

HOME ABOUT PARTICIPATE COURSES SPEAKERS DAILY SCHEDULE VENUE SOCIAL ACTIVITIES EXPLORE ATHENS

13-19 July 2017, Athens, Greece

1st ACM Europe Summer School | Data Science

Social Network Analysis

Dino Pedreschi



UNIVERSITÀ DI PISA

Università di Pisa & ISTI-CNR

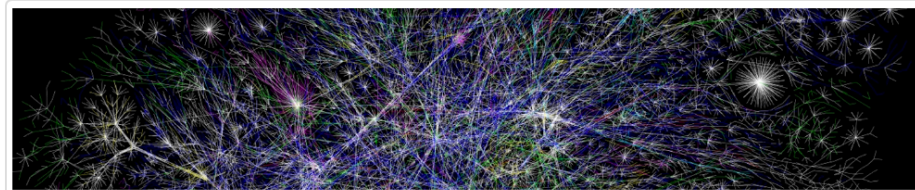


ISTITUTO DI SCIENZA E TECNOLOGIE
DELL'INFORMAZIONE "A. FAEDO"

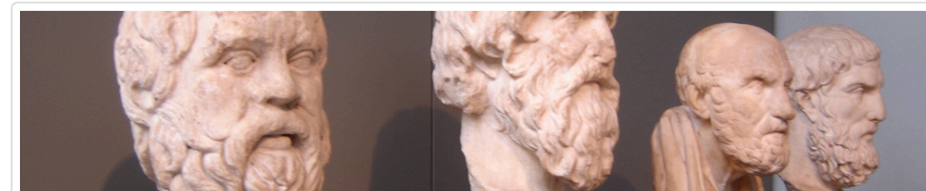
<http://kdd.isti.cnr.it>



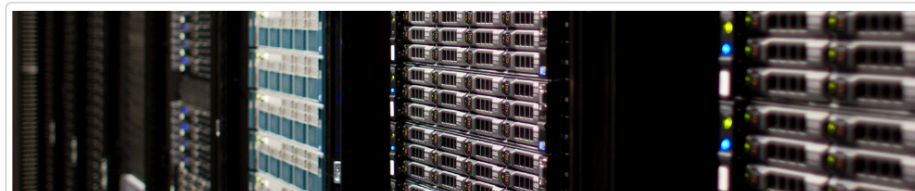
Mobility Data Mining for Science of Cities



Social Network Analysis and Visual Analytics



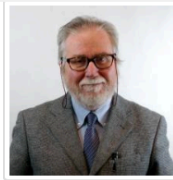
Ethical Data Mining



Analytical Platforms and Infrastructures for Social Mining



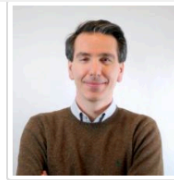
Dino Pedreschi
Full Professor
pedre@di.unipi.it



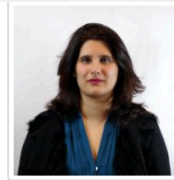
Franco Turini
Full Professor
turini@di.unipi.it



Fosca Giannotti
Director of Research
fosca.giannotti@isti.cnr.it



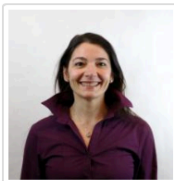
Salvatore Ruggieri
Full Professor
ruggieri@di.unipi.it



Anna Monreale
Assistant Professor
annam@di.unipi.it



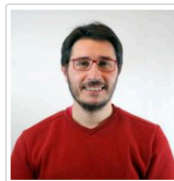
Alina Sirbu
Assistant Professor
alina.sirbu@unipi.it



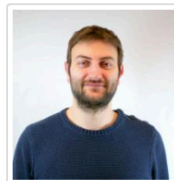
Barbara Furletti
Researcher
barbara.furletti@isti.cnr.it



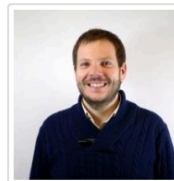
Mirco Nanni
Researcher
mirco.nanni@isti.cnr.it



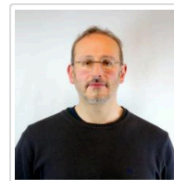
Salvatore Rinzivillo
Researcher
rinzivillo@isti.cnr.it



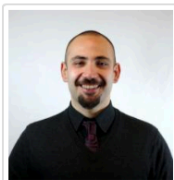
Roberto Trasarti
Researcher
roberto.trasarti@isti.cnr.it



Paolo Cintia
Post Doc
paolo.cintia@isti.cnr.it



Valerio Grossi
Post Doc
vgrossi@di.unipi.it



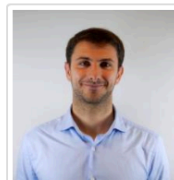
Luca Pappalardo
Post Doc
luca.pappalardo@isti.cnr.it



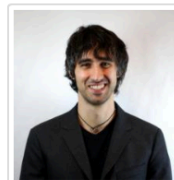
Giulio Rossetti
Post Doc
giulio.rossetti@isti.cnr.it



Alessandro Baroni
Ph.D. Student
baroni@di.unipi.it



Lorenzo Gabrielli
Ph.D. Student
lorenzo.gabrielli@isti.cnr.it



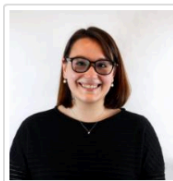
Riccardo Guidotti
Ph.D. Student
riccardo.guidotti@isti.cnr.it



Ioanna Miliou
Ph.D. Student
ioanna.miliou@for.unipi.it



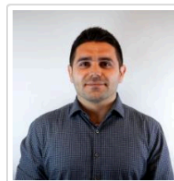
Letizia Milli
Ph.D. Student
letizia.milli@isti.cnr.it



Laura Pollacci
Ph.D. Student
laura.pollacci@di.unipi.it



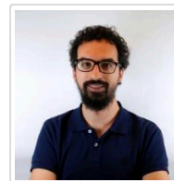
Francesca Pratesi
Ph.D. Student
francesca.pratesi@isti.cnr.it



Farzad Vaziri
Ph.D. Student
vaziri@di.unipi.it



Viola Bachini
Collaborator
viola.bachini@isti.cnr.it



Daniele Fadda
Collaborator
daniele.fadda@isti.cnr.it





Social Mining & Big Data Analytics

H2020 - www.sobigdata.eu

September 2015- August 2019



ISTITUTO DI SCIENZA E TECNOLOGIE
DELL'INFORMAZIONE "A. FAEDO"



UNIVERSITÀ DI PISA



SCUOLA
NORMALE
SUPERIORE

Institutions
Markets
Technologies

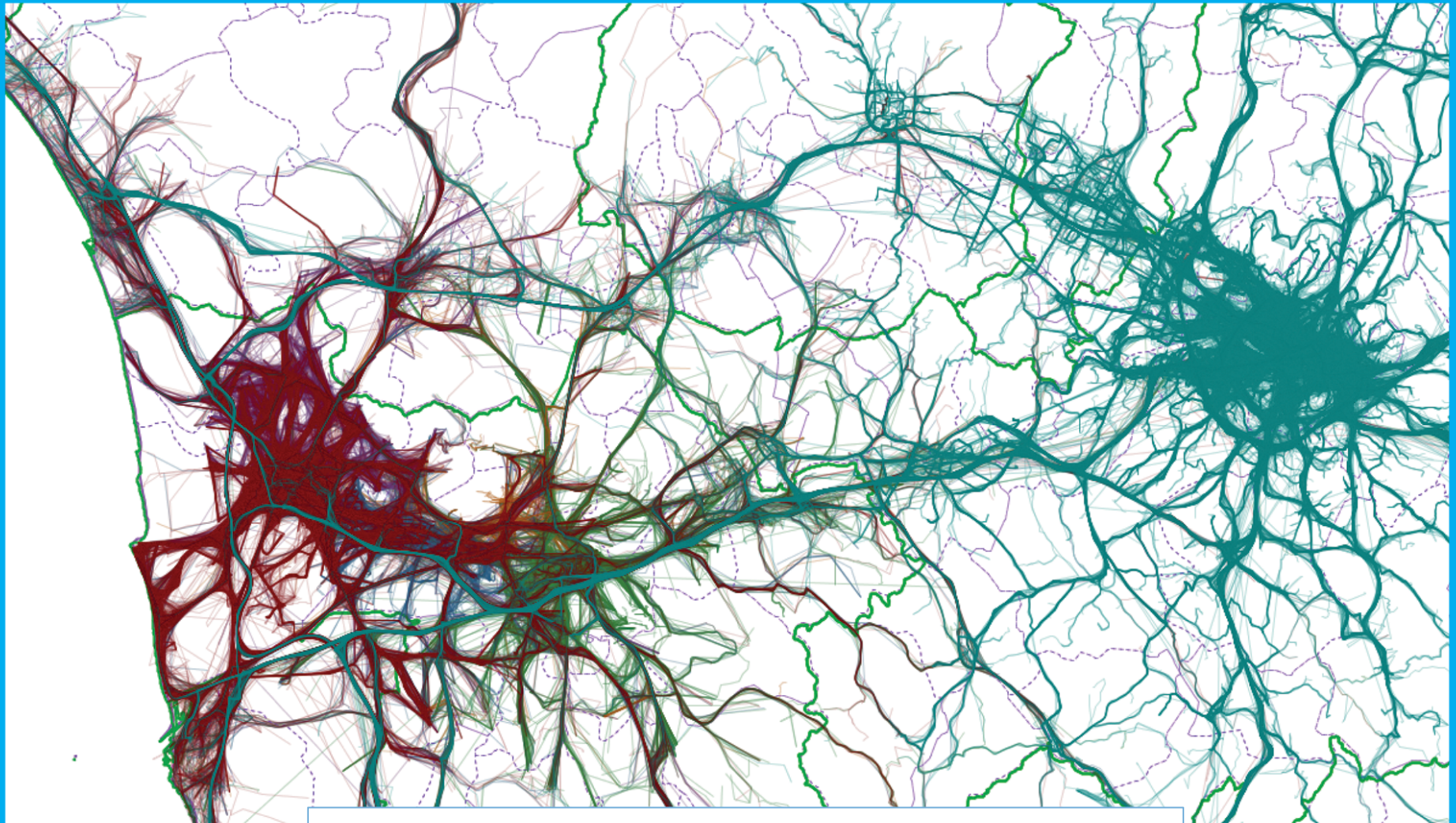
IMT

INSTITUTE
FOR ADVANCED
STUDIES
LUCCA

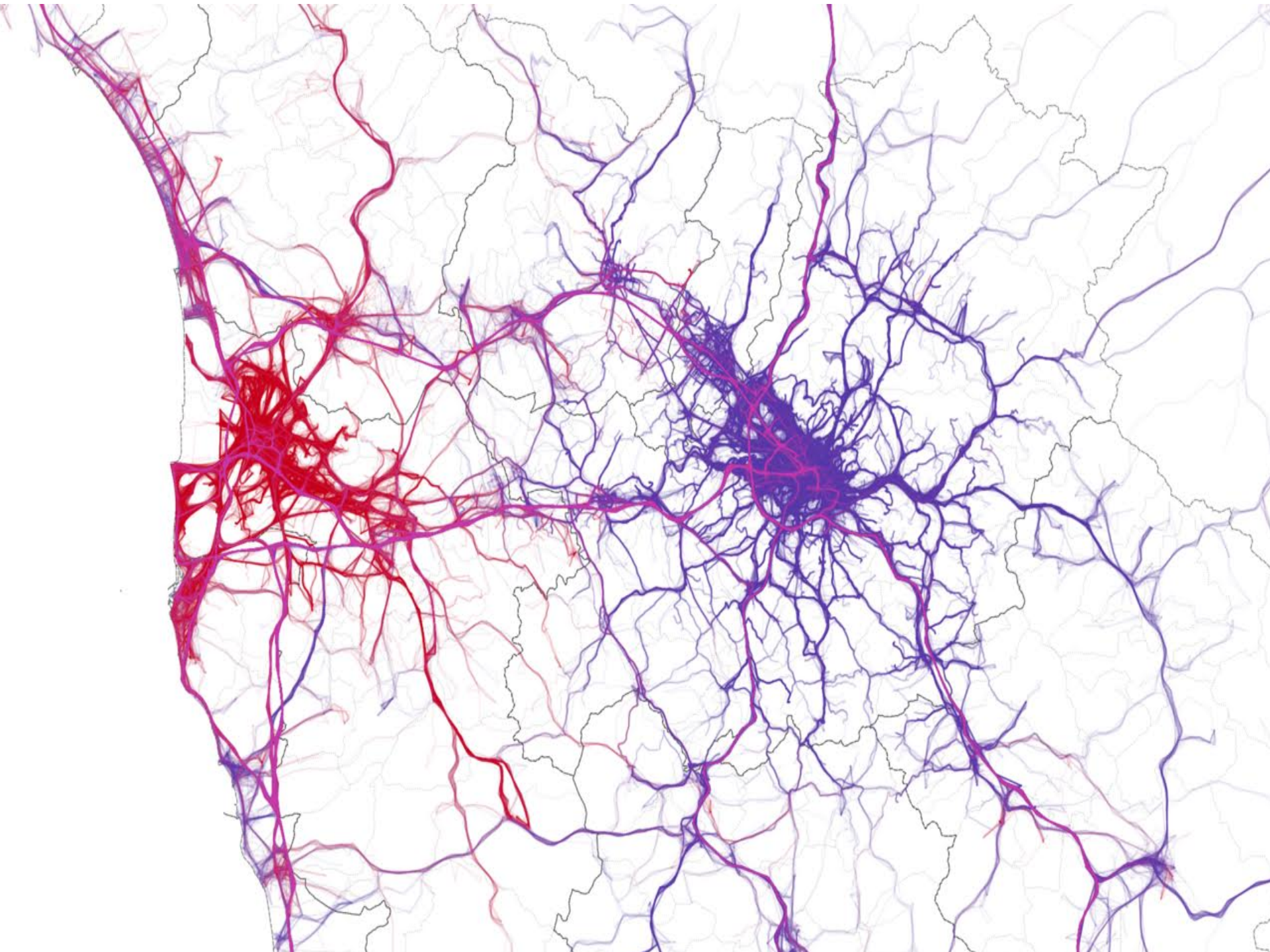




Exploratory: Big Data for City of Citizens



Personal Mobility, Social + Mobility, Personal Sensing



Pisa

Pisa



NUMBER OF VEHICLES

5,615

Source: OCTO Telemati...



RESIDENTS

89,694

Source: census 2001



DAYS

31

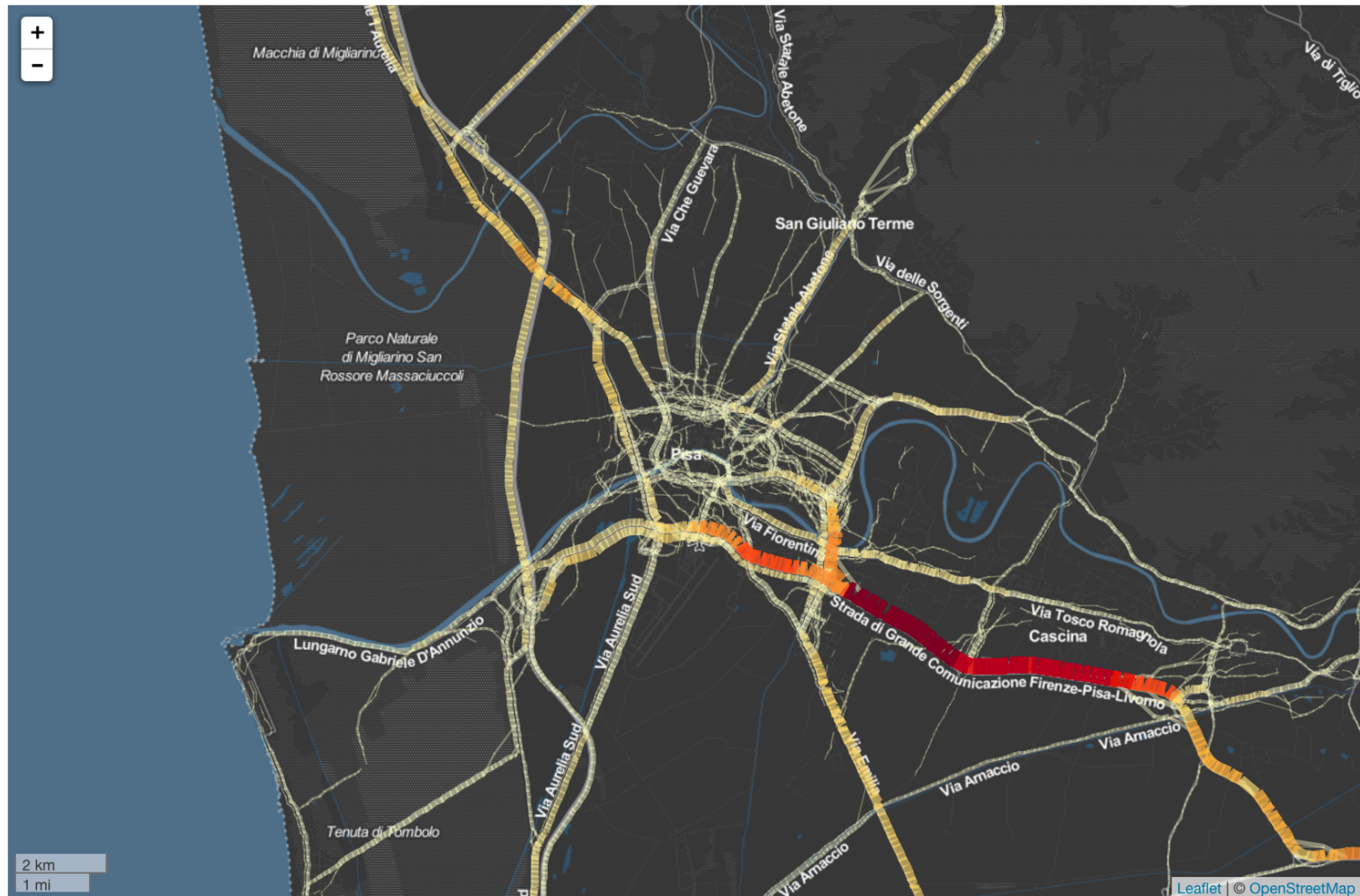
May 1-31, 2011



SYSTEMATICS TRAJs

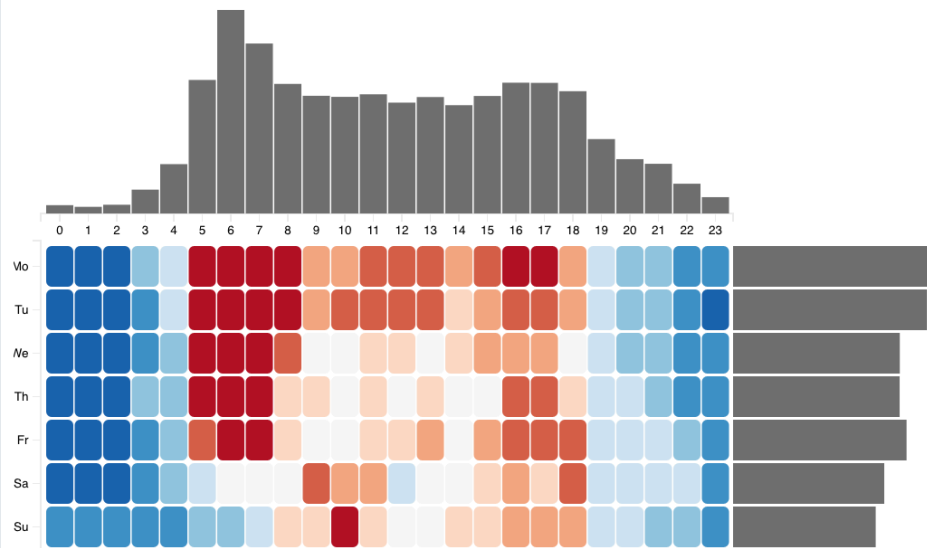
38.88%

Flows and Traffic

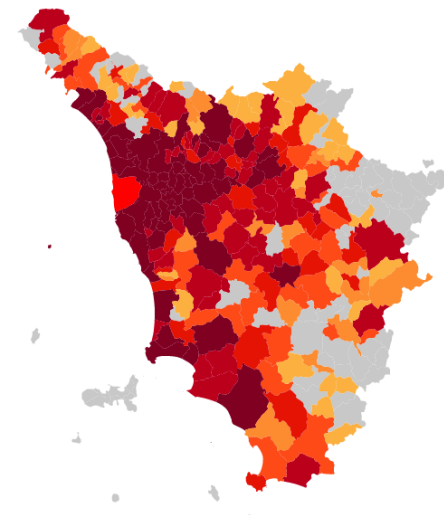


Flows of traffic exiting from the city.

