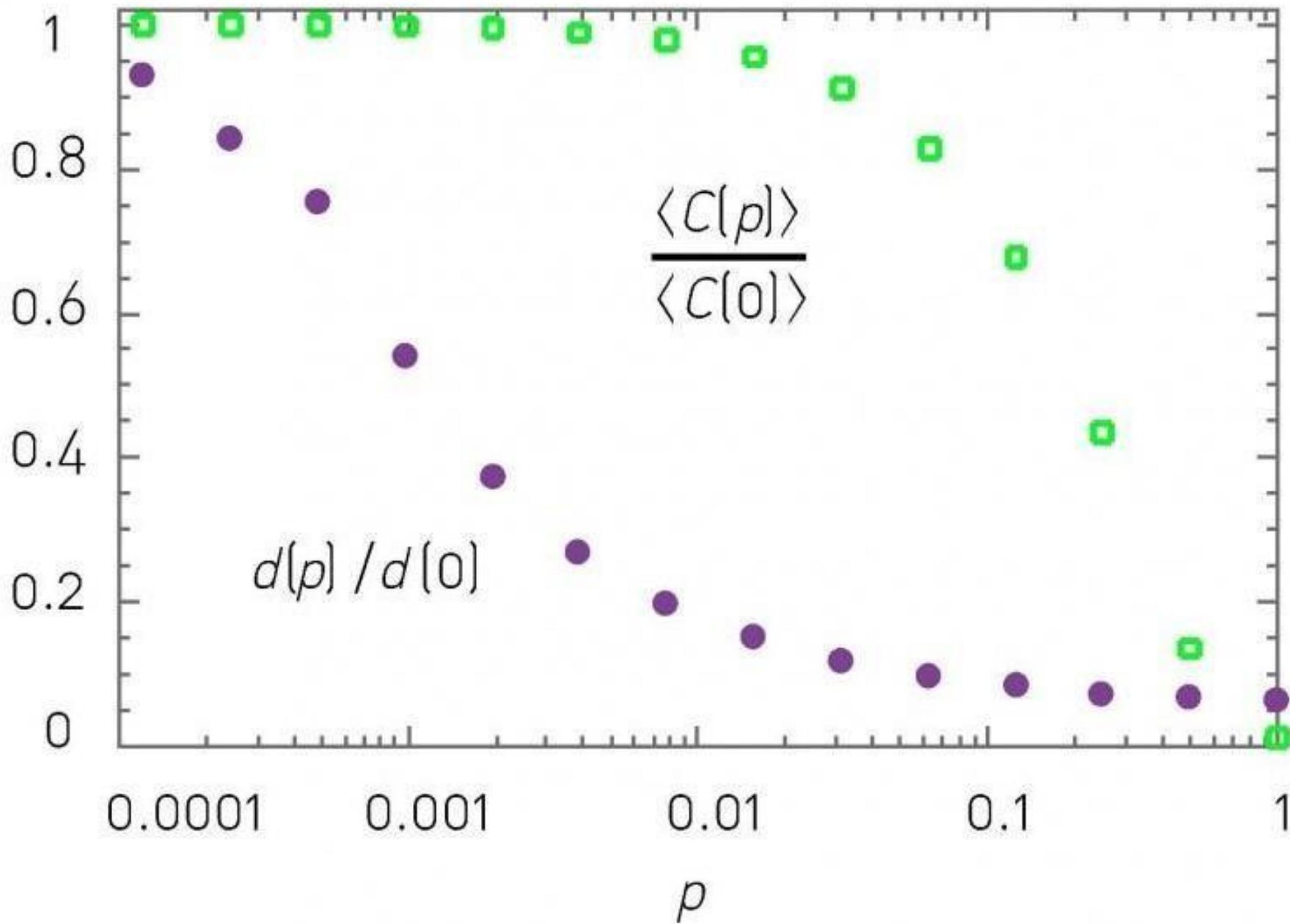
SMALL-WORLD

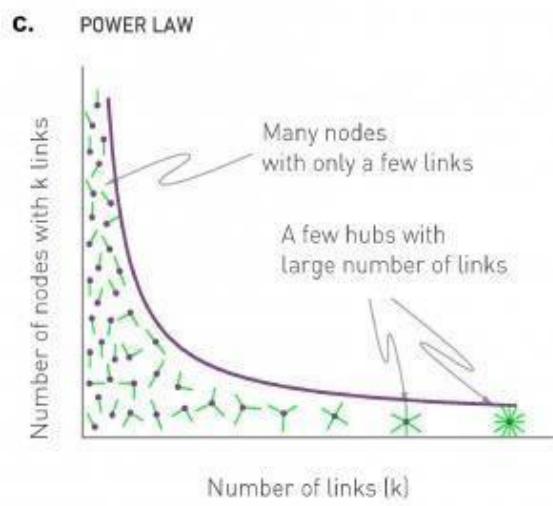
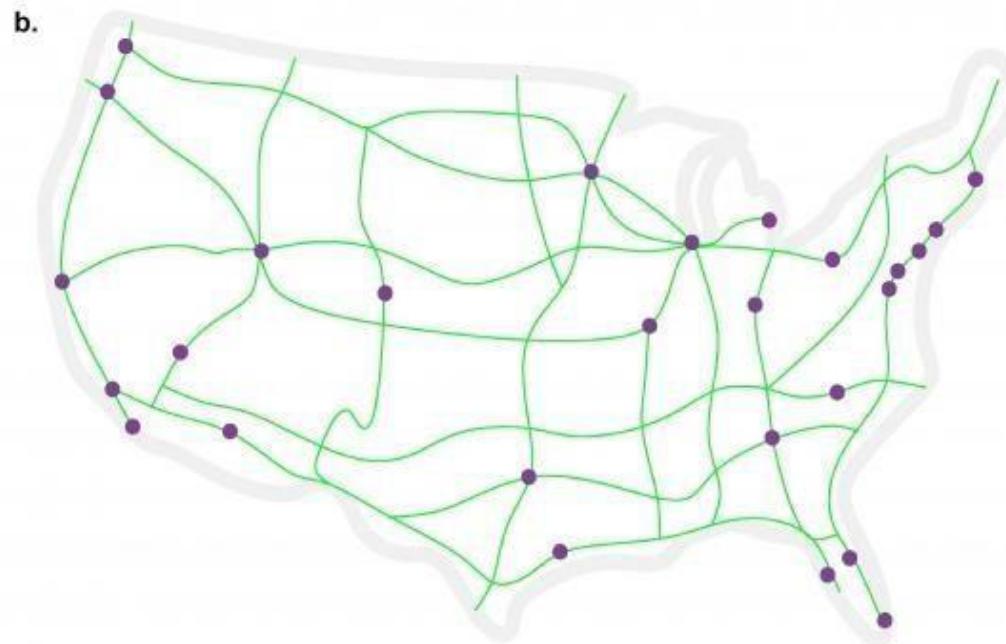
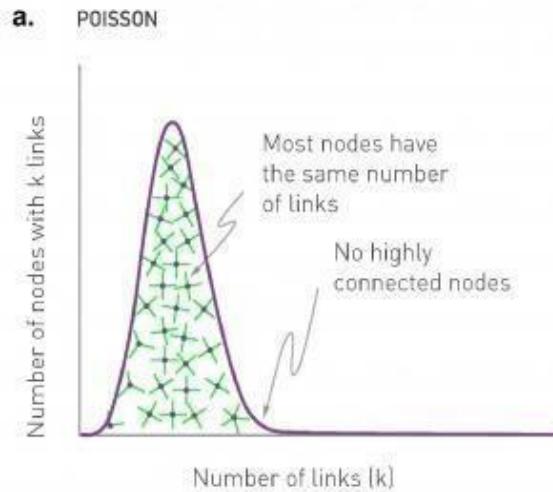
**c.**

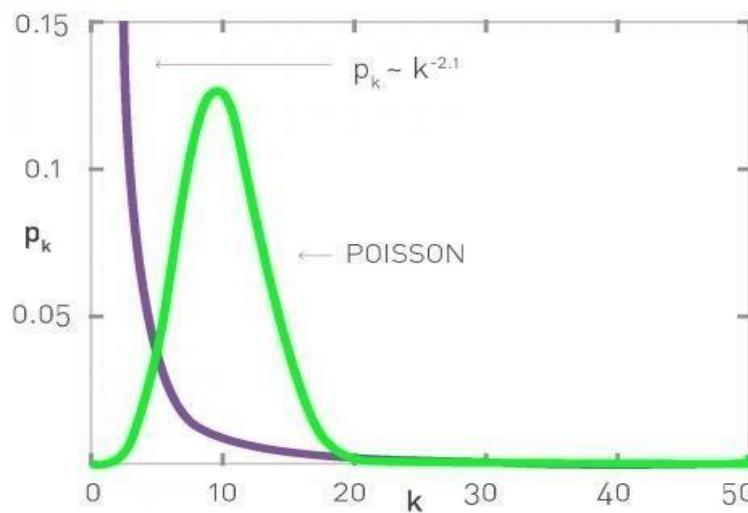
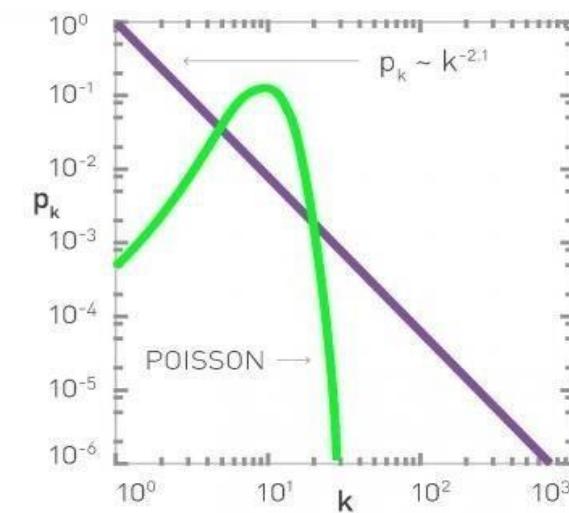
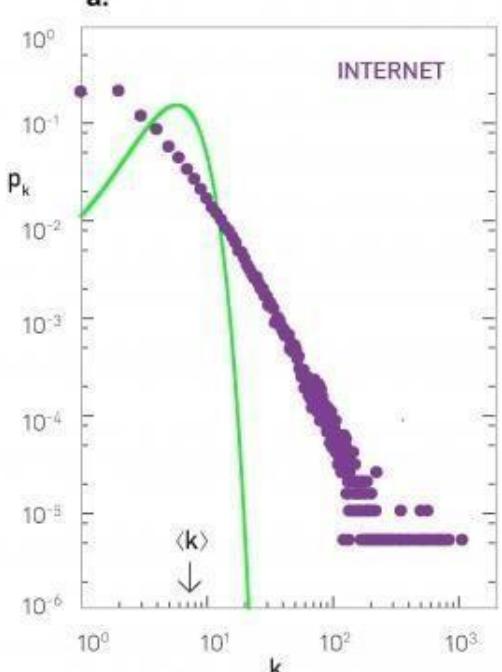
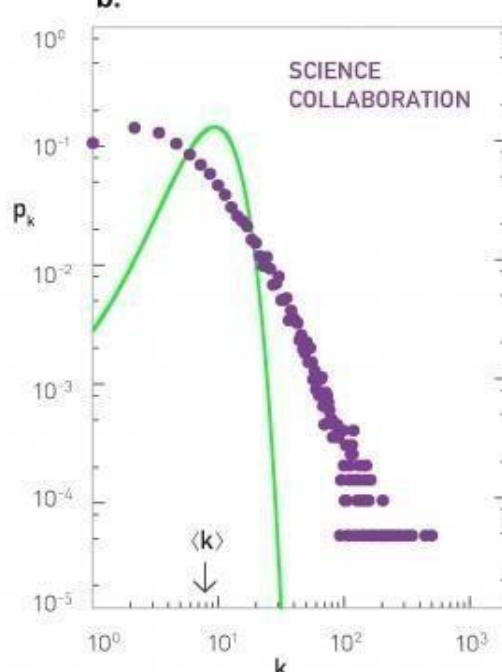
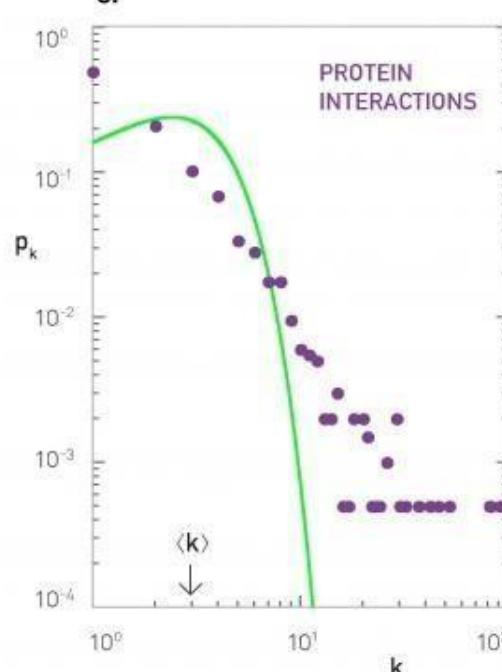
RANDOM