Temporal Matrix



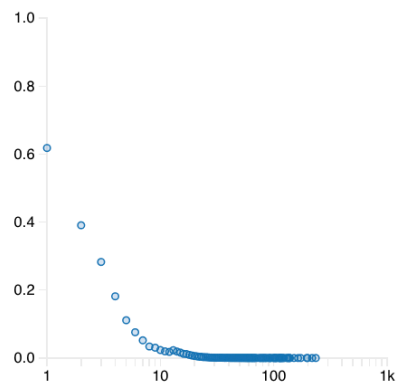
OD Map



Map of origins and destinations

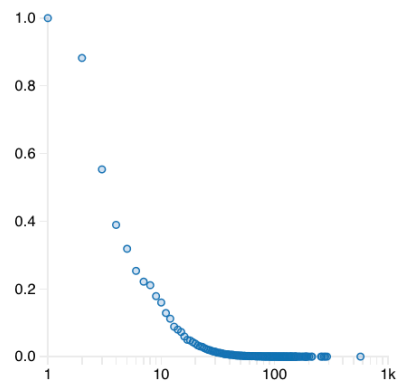
Lengths of Trajs

set of buttons here



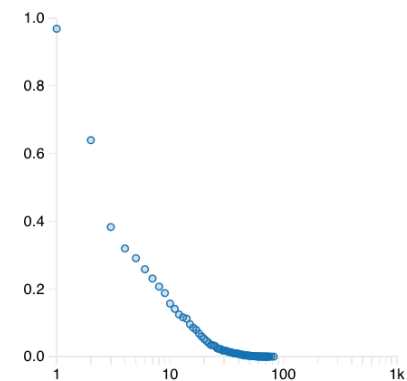
Lengths of Trajs

Duration of Trajs



Distribution of durations

Speeds of Trajs



Distribution of speeds



**Analizzare l'influenza dei grandi attrattori sulla
mobilità dei territori circostanti**

CASO STUDIO:

gli aeroporti di Firenze e Pisa e la propensione dei
residenti toscani all'uso delle due infrastrutture.



ACCESS POINT PISA

principali access point alla città di Pisa
degli utenti che si recano all'aeroporto Galilei
usando l'automobile



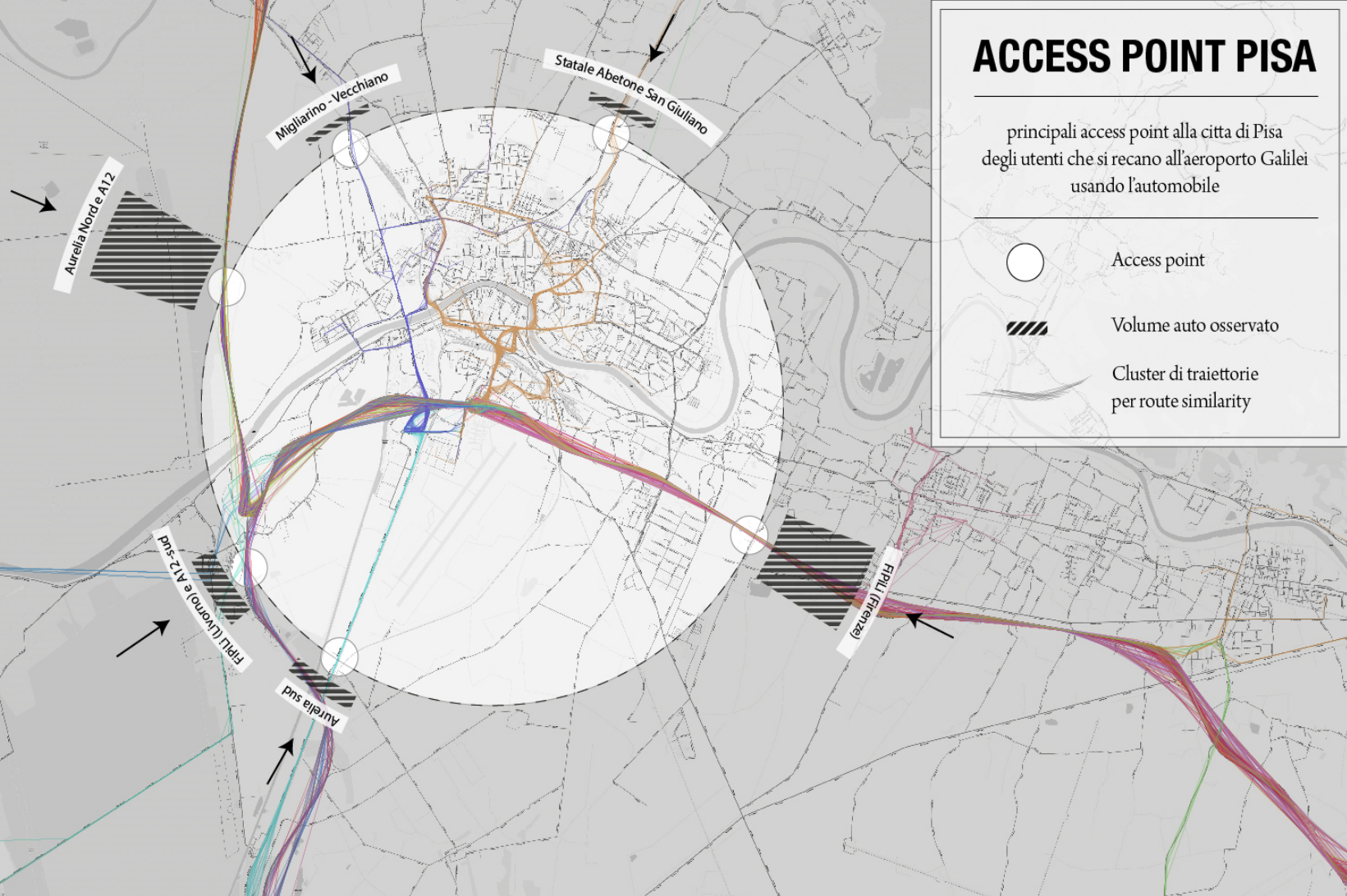
Access point

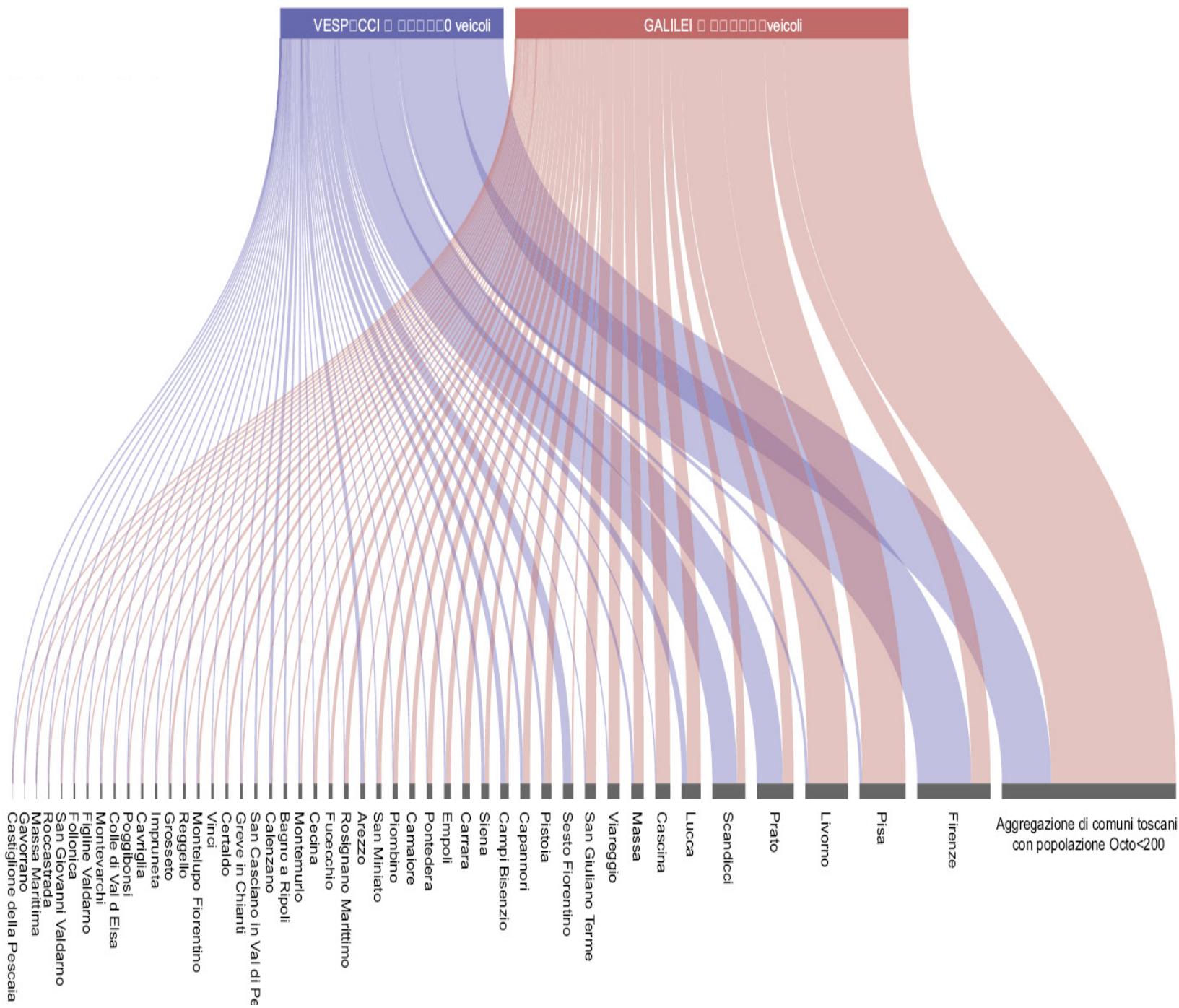


Volume auto osservato

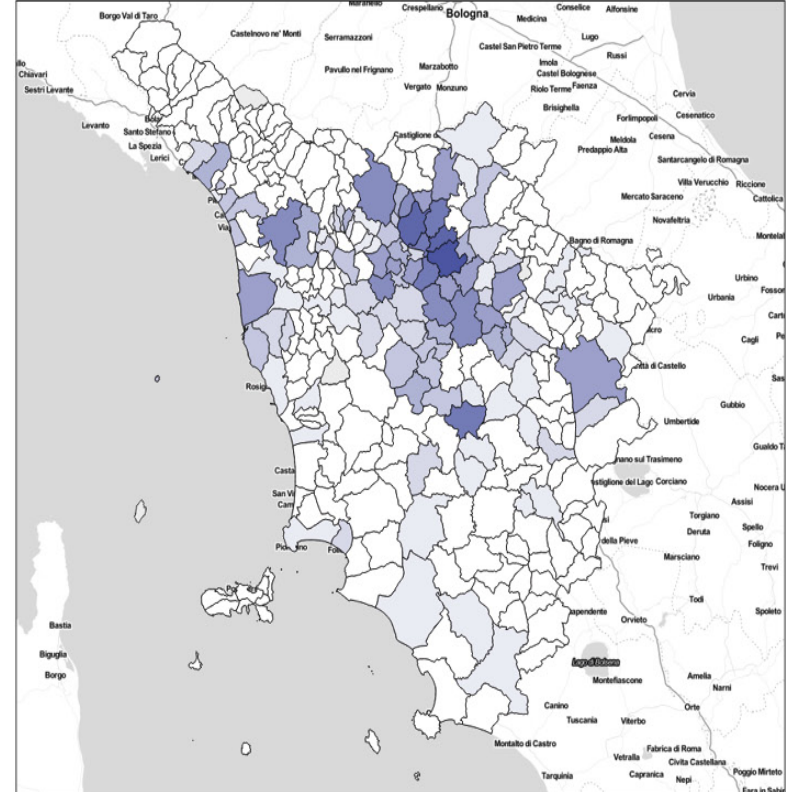
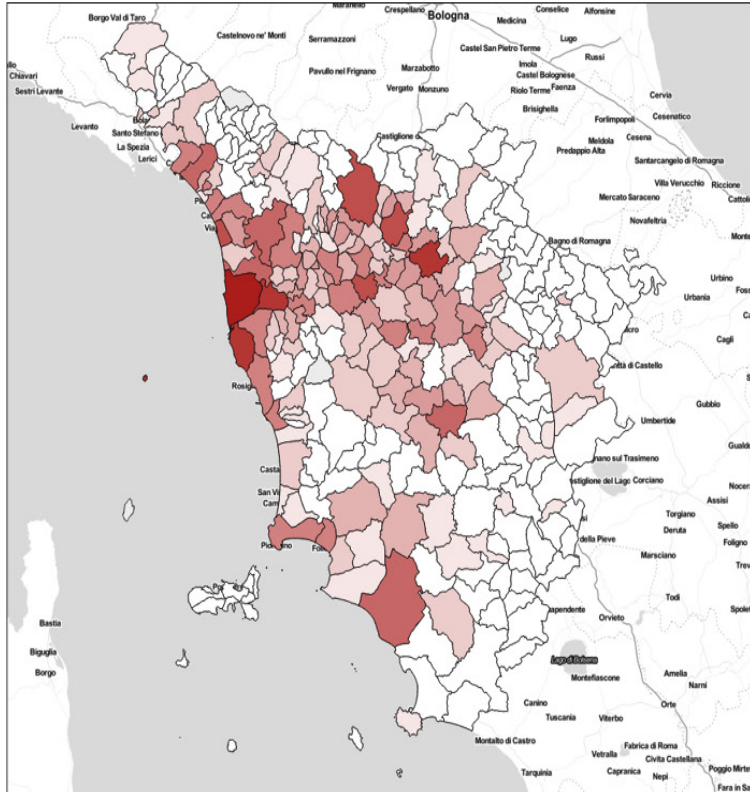


Cluster di traiettorie
per route similarity





Attractiveness of Galilei vs. Vespucci



Modeling Investments and Attractiveness on Tuscan Airports

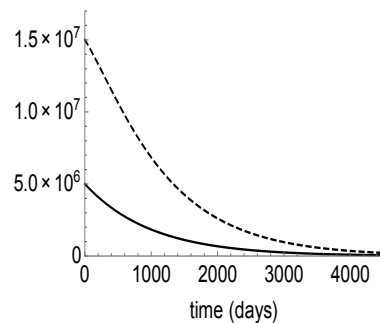
An intertwined system based on investment and attractiveness

$$\begin{aligned} \frac{d}{dt}A &= s(mF - (k + e)A), & A &\rightarrow \text{Attractiveness of airport} \\ \frac{d}{dt}F &= -rF + re\frac{bA}{1 + bhA}. & F &\rightarrow \text{Number of passengers served} \end{aligned}$$

Attractiveness is proportional to the cost of operating the airport (k) and the extra investments (e)

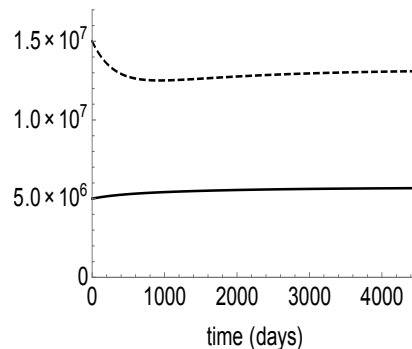
Simple case: non spatial model

No extra investments ($e=0$)