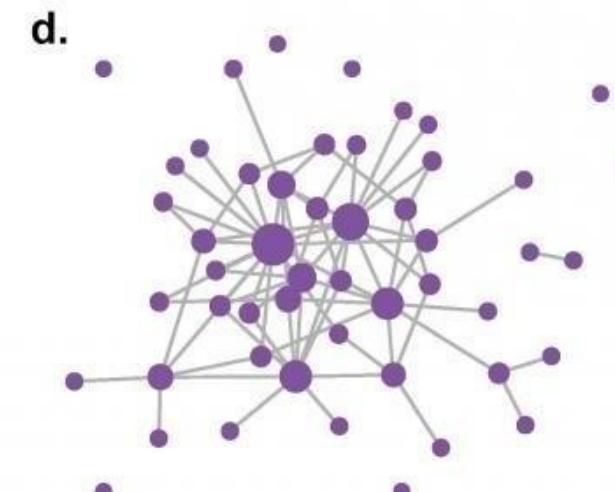
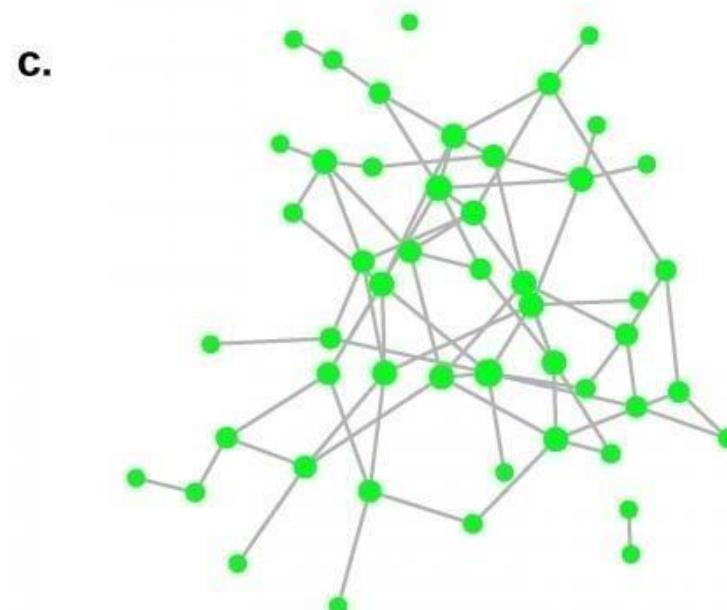
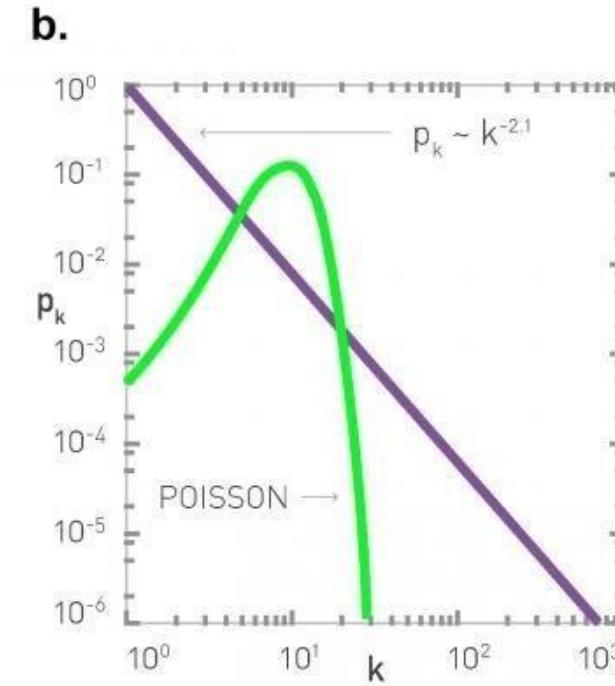
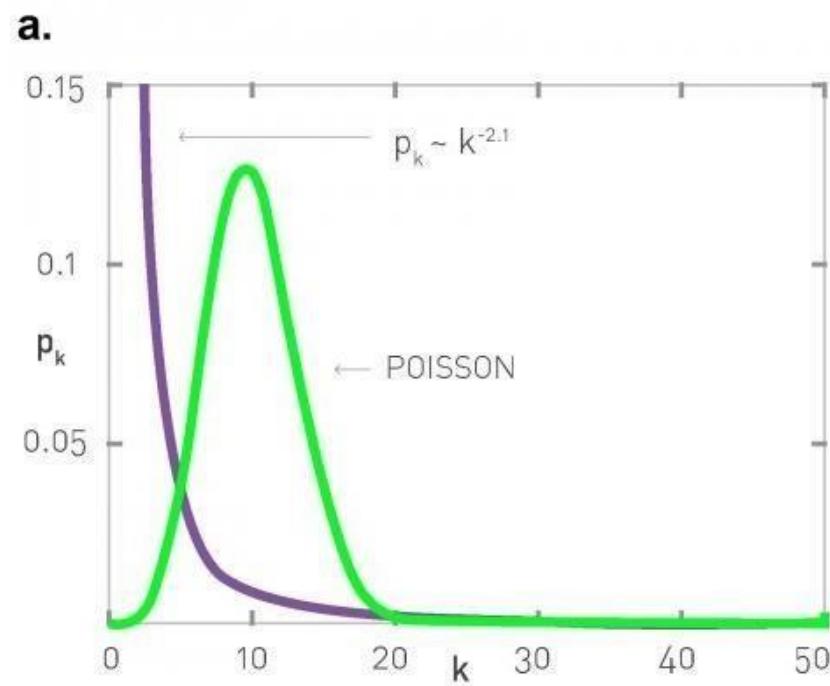
 $p = 0$  $p = 1$