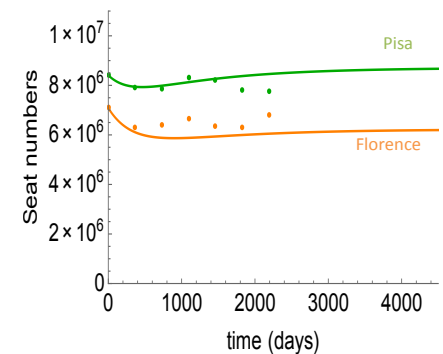
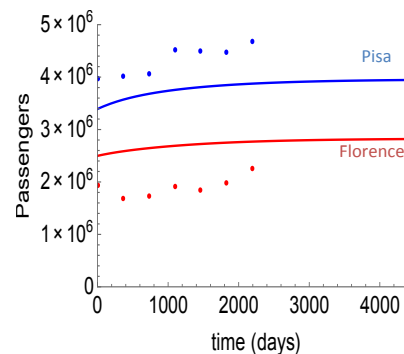
(a) $e = 0$

Structural investments ($e=0.05$)



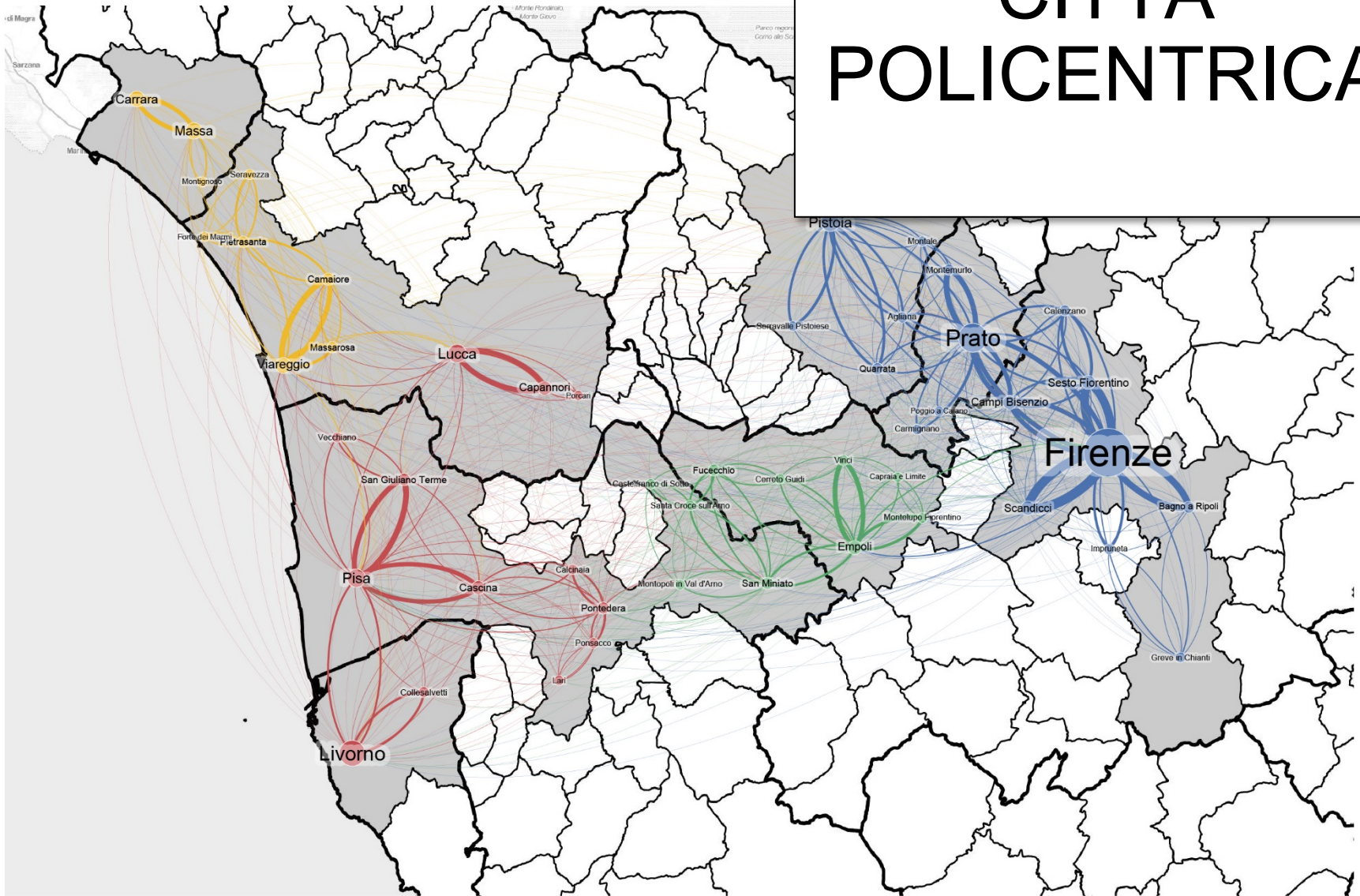
(b) $e = 0.05$

Spatial model: two airports, two populations



The two airports reach an equilibrium: neither of the two is overwhelming the other

CITTÀ POLICENTRICA

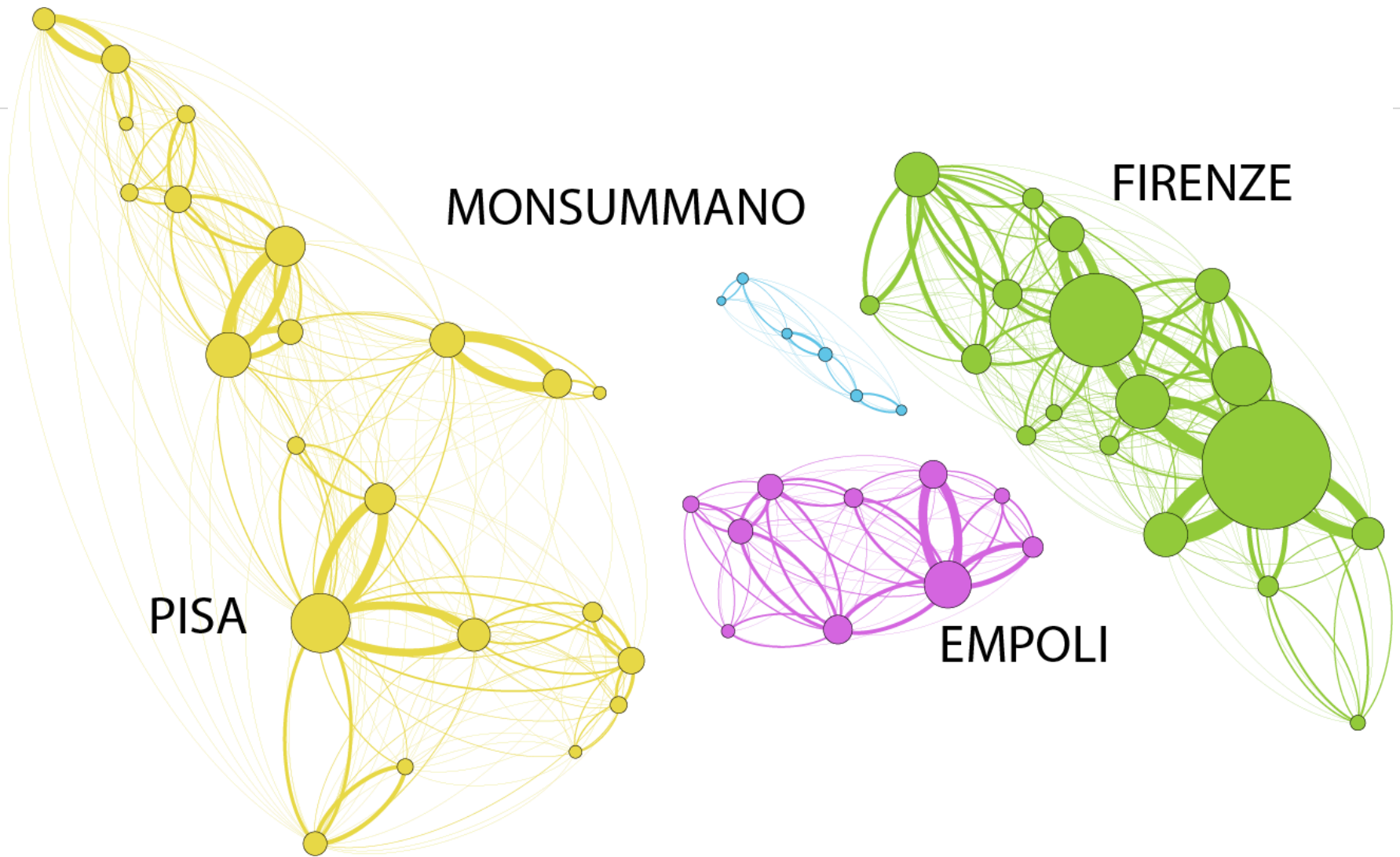


MONSUMMANO

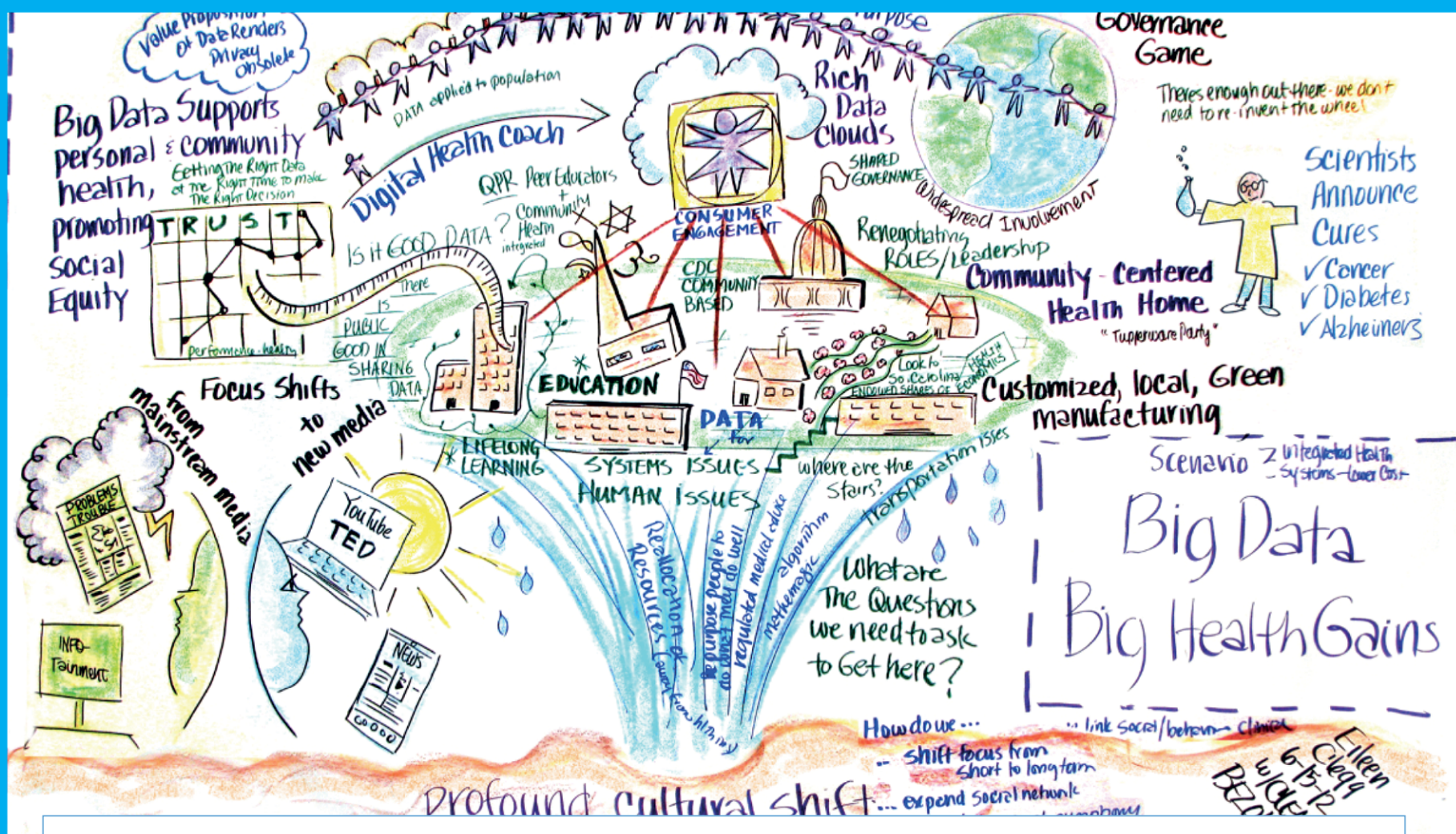
FIRENZE

PISA

EMPOLI



Exploratory: Big Data for Societal Debates



Polarization, controversy and topic trends on societal debates through social media

LONGER IN

NHS

IMMIGRATION

HOW BRITAIN TALKED BREXIT

ECONOMY

LAW

SOVEREIGN



3 Million Brexit Tweets Reveal Leave Voters Talked About Immigration More Than Anything Else

Groundbreaking analysis shows immigration, not sovereignty or the NHS, dominated the conversation – and making British judges responsible for British law was a key theme for Leave supporters.



James Ball

BuzzFeed Special
Correspondent



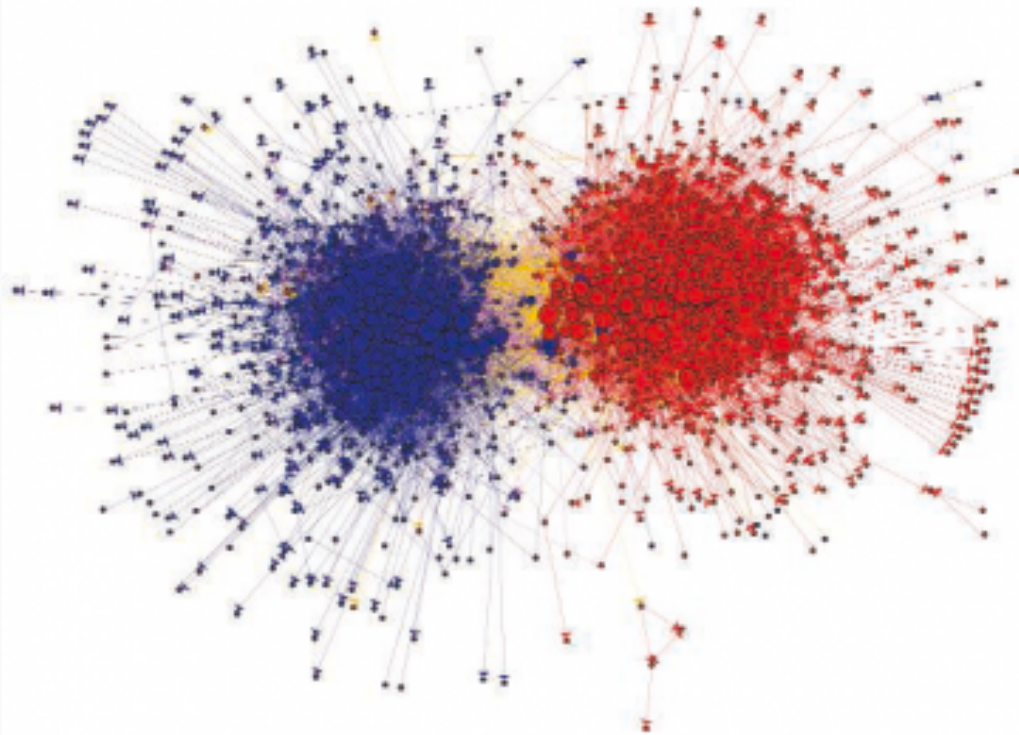
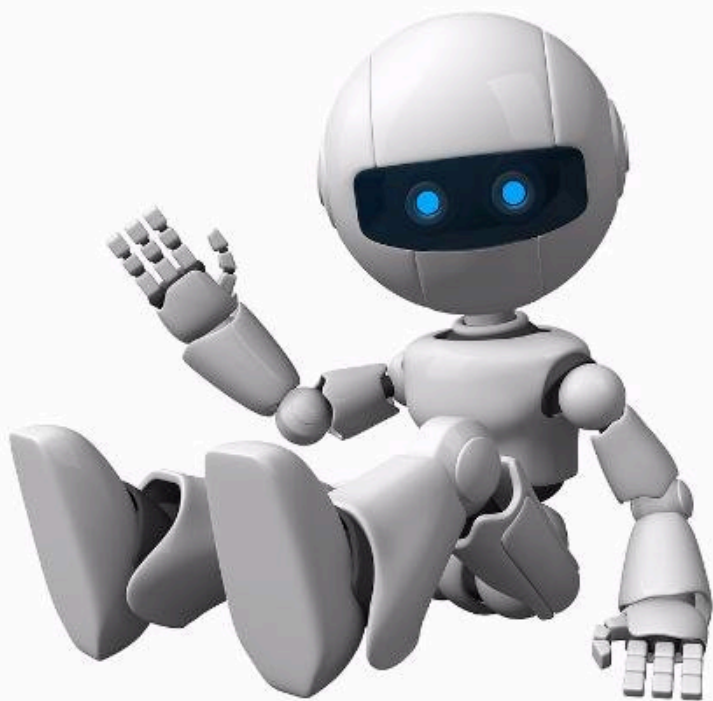
Chris Applegate

Editorial Developer, UK

posted on Dec. 9, 2016, at 2:03 p.m.

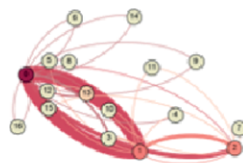
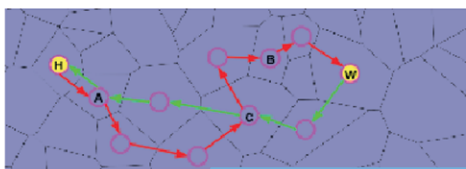
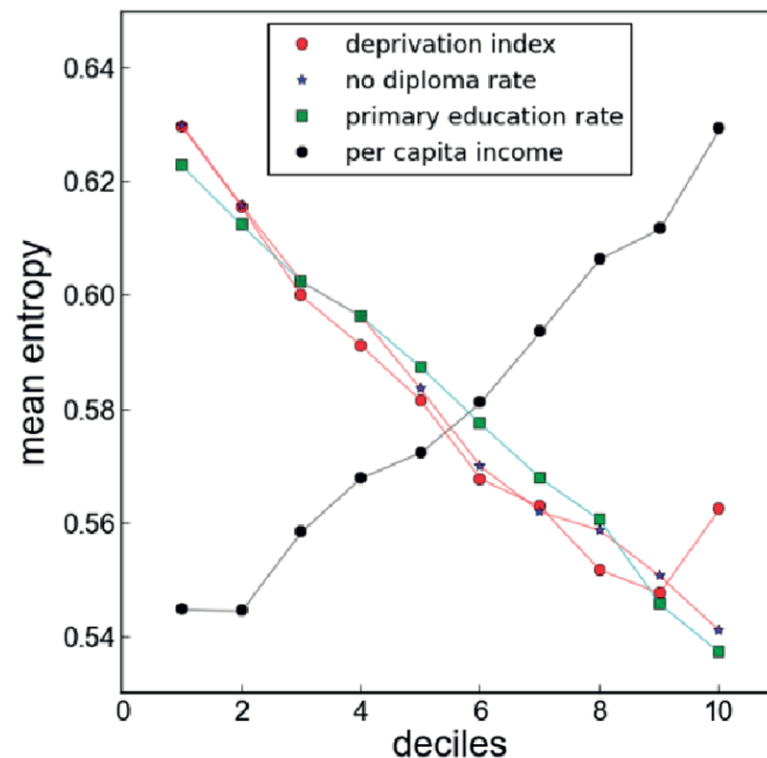
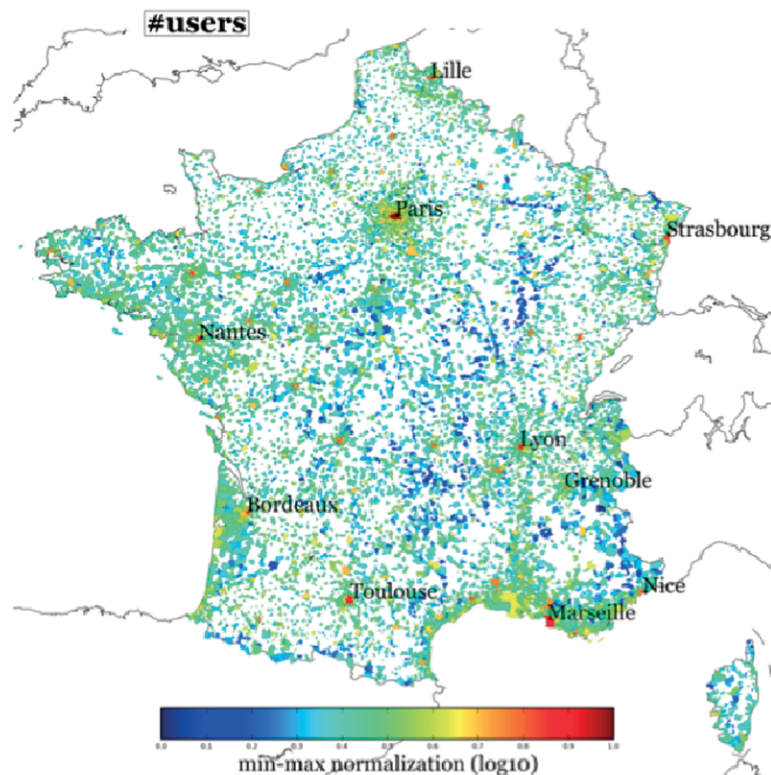


https://www.buzzfeed.com/jamesball/3-million-brexit-tweets-reveal-leave-voters-talked-about-imm?utm_term=.jmdQE9JNR#.fuOOrb145



Exploratory:

Big Data for Well Being and Economic Performance



$$d_i^{(n)} = \sum_{j=1}^{|V|} \frac{1}{k_j} M_{ij} p_j^{(n-1)} \forall i$$

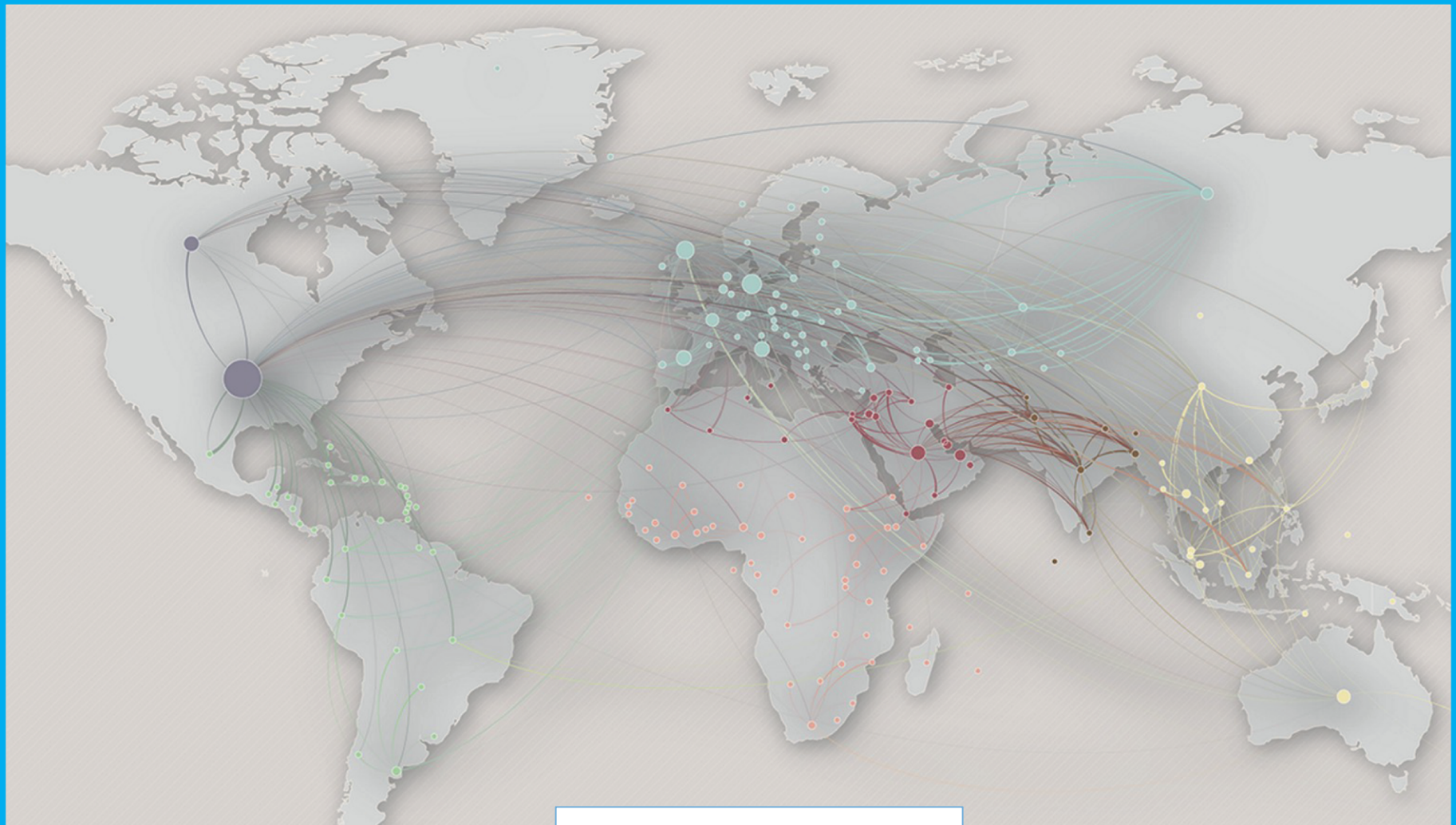
$$p_j^{(n)} = \sum_{i=1}^{|U|} \frac{1}{k_i} M_{ij} d_i^{(n-1)} \forall j$$

Deprivation Index (in France) predicted with Mobile Phone traces

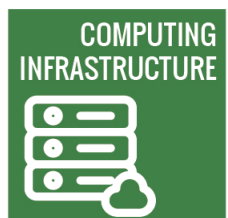


Next Exploratory:

Big Data for Migration Studies



Human Migration Flows



New economic growth: the role of science, technology, innovation and infrastructure



Policy recommendations

G7 Academies of Science urge governments to:

- i. expand investment and capabilities in science and pre-competitive technologies;

Maryse Lassonde
ROYAL SOCIETY OF CANADA

Sébastien Candel
ACADÉMIE DES SCIENCES

Jörg Hacker
LEOPOLDINA NATIONALE AKADEMIE
DER WISSENSCHAFTEN

Alberto Quadrio-Curzio
ACCADEMIA NAZIONALE DEI LINCEI

Takashi Onishi
SCIENCE COUNCIL OF JAPAN

Venki Ramakrishnan
ROYAL SOCIETY

Marcia McNutt
NATIONAL ACADEMY OF SCIENCES

G7 Academies' Joint Statements 2017

May

Attention should be given to emerging technologies in light of their potential to impact virtually all economic activities:

- Data Science, thanks to the ability to extract new knowledge and policy capability by the integrated algorithmic analysis of highly diverse data generated today at exponentially growing pace.

The Future of Jobs

Employment, Skills and Workforce Strategy for the Fourth Industrial Revolution

January 2016

New and Emerging Roles

Our research also explicitly asked respondents about new and emerging job categories and functions that they expect to become critically important to their industry by the year 2020, and where within their global operations they would expect to locate such roles.

Two job types stand out due to the frequency and consistency with which they were mentioned across practically all industries and geographies. The first are data analysts, as already frequently mentioned above, which companies expect will help them make sense and derive insights from the torrent of data generated by the technological disruptions referenced above. The second



http://www3.weforum.org/docs/WEF_Future_of_Jobs.pdf



ENROLLMENT FOR
FOREIGN STUDENTS

Magistrale in Data Science and
Business Informatics

- » Insegnamenti
- » Docenti
- » Orario
- » Calendario AA 2016/2017
- » Calendario AA 2017/2018
- » Calendario appelli
- » Piani di studio
- » Servizio di tutorato
- » Progetto formativo
- » Lauree
- » Internazionale
- » Valutazione



Presentazione

Il Corso di Laurea Magistrale in **Data Science and Business Informatics** (fino all'A.A. 2016/17, Business Informatics) è progettato, a partire dal 2002, per preparare laureati magistrali in grado di padroneggiare sia le tecnologie informatiche che di

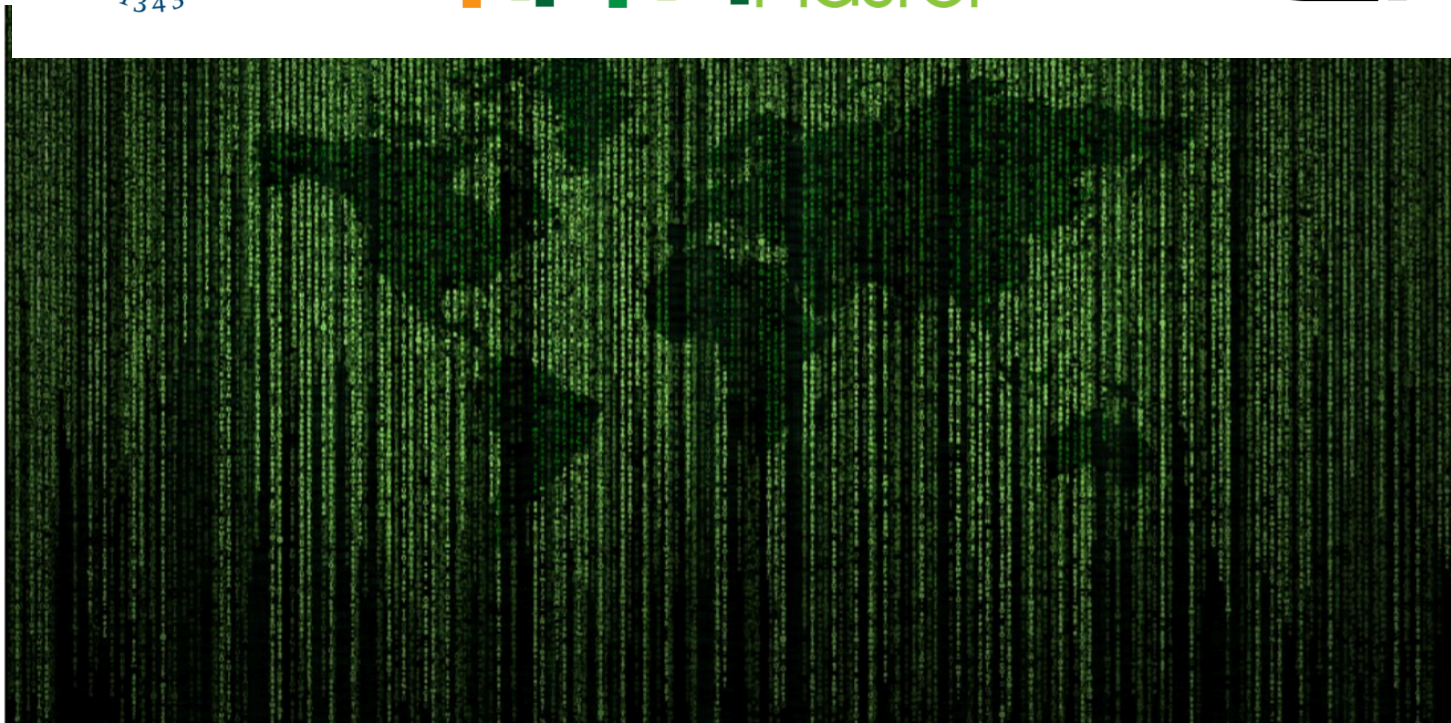
NOTIZIE DIPARTIMENTO

- » 12 PhD positions in Computer Science at University of Pisa
- » Dottorato di ricerca