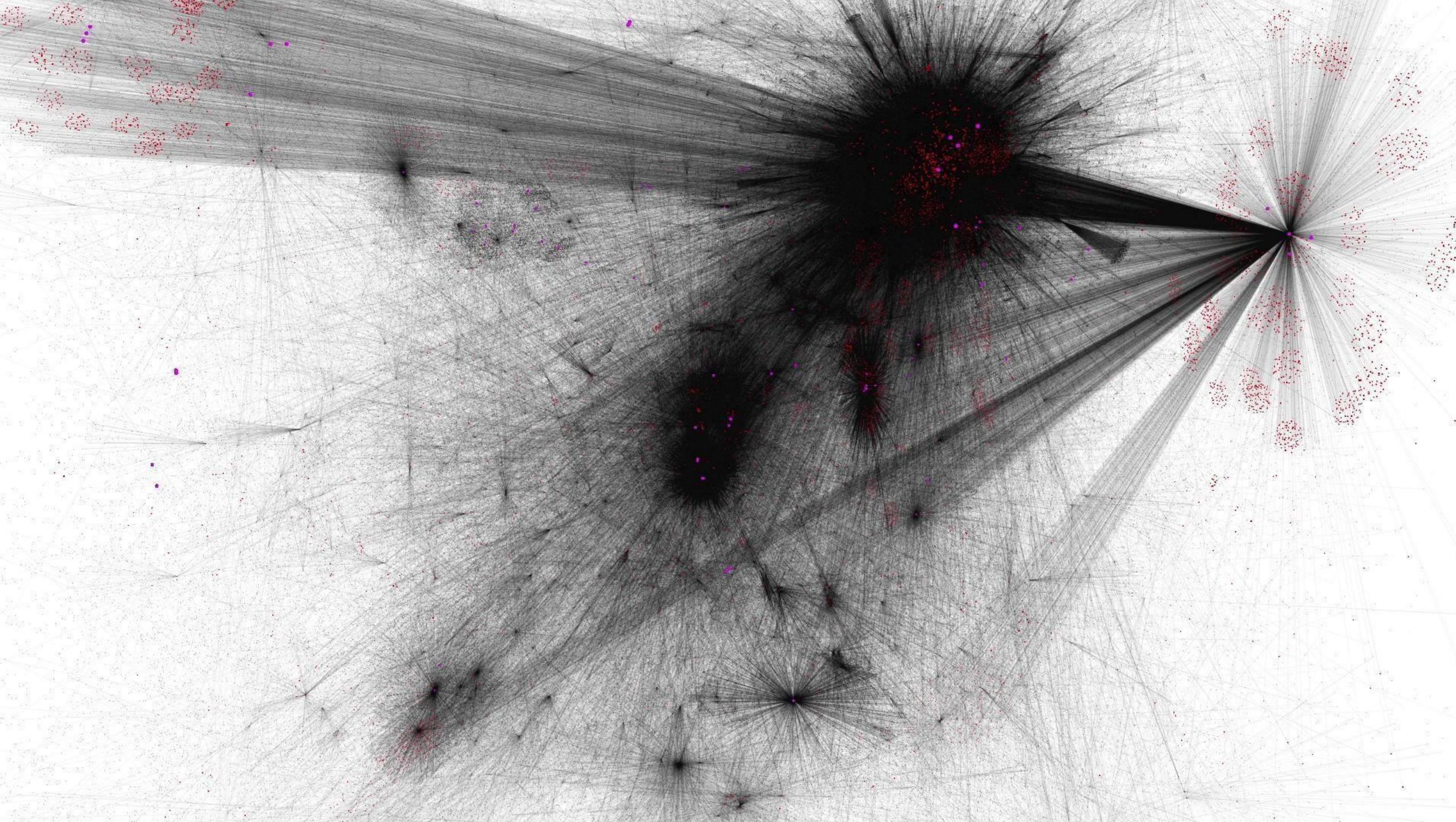
Increasing randomness

d.



a.**b.****a.****b.****c.**

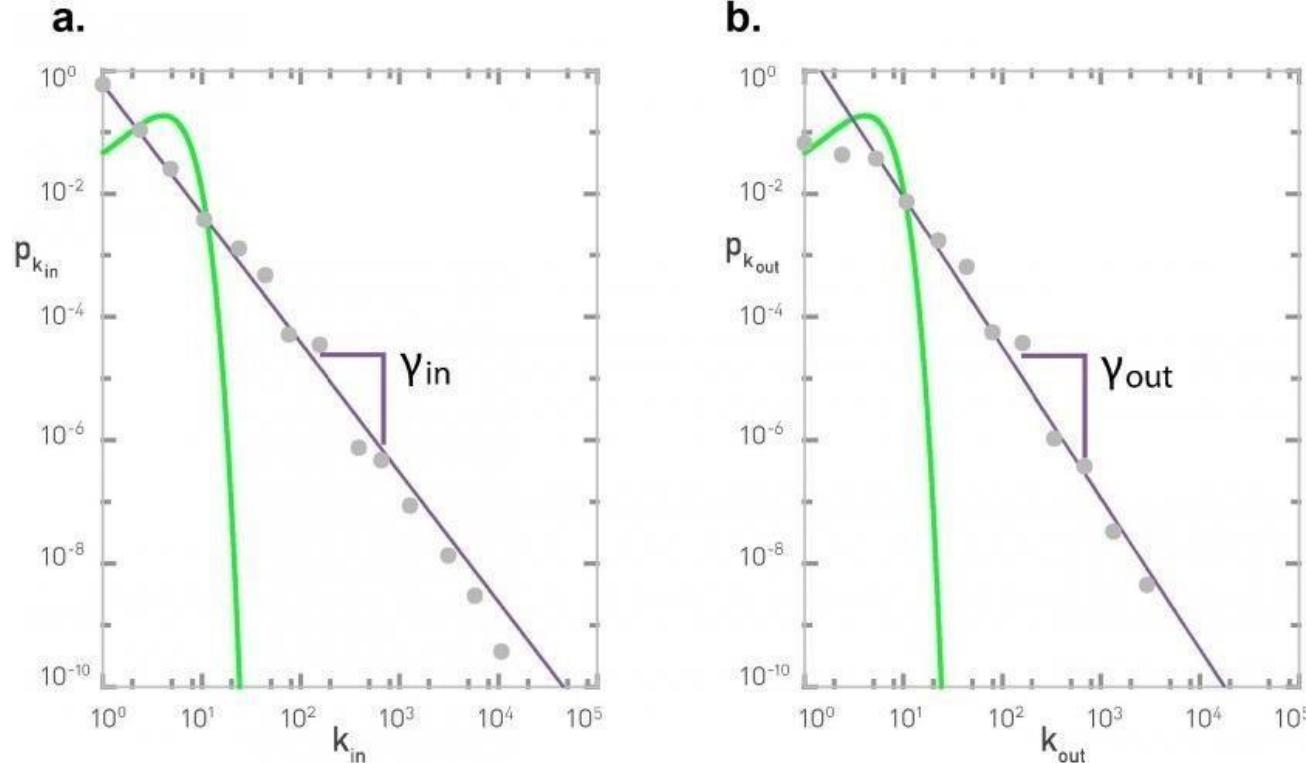




Scale-Free Networks

If the network of WWW was to be a random network than the degree distribution should follow Poisson distribution.