<http://masterbigdata.it>



Big Data Analytics & Social Mining
SoBigData
Master



Febbraio
01
2017

Master Big Data 2017
Inaugurazione



SCUOLA
NORMALE
SUPERIORE



Sant'Anna
Scuola Universitaria Superiore Pisa

Dal 2014



Ph.D. in Data Science

Start: academic year 2017-2018

<http://phd.sns.it/data-science/>



SCUOLA
ALTI STUDI
LUCCA





www.sobigdata.eu

**H2020 excellent science
research infrastructure**



SCUOLA
NORMALE
SUPERIORE



SCUOLA
ALTI STUDI
LUCCA



Consiglio
Nazionale delle
Ricerche

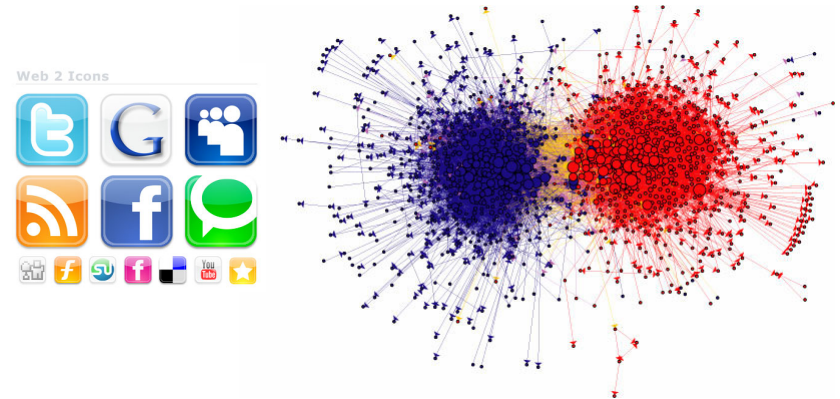


Big data proxies of social life

Shopping patterns & lifestyle



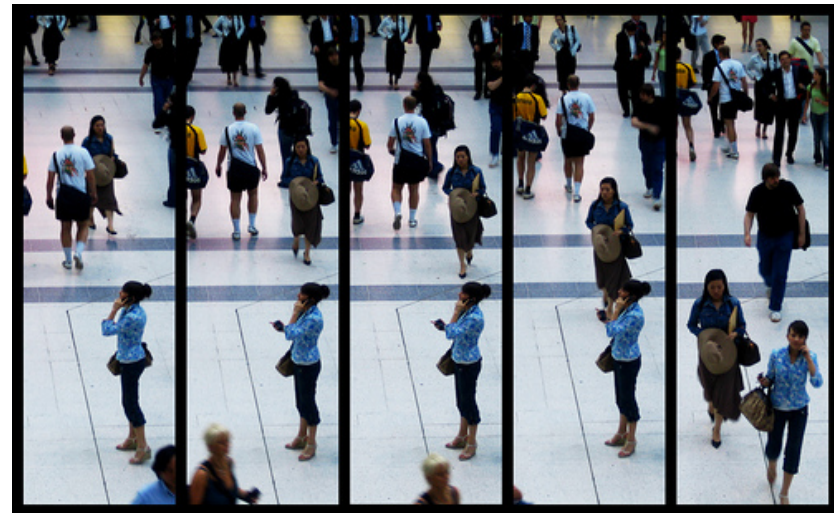
RELATIONSHIPS & SOCIAL TIES

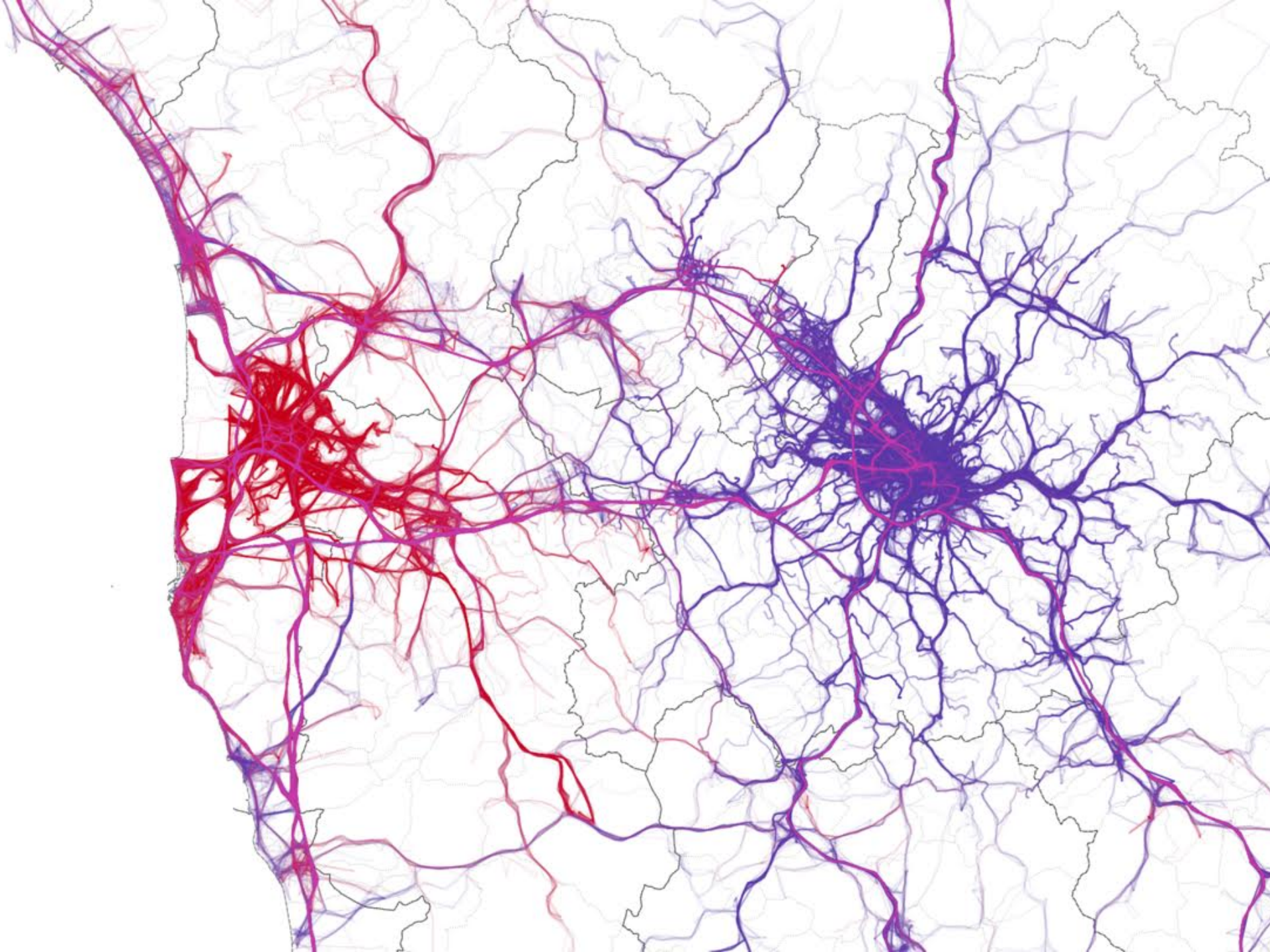


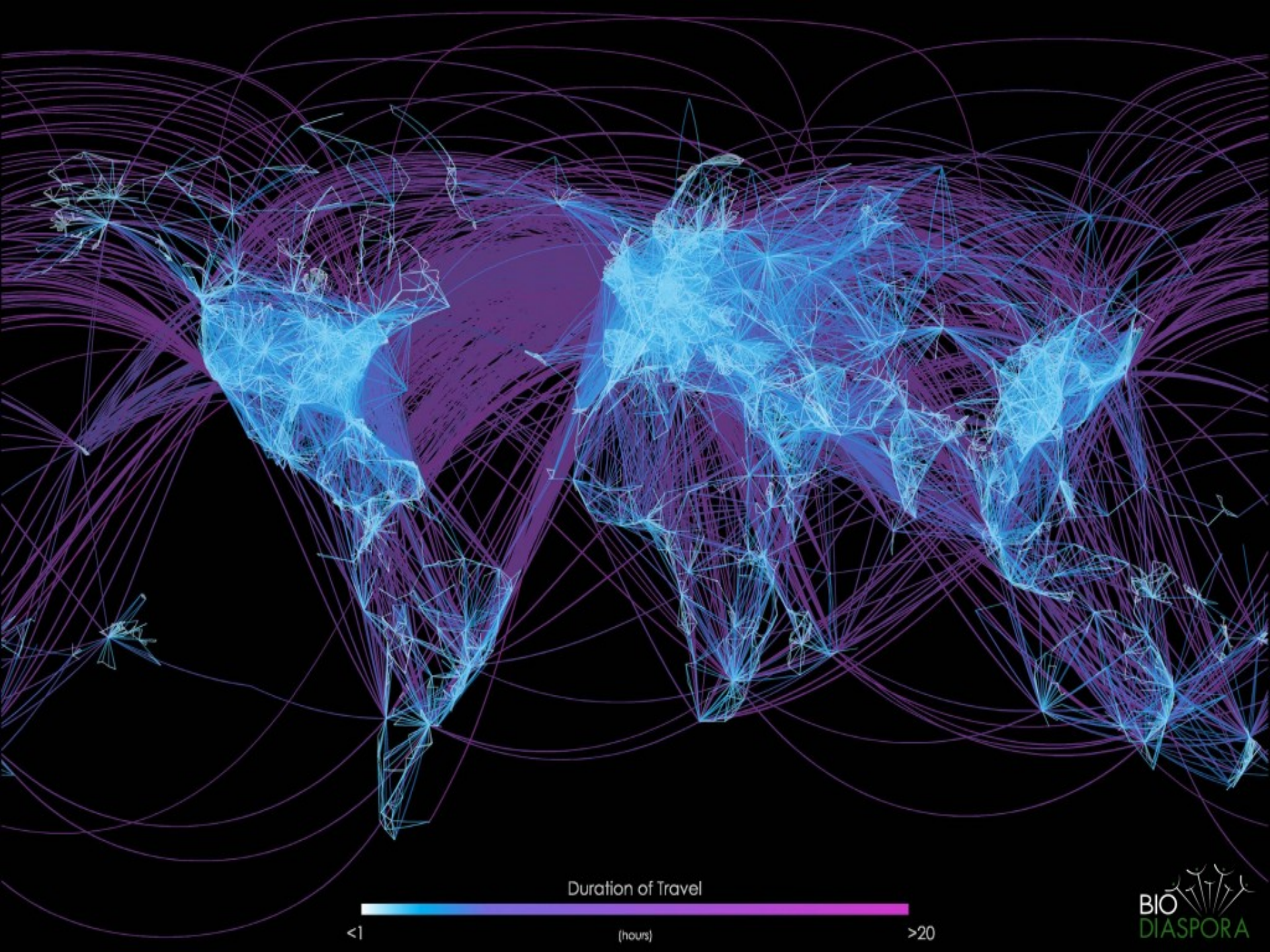
DESIRES, OPINIONS, SENTIMENTS



MOVEMENTS







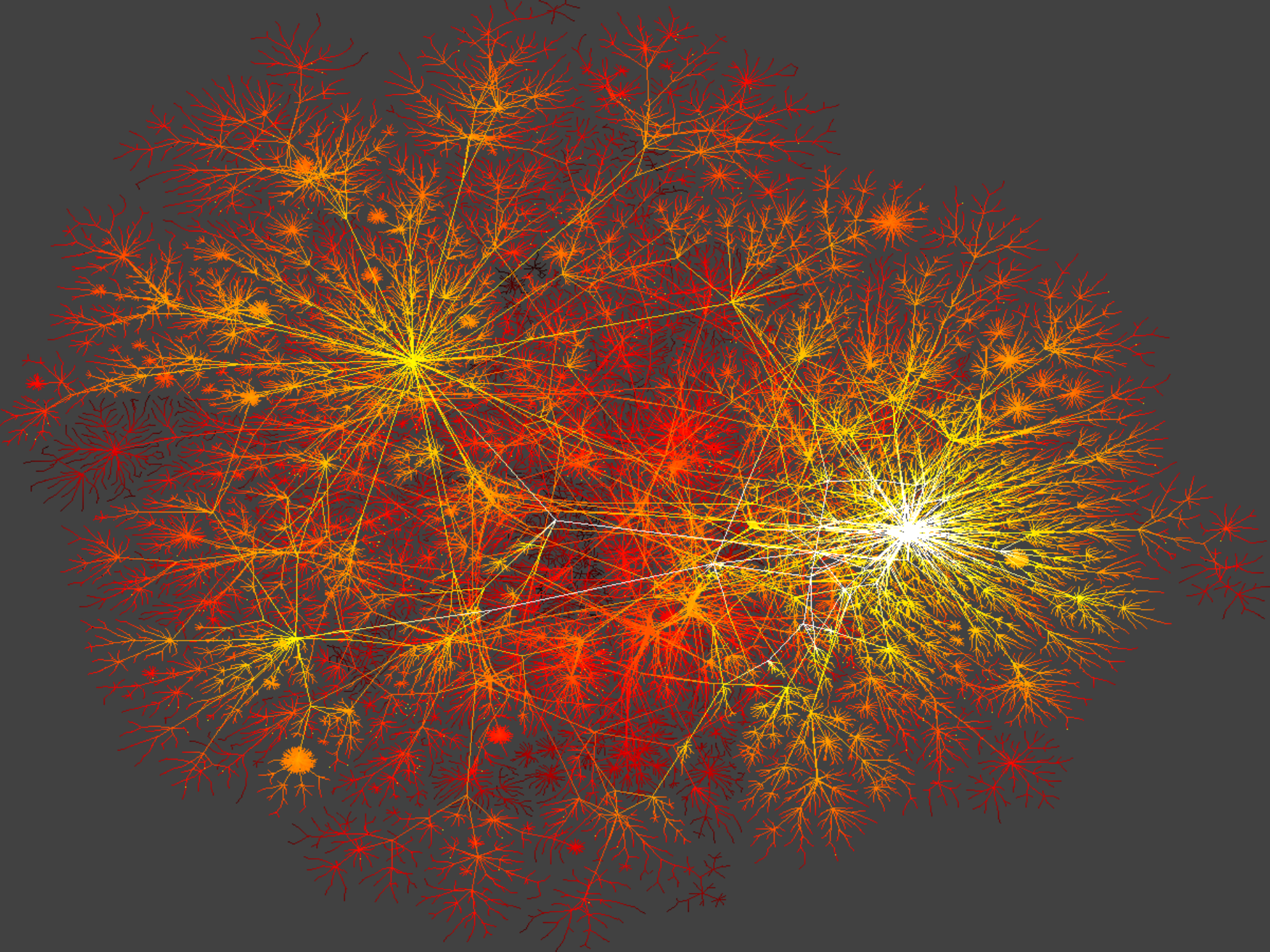
Duration of Travel

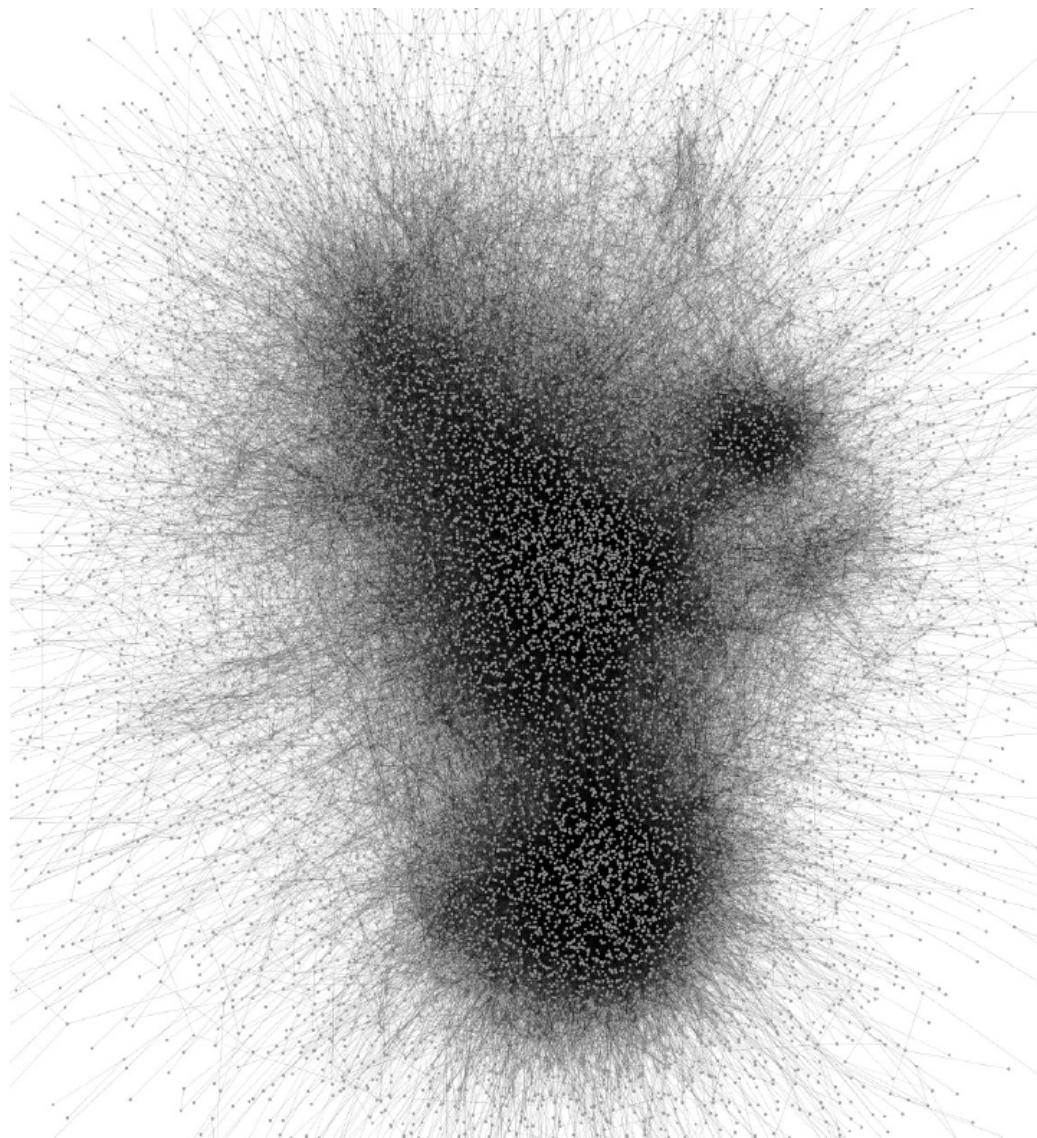
<1

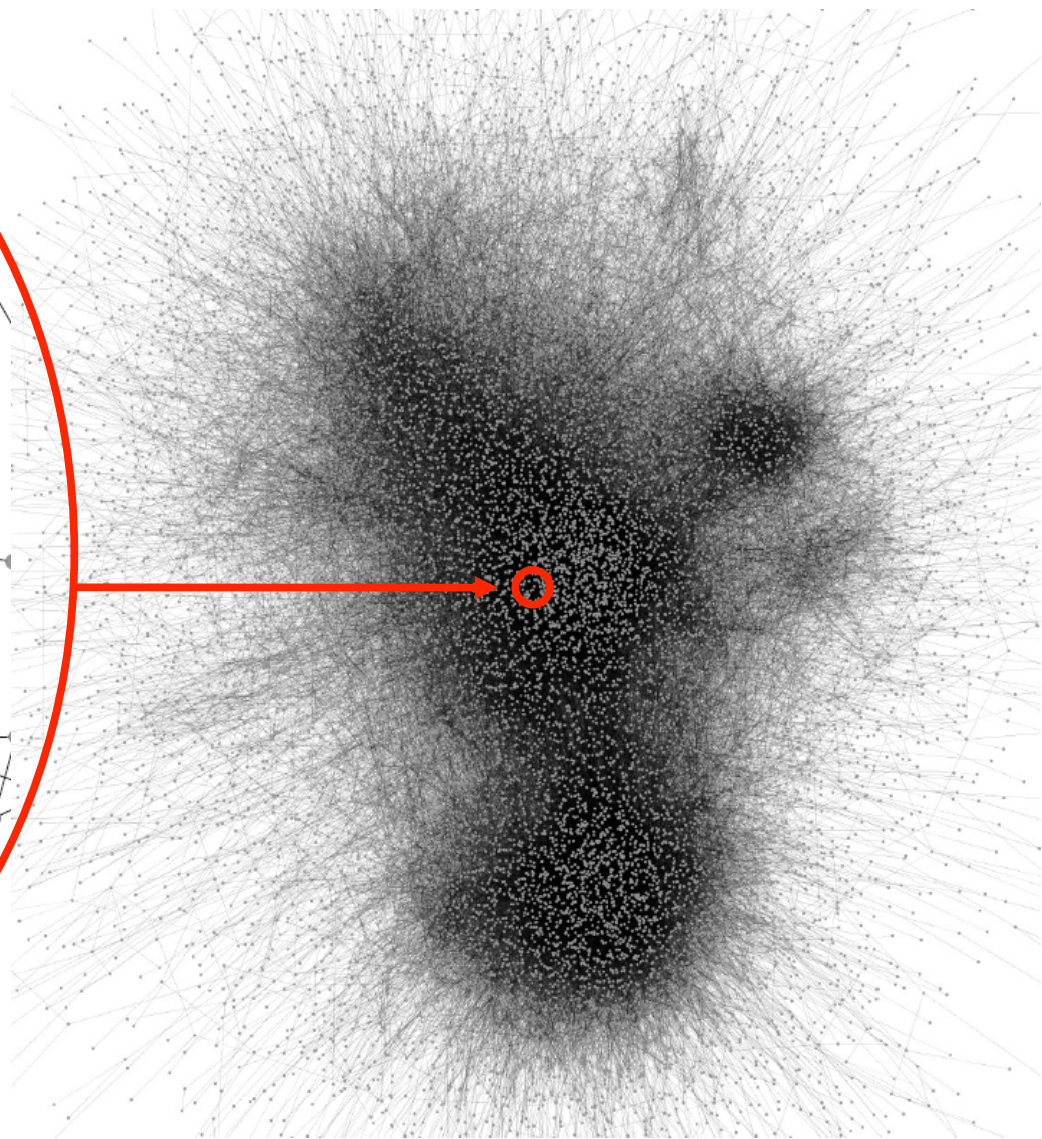
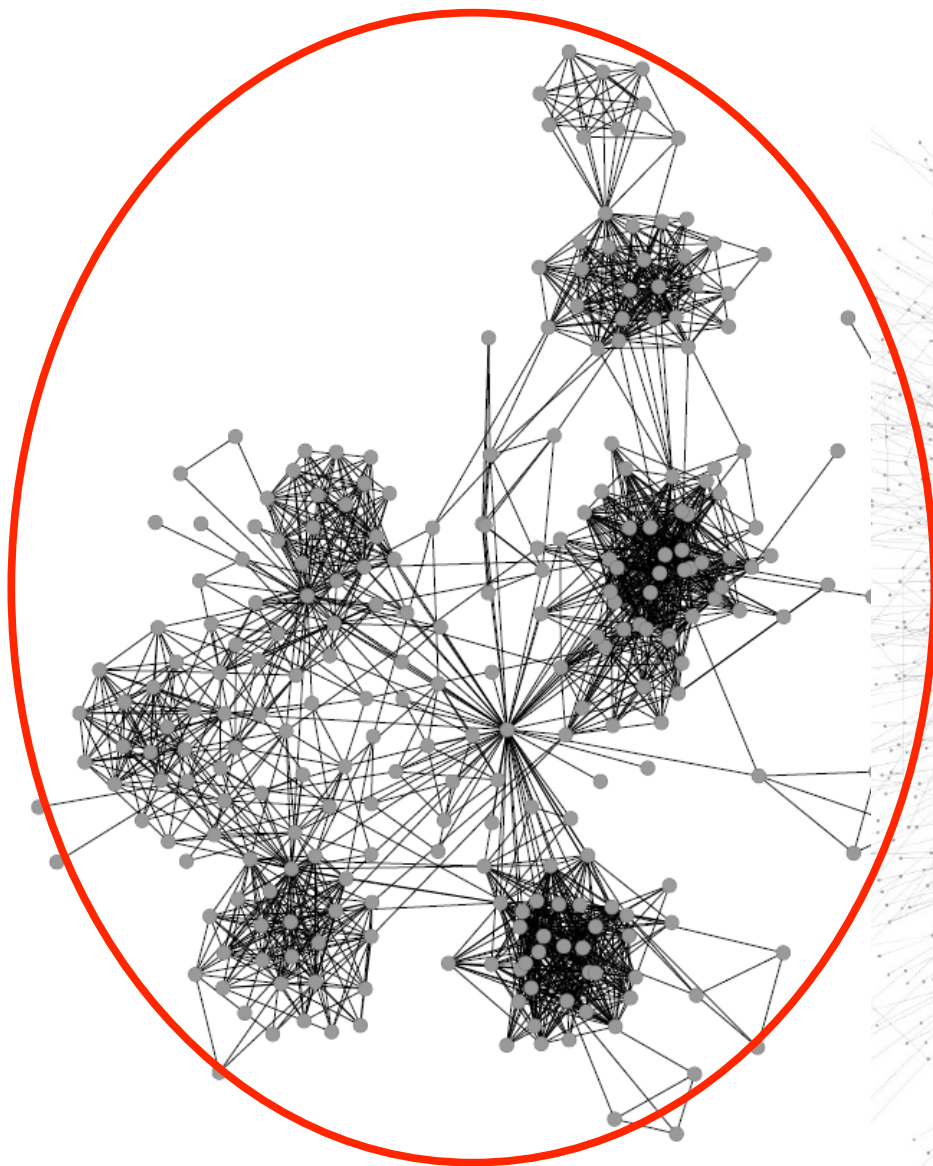
(hours)

>20











Complex (Social) Networks

- Big graph data and social, information, biological and technological networks
- The architecture of complexity and how real networks differ from random networks:
 - node degree and long tails,
 - social distance and small worlds,
 - clustering and triadic closure.
- Comparing real networks and random graphs.
- The main models of network science: small world and preferential attachment.



Complex (Social) Networks

- Strong and weak ties, community structure and long-range bridges.
- Robustness of networks to failures and attacks.
- Cascades and spreading. Network models for diffusion and epidemics. The strength of weak ties for the diffusion of information. The strength of strong ties for the diffusion of innovation.



Complex (Social) Networks

- Textbooks
 - Albert-Laszlo Barabasi. *Network Science* (2016)
 - <http://barabasi.com/book/network-science>
 - Easley, Kleinberg: *Networks, Crowds, and Markets* (2010)
 - <http://www.cs.cornell.edu/home/kleinber/networks-book/>
- Network Analytics Software:
 - Cytoscape: <http://www.cytoscape.org/>
 - Gephi: <http://gephi.github.io/>
- Network dynamics simulation :
 - NetLogo: <https://ccl.northwestern.edu/netlogo/>
- Network Data Repository
 - <http://networkrepository.com/>

Wiki of the course

- <http://didawiki.di.unipi.it/doku.php/wma/acm-athens-july2017>
- Special thanks to
 - Fosca Giannotti, ISTI-CNR Pisa
 - Albert-Laszlo Barabasi, Northeastern Univ. Boston
 - Giulio Rossetti, University of Pisa
 - Jure Leskovec, Stanford Univ.



The architecture of complexity

Lecture 1



Complex

[adj., v. kuh m-pleks, kom-pleks; n. kom-pleks]

—adjective

1.

composed of many interconnected parts; compound; composite: a complex highway system.

2.

characterized by a very complicated or involved arrangement of parts, units, etc.: complex machinery.

3.

so complicated or intricate as to be hard to understand or deal with: a complex problem.

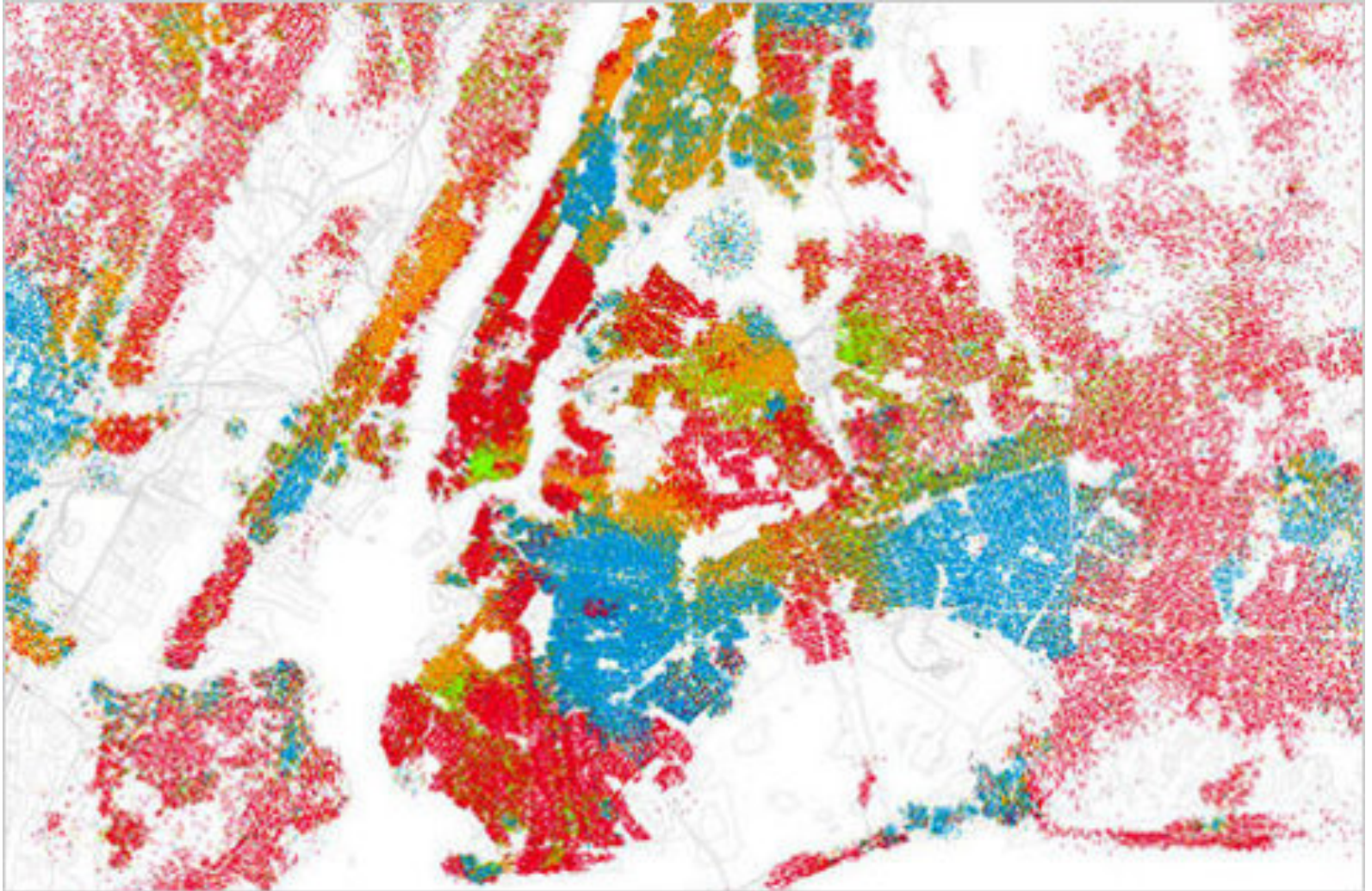
Source: Dictionary.com

Complexity, a **scientific theory** which asserts that some systems display behavioral phenomena that are completely inexplicable by any conventional analysis of the systems' constituent parts. These phenomena, commonly referred to as **emergent behaviour**, seem to occur in many complex systems involving living organisms, such as a stock market or the human brain.

Source: John L. Casti, Encyclopædia Britannica

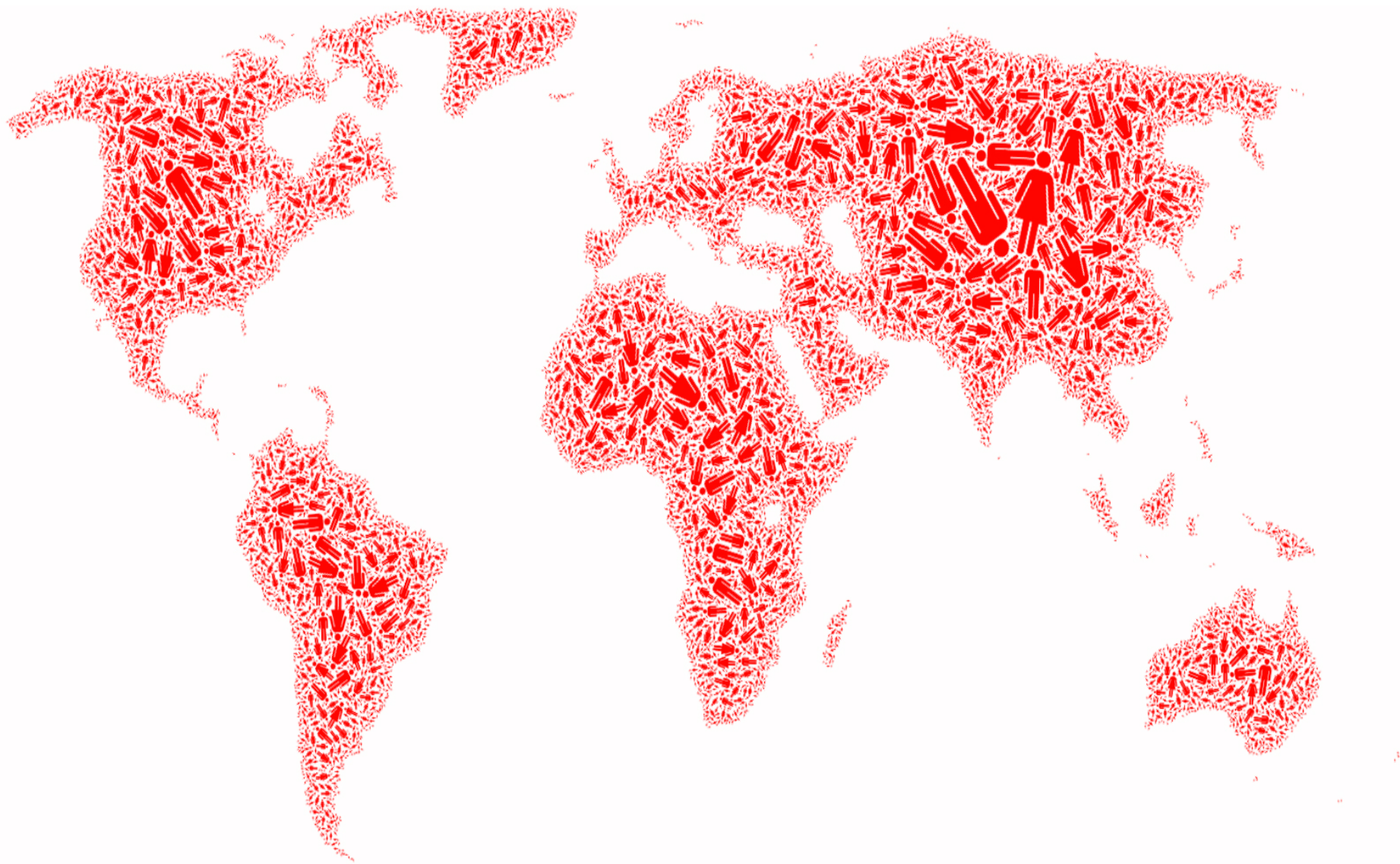
Complexity

Emergent behavior: segregation

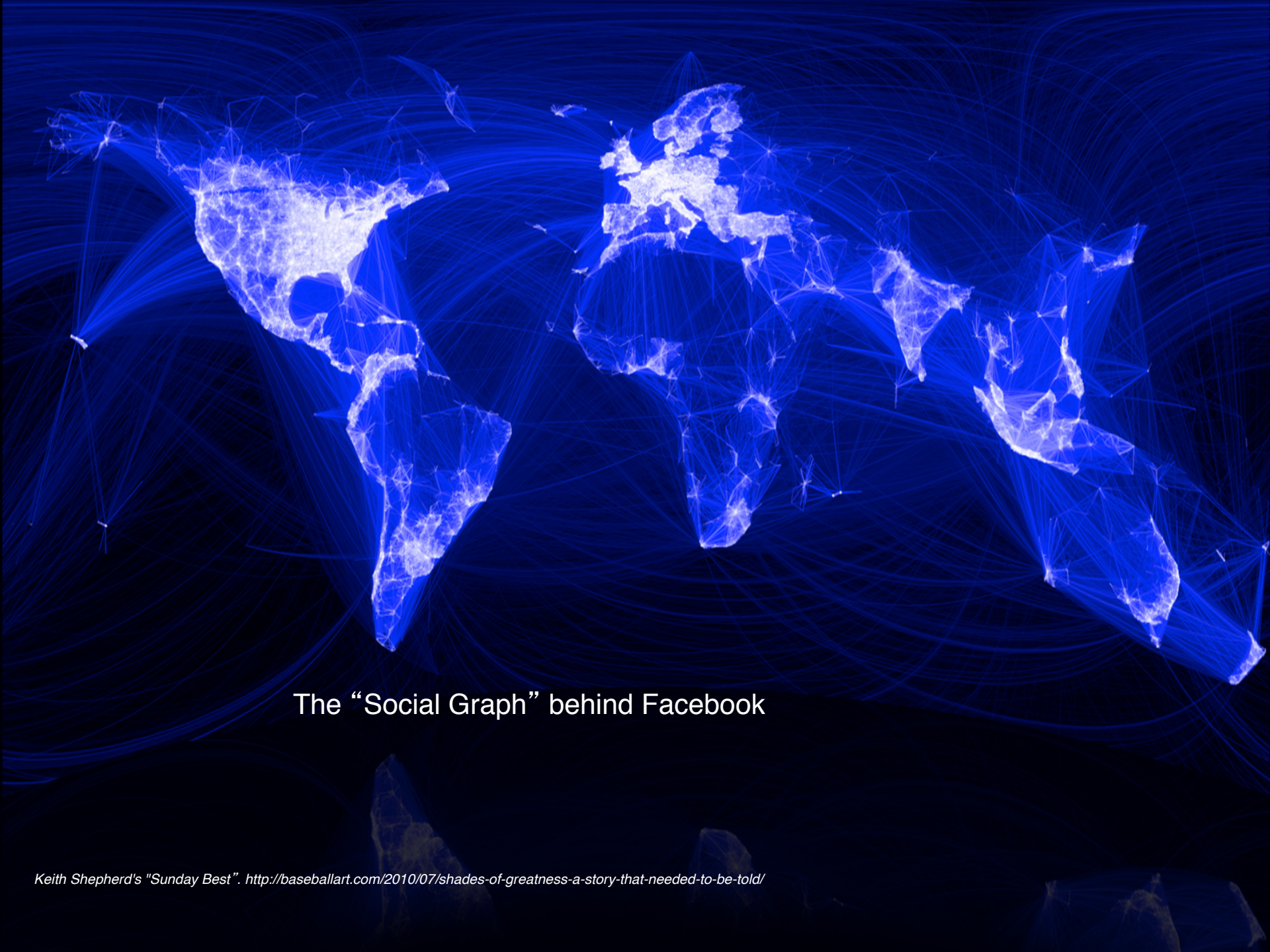


Behind each complex system there is a **network**, that defines the interactions between the components.