But it follows a power law distribution, with the exponent of gamma: $p_k = k^{-\gamma}$
Take the logarithm, we obtain: $\log p_k = -\gamma \log k$



Why do we have Hubs?

The Barabasi-Albert Model

For to every one who has will more be given, and he will have abundance; but from him who has not, even what he has will be taken away.

— Matthew 25:29, RSV



<https://www.maszol.ro/>

Growth

Real networks are getting larger and larger of a result of a growth process. While in random networks we usually have fix N.

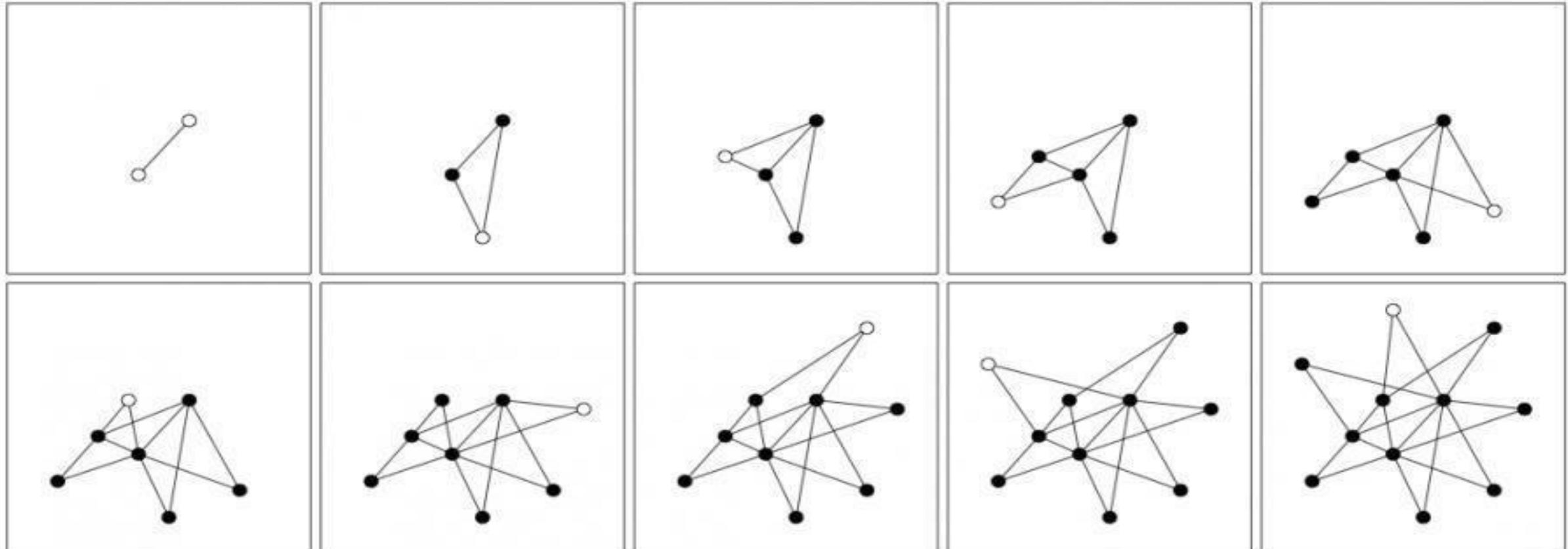
Preferential Attachment

In real networks nodes tend to connect to the more connected nodes.

Why do we have Hubs?

At each step we add new node to the network, with links connect to nodes already in the network.

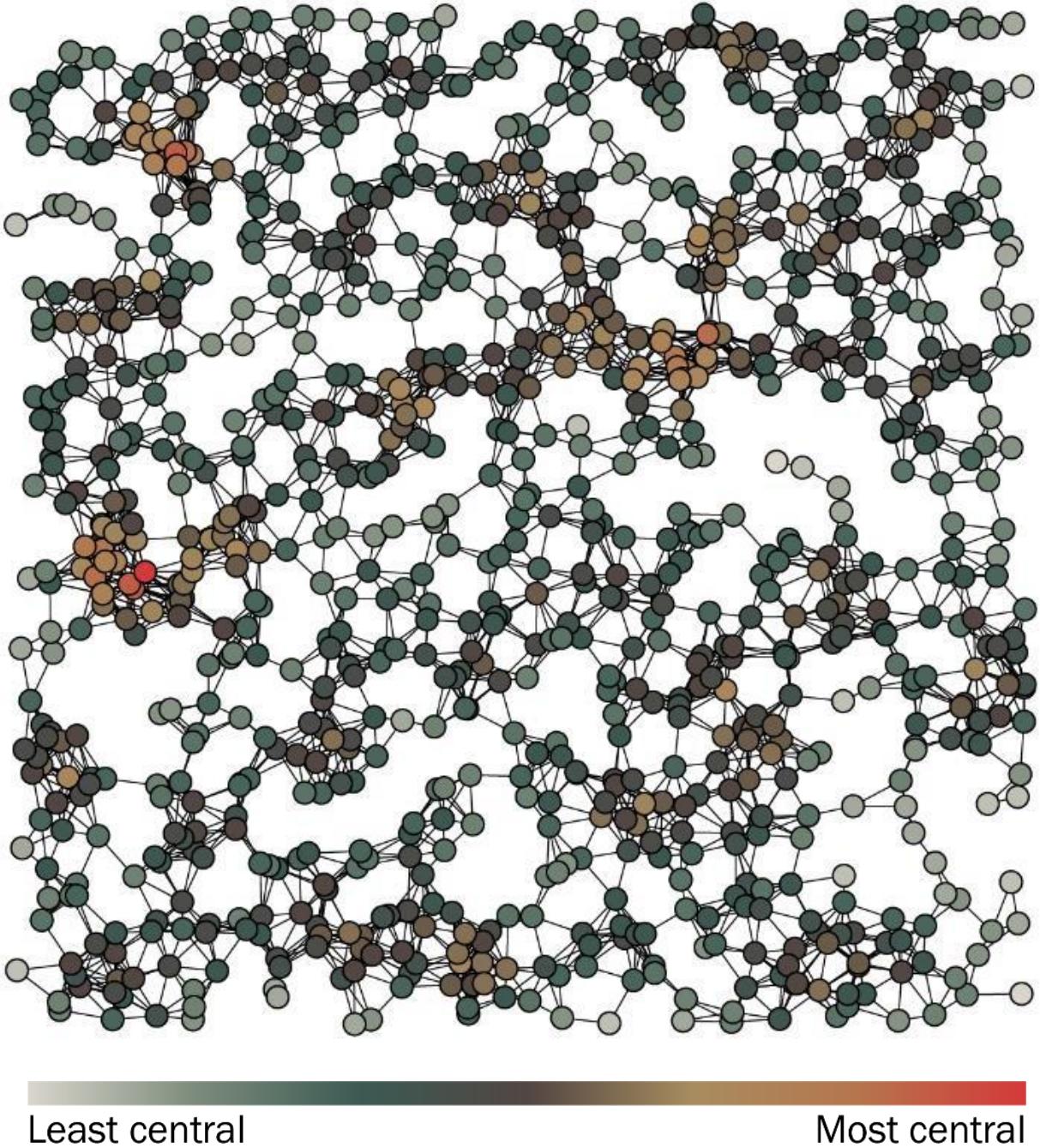
The probability of that a link of a new node connected to a present node is depends on the degree of the node: $\Pi(k_i) = \frac{k_i}{\sum_j k_j}$