Social, informational,
technological, biological networks

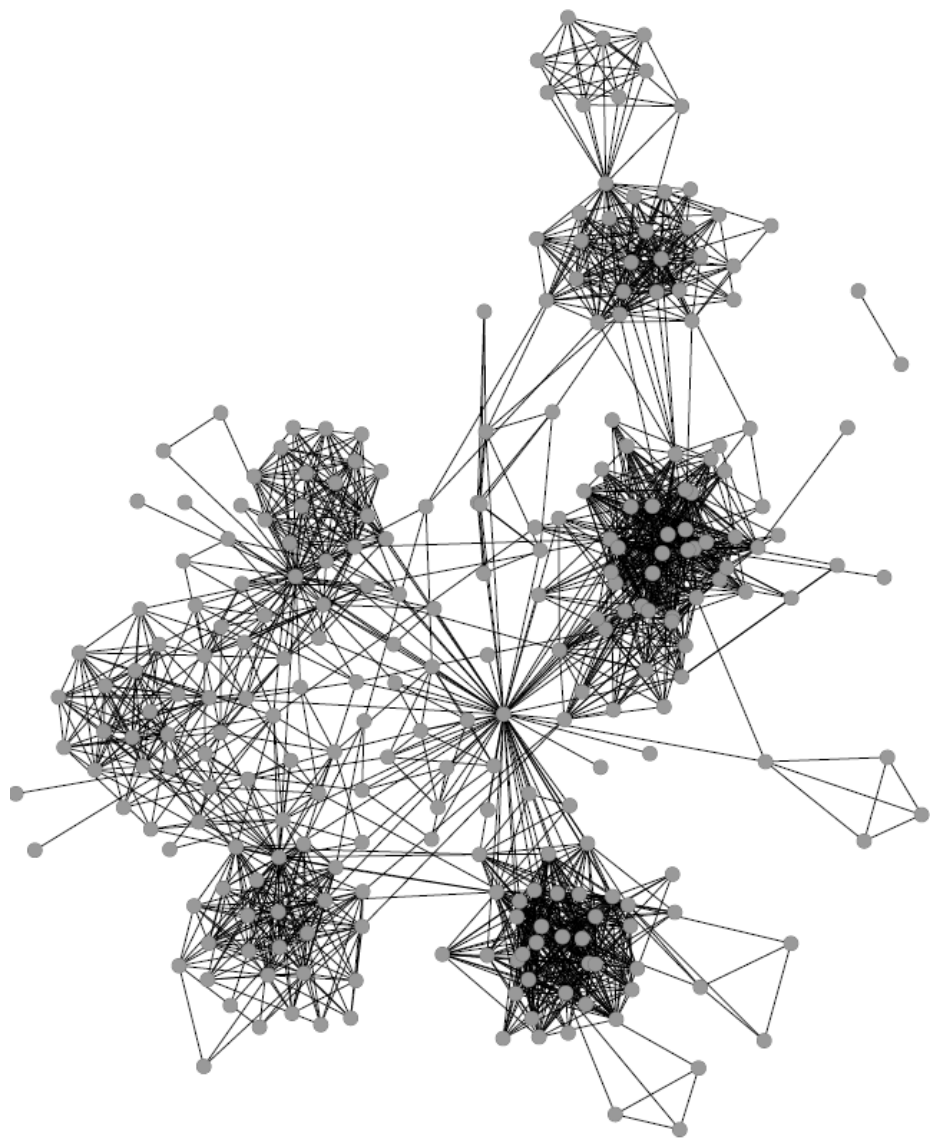
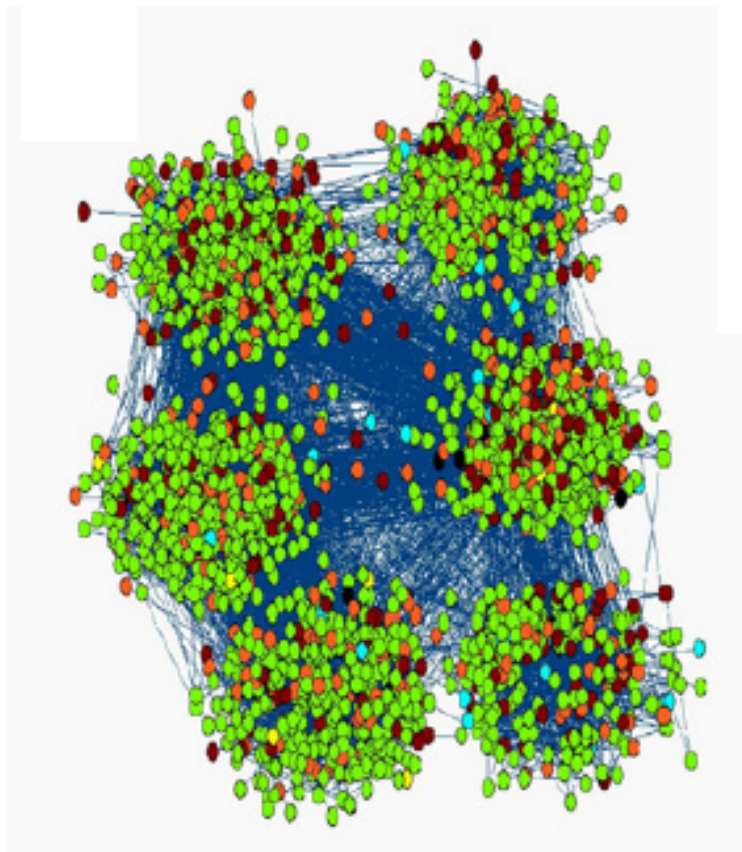


The "Day of 7 Billion" has been in October 2011

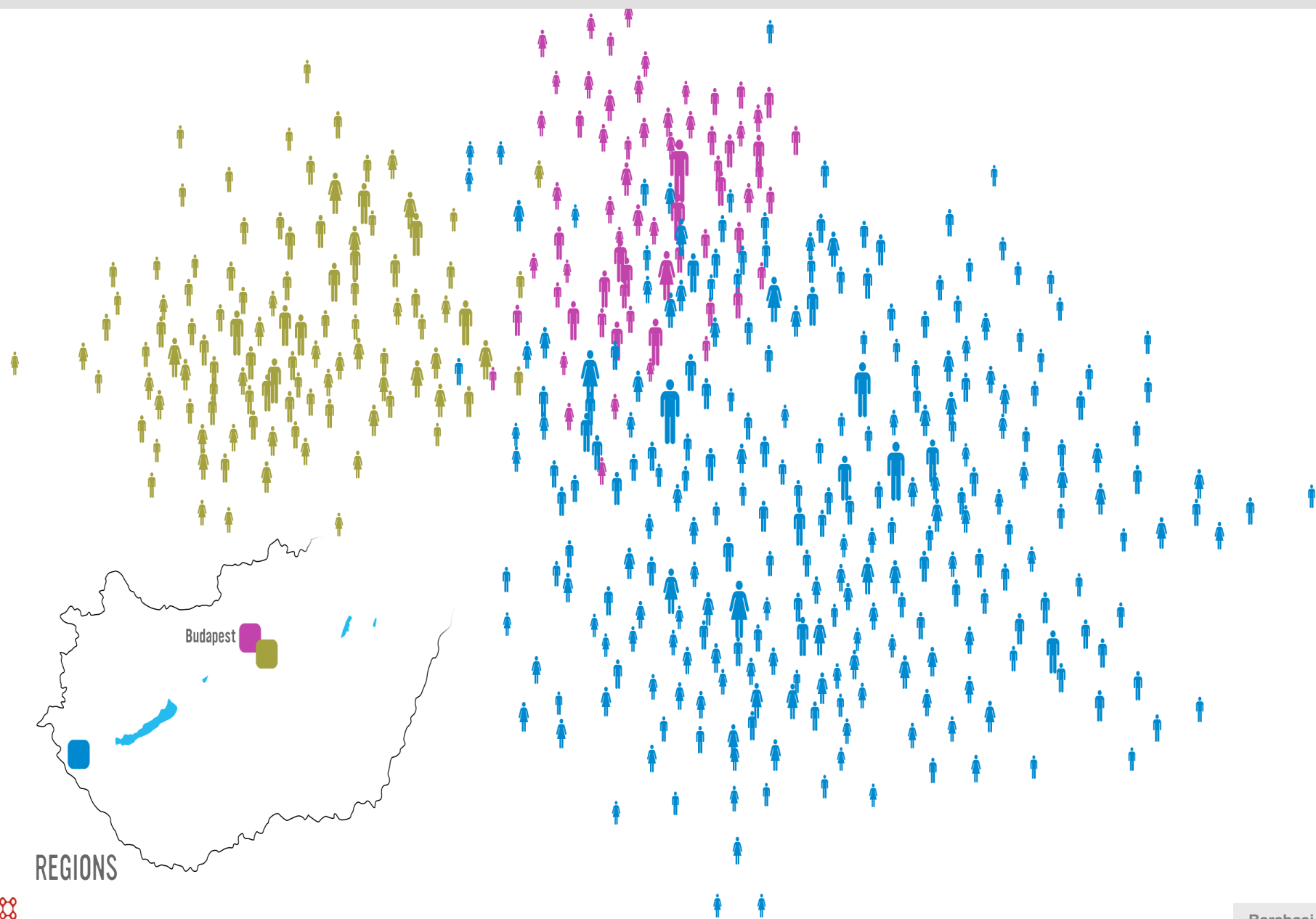


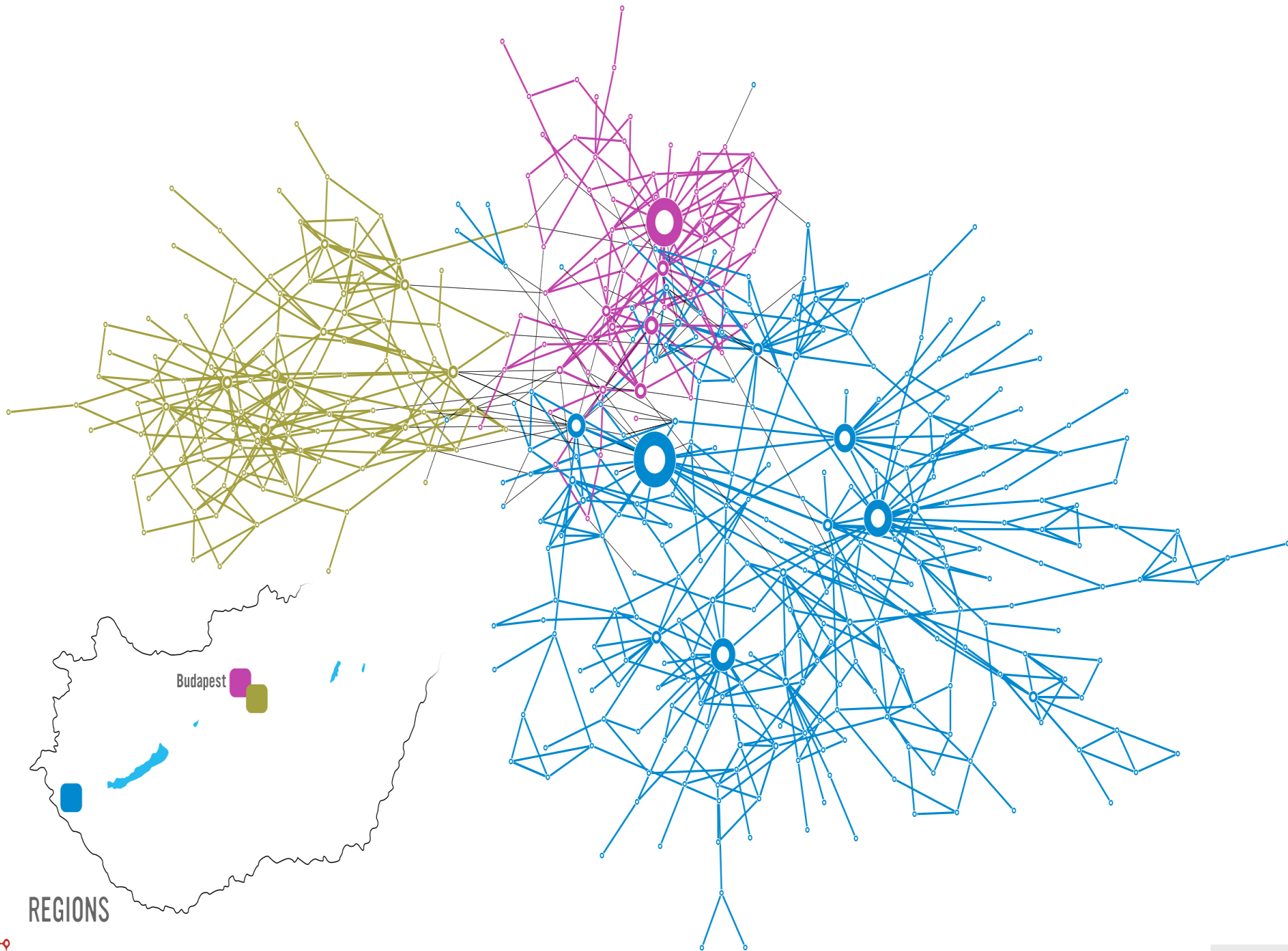
The “Social Graph” behind Facebook

Keith Shepherd's "Sunday Best". <http://baseballart.com/2010/07/shades-of-greatness-a-story-that-needed-to-be-told/>

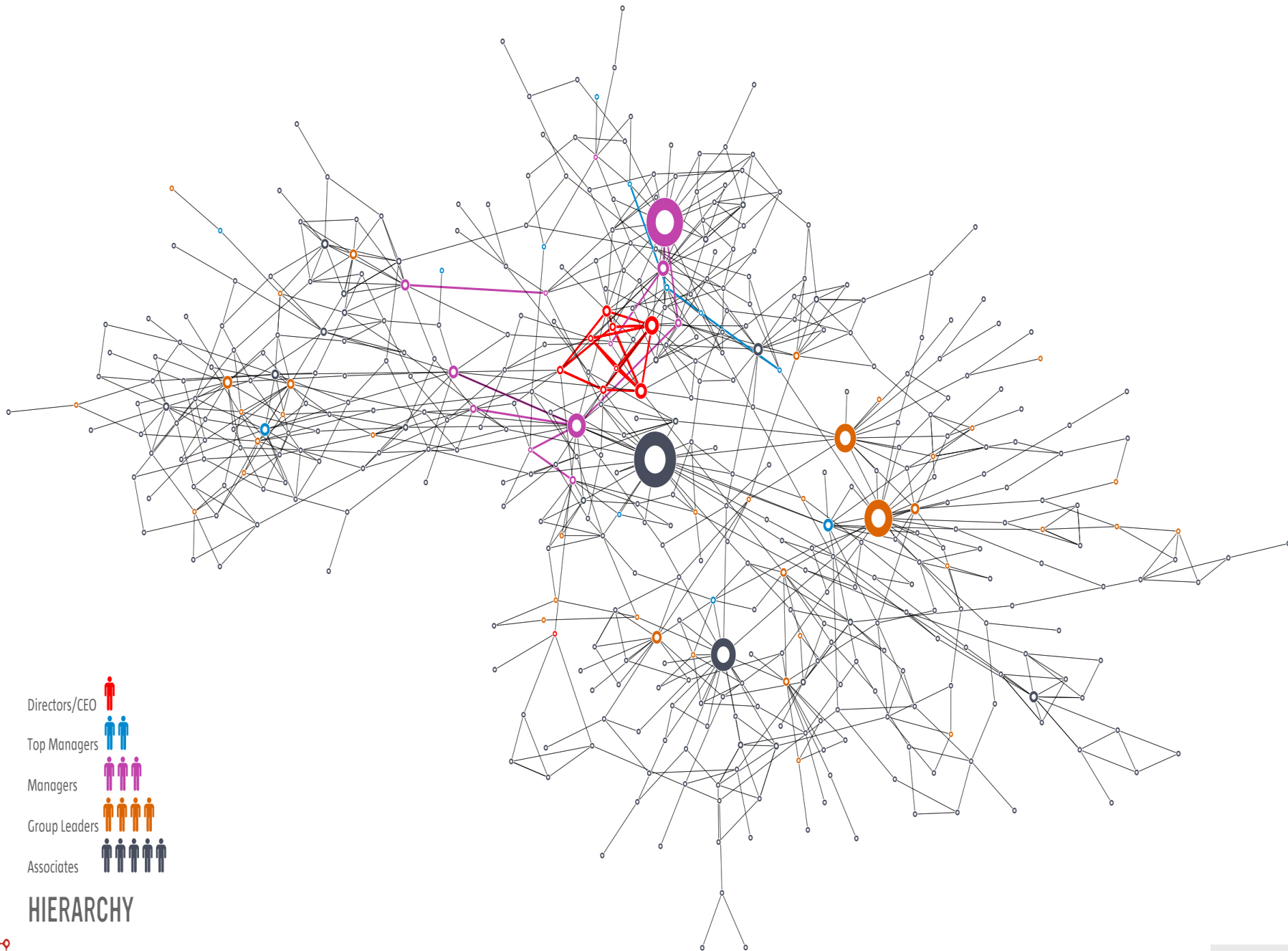


Mapping Organizations



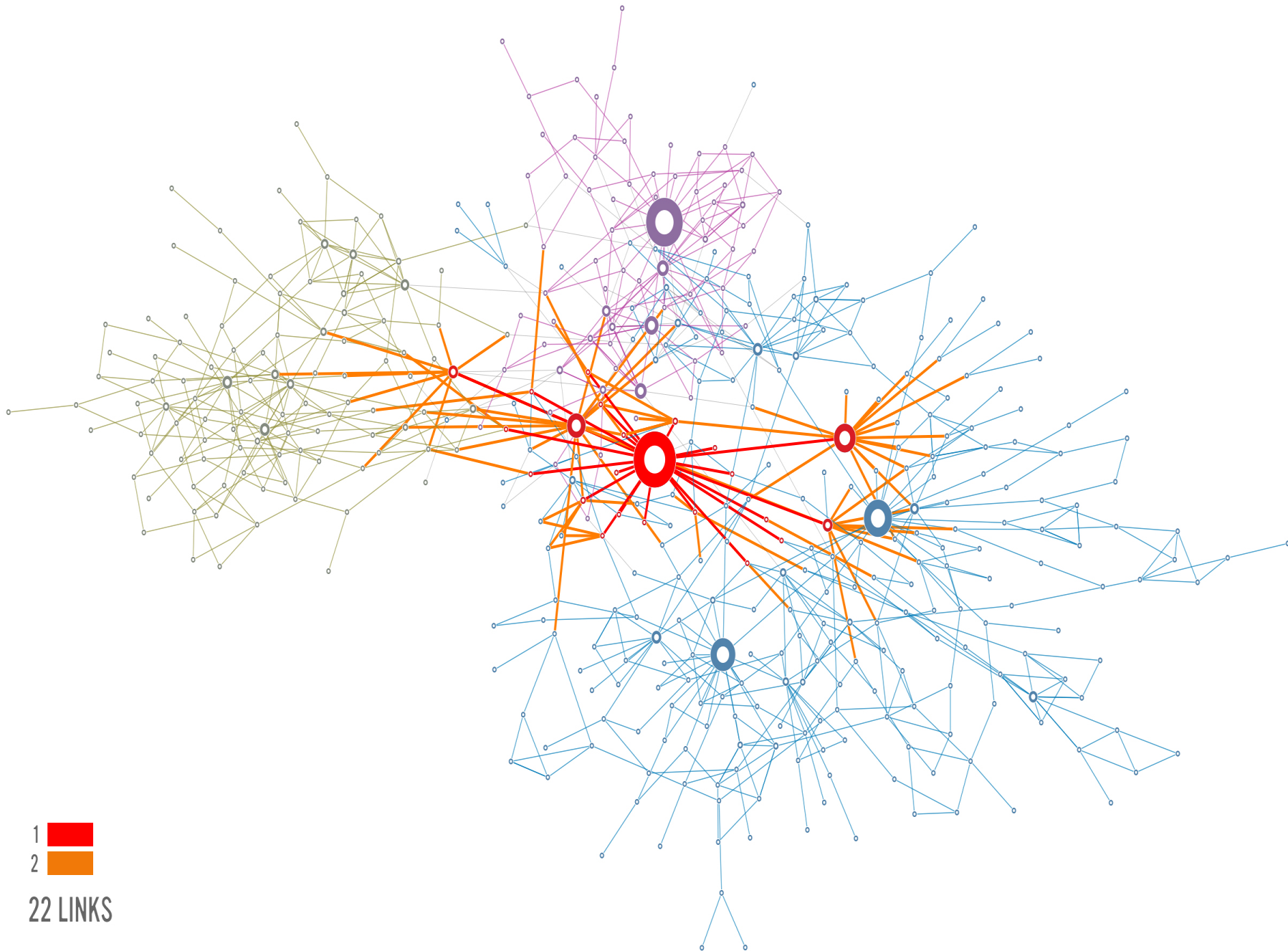


REGIONS



- Directors/CEO 
- Top Managers 
- Managers 
- Group Leaders 
- Associates 

HIERARCHY



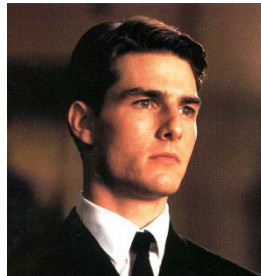
1
2

22 LINKS

COLLABORATION NETWORKS: ACTOR NETWORK

Nodes: actors

Links: cast jointly



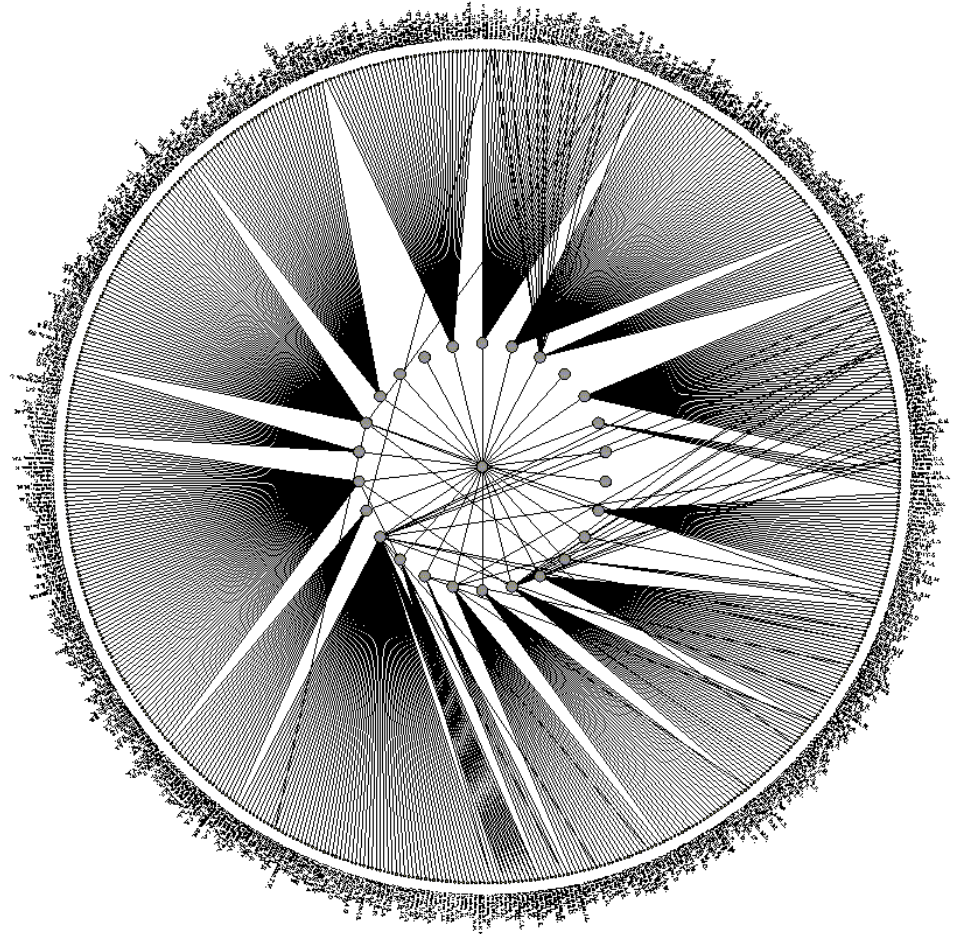
Days of Thunder (1990)
Far and Away (1992)
Eyes Wide Shut (1999)



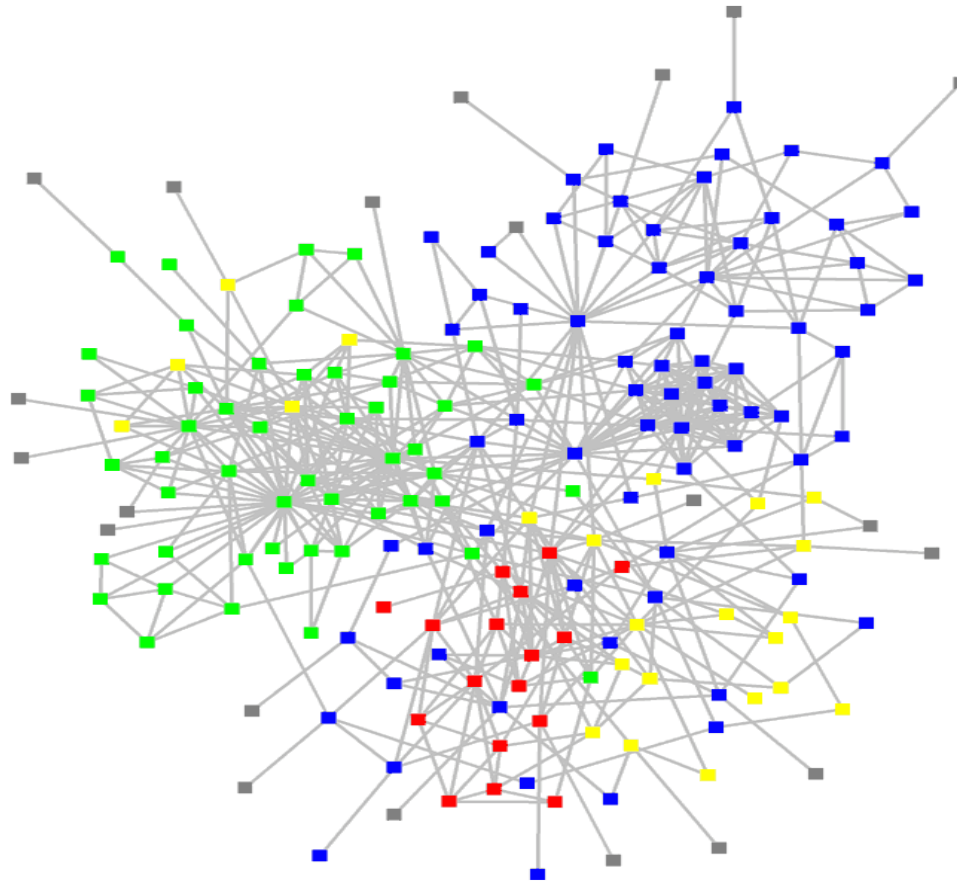
$N = 212,250$ actors $\langle k \rangle = 28.78$

Nodes: scientist (authors)

Links: write paper together



STRUCTURE OF AN ORGANIZATION



www.orgnet.com

BUSINESS TIES IN US BIOTECH-INDUSTRY

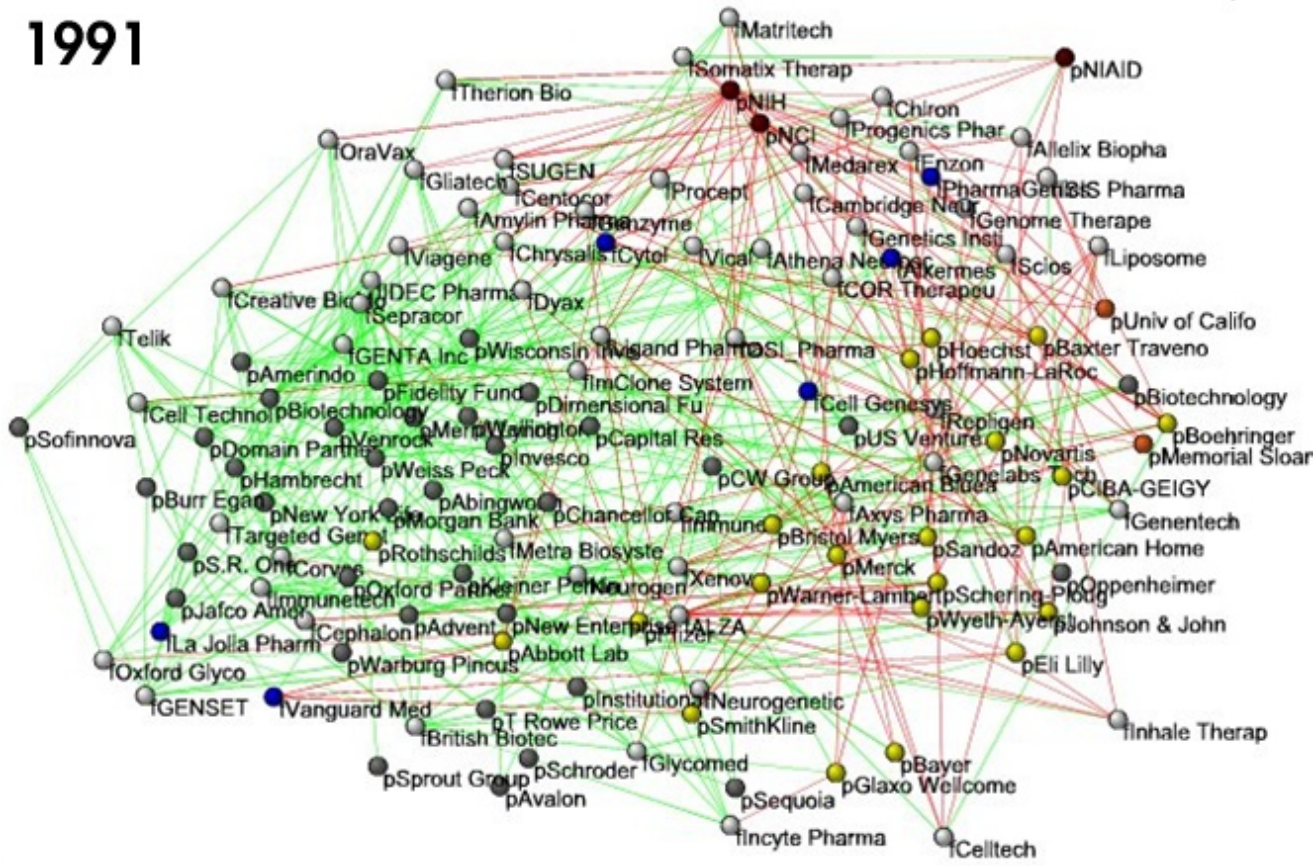
1991

Nodes:

- Companies
- Investment
- Pharma
- Research Labs
- Public
- Biotechnology

Links:

- Collaborations
- Financial
- R&D

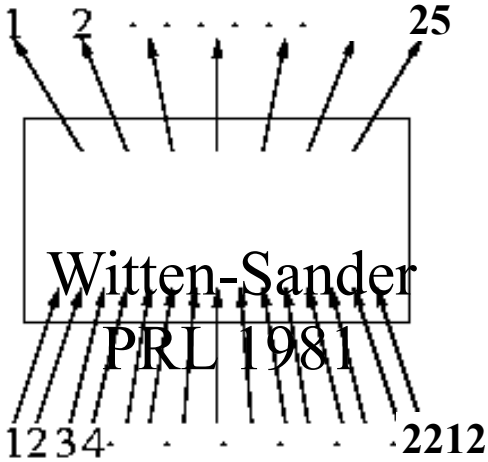


<http://ecclectic.ss.uci.edu/~drwhite/Movie>

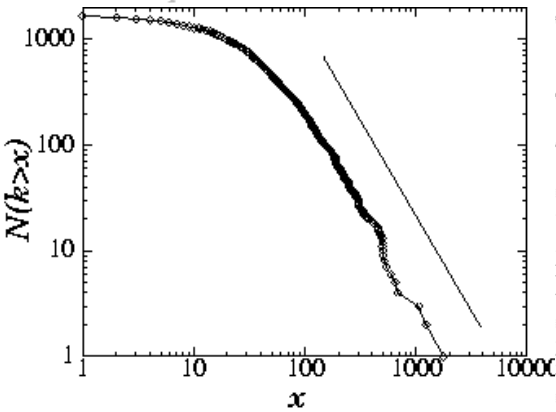
Information networks: the Web and Science Citation Indexes

1,000 Most Cited Physicists
Out of over 500,000 E
(see <http://www.esl.nu>)

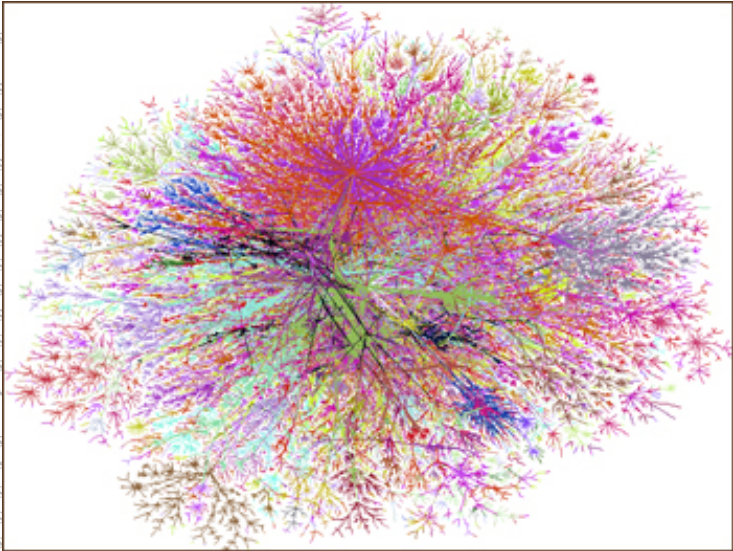
Author name	Institution	Country	Field
Witten	Princeton (U)	USA, NJ	High
Gossard	UCSB (U)	USA, CA	Sem
Cava	Princeton (U)	USA, NJ	Sup
Ballogg	Princeton (U)	USA, NJ	Sup
Ploog	Max-Planck (NL)	Germany	Sem
	Nuclear Cent.	Switzerland	Astr
	State (U)	USA, FL	Solid
	Frank (NL)	Germany	Sem
	anck (U)	USA, TX	High
	(U)	USA, CA	Poly
	on (U)	USA, NJ	Solid
	Western (U)	USA, IL	S
	Univ. (U)	Switzerland	S
	bs (I)	USA, NJ	S
	I/NL)	USA, CA	C
	L (U)	USA, IL	S
	d (U)	USA, CA	S
	n Univ. (U)	USA, TX	S
)	Switzerland	S
	BL (U/NL)	USA, CA	S
	n Univ. (U)	USA, TX	S



1736 PRL papers (1988)

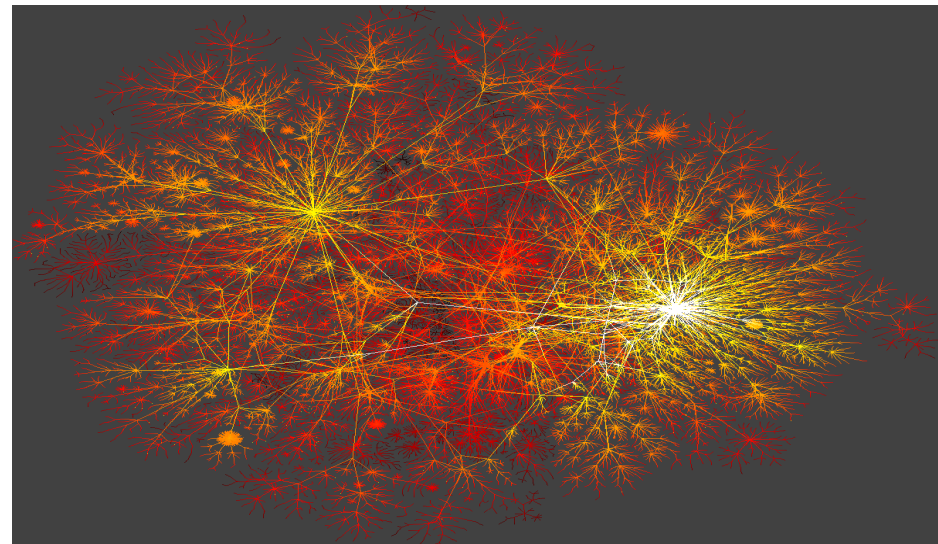