Which node is a hub?

Node centrality

- Degree centrality



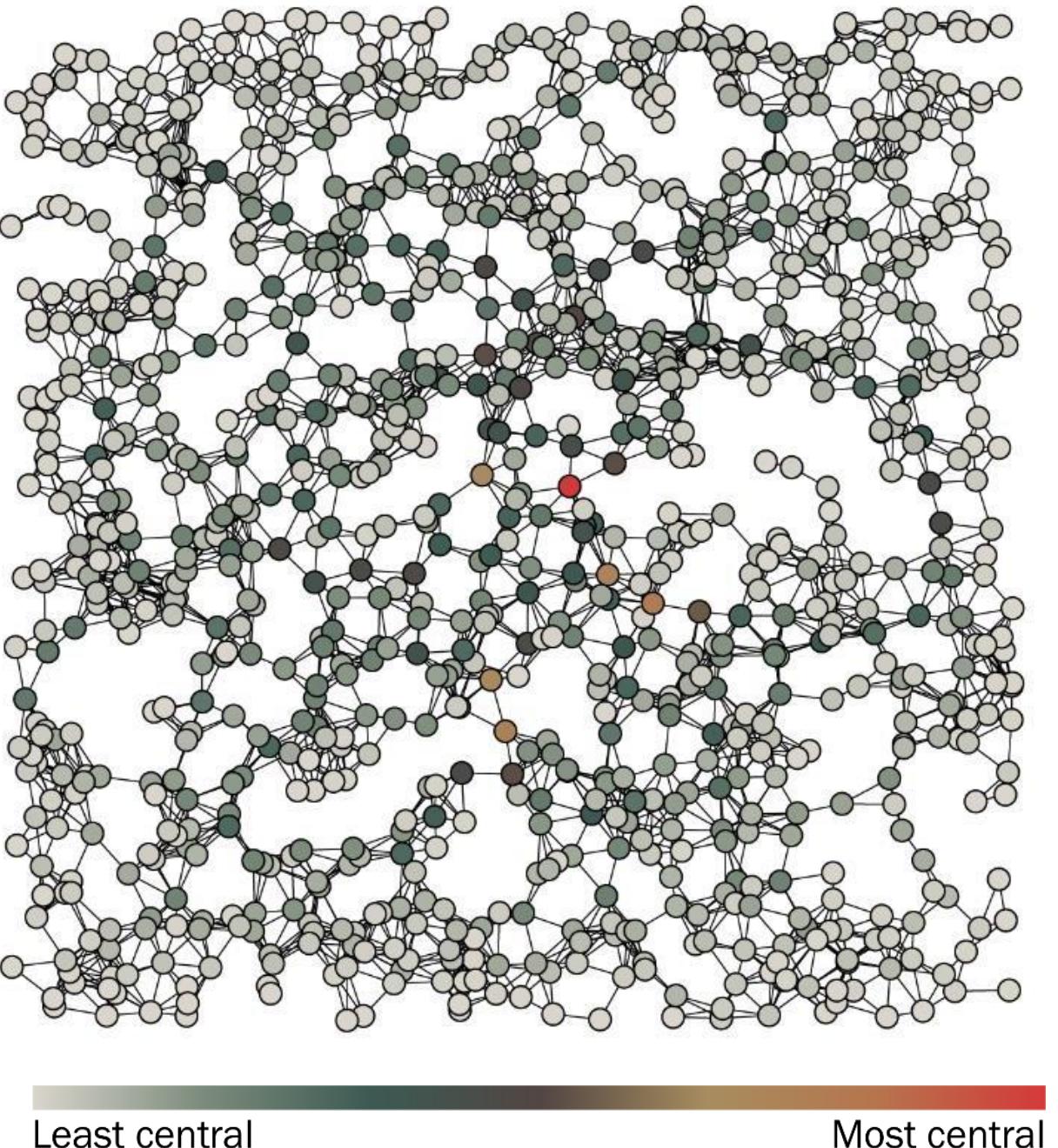
Which node is a hub?

Node centrality

- Betweenness centrality

Captures the node's role as a bridge between groups of nodes. Betweenness is about how critical a node is to the network's functioning as a bridge point between other parts of the network.

$$C_B(i) = \sum_{i \neq j \neq k} \sigma_{jk}(i)/\sigma_{jk}$$



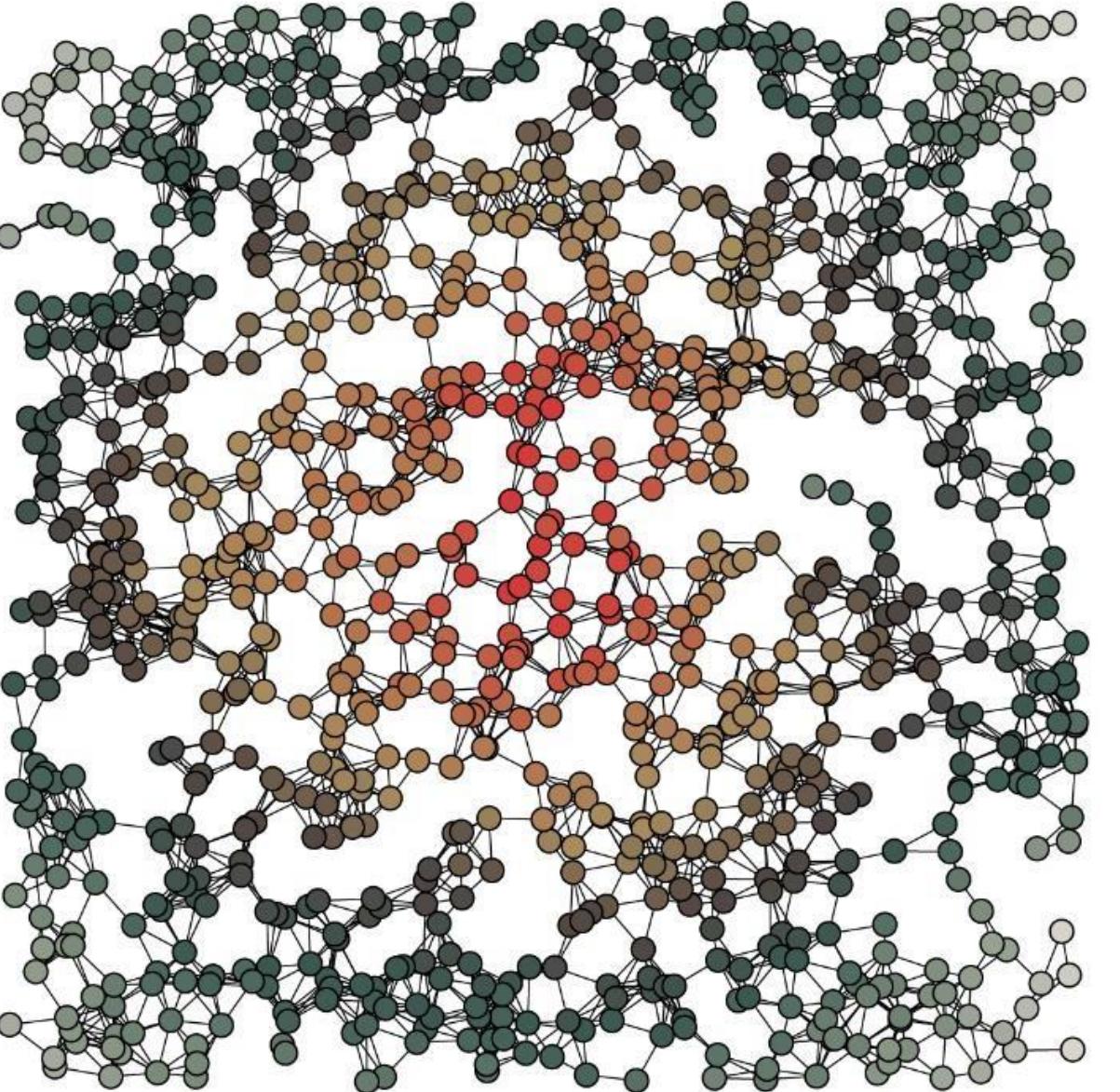
Which node is a hub?

Node centrality

- Closeness centrality

Closeness captures how close a node is to any other randomly selected node in the network. That is how quickly can the mode reach other nodes in the network.

$$C_C(i) = 1 / \sum_{i \neq j} d_{ij}$$



Least central

Most central

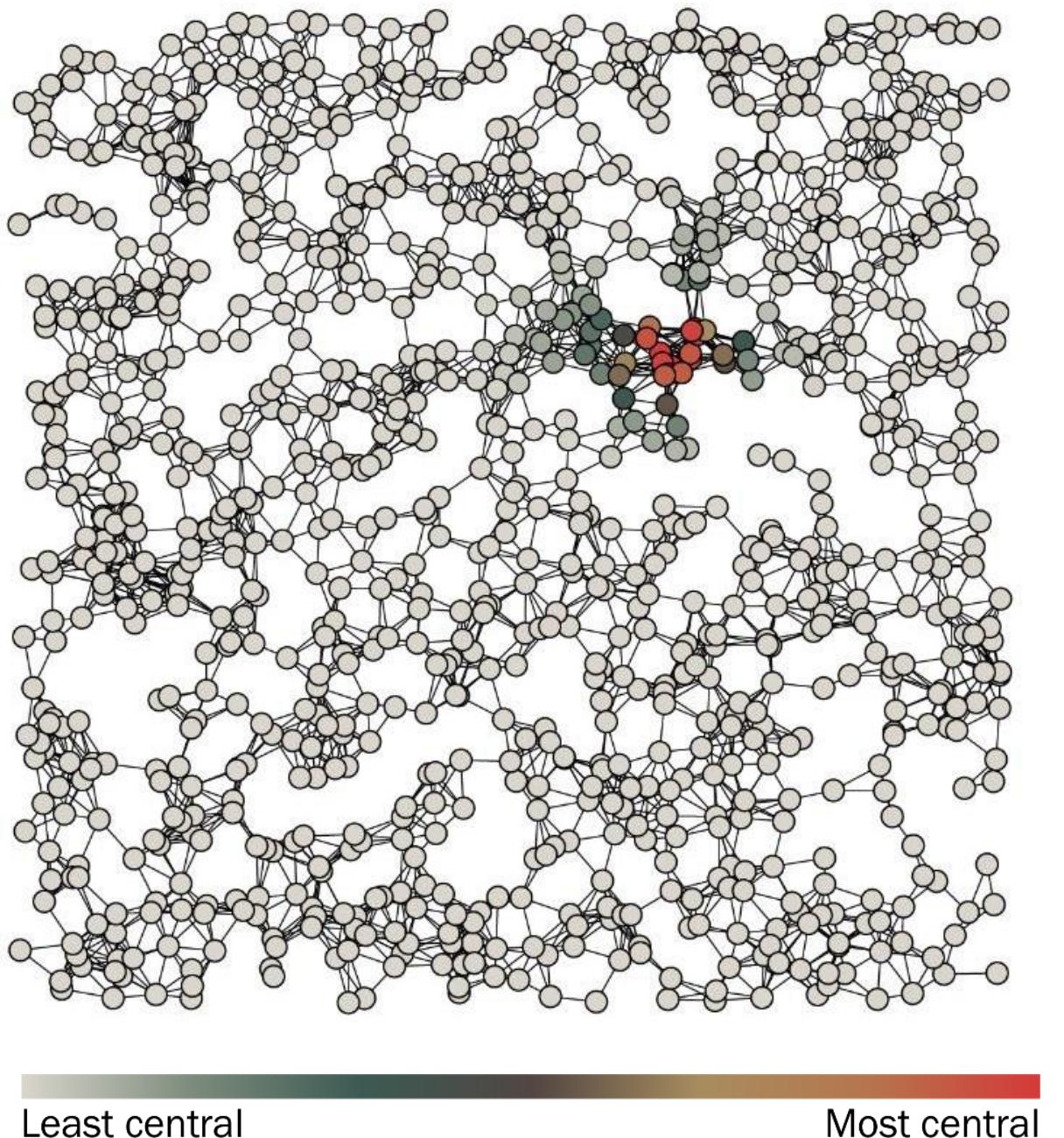
Which node is a hub?

Node centrality

- Eigenvector centrality

How important a node in a network based on the significance of the nodes that given node is connected to. So instead of looking on the number of connections it represents the value of the connections.

$$C_E(i) = x_i = \frac{1}{\lambda} \sum_{j \in M(i)} x_j = \frac{1}{\lambda} \sum_{j \in G} A_{ij} x_j$$



Which node is a hub?

Betweenness

Captures the node's role as a bridge between groups of nodes. Betweenness is about how critical a node is to the network's functioning as a bridge point between other parts of the network.

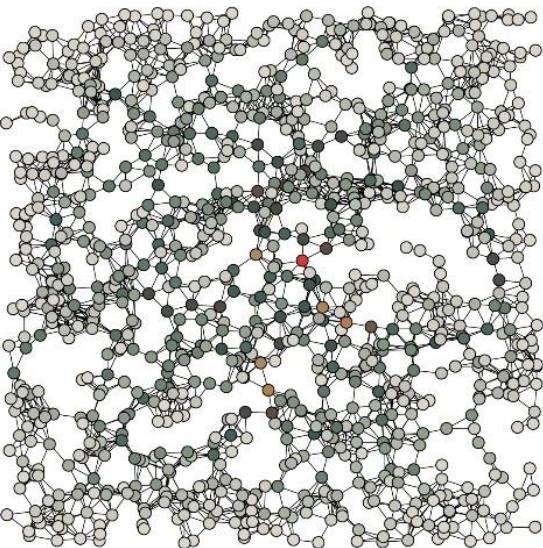
Closeness

Closeness captures how close a node is to any other randomly selected node in the network. That is how quickly can the node reach other nodes in the network.

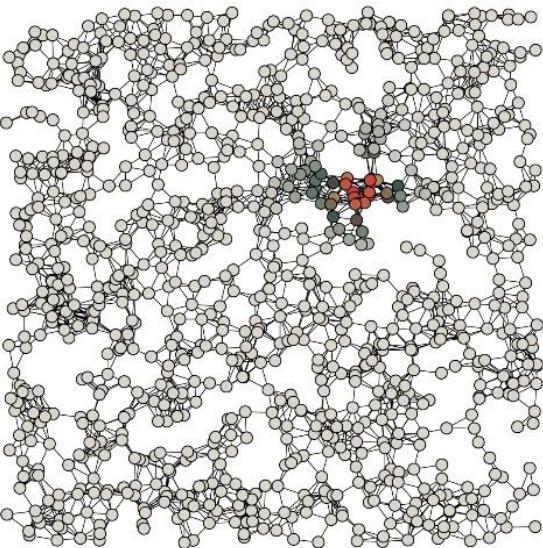
Eigenvector centrality

How important a node in a network based on the significance of the nodes that given node is connected to. So instead of looking on the number of connections it represents the value of the connections.

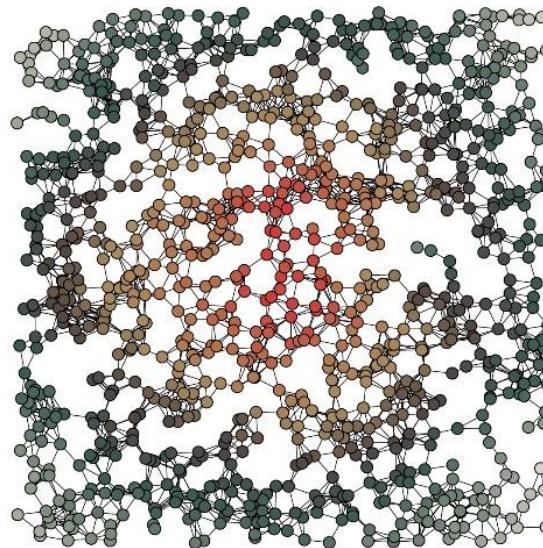
Degree centrality



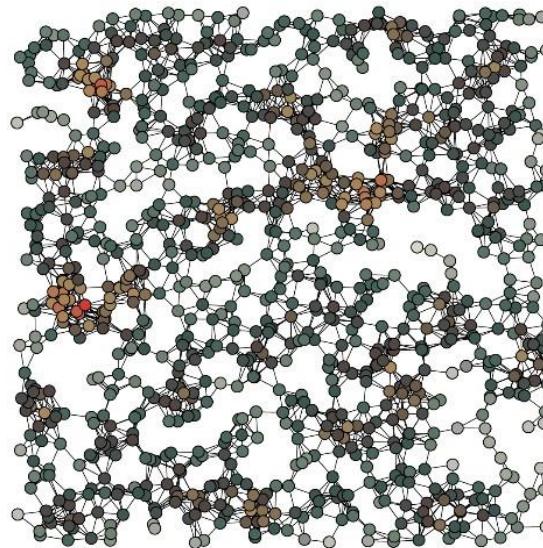
A Betweenness



C Eigenvector



B Closeness



D Degree

Least central

Most central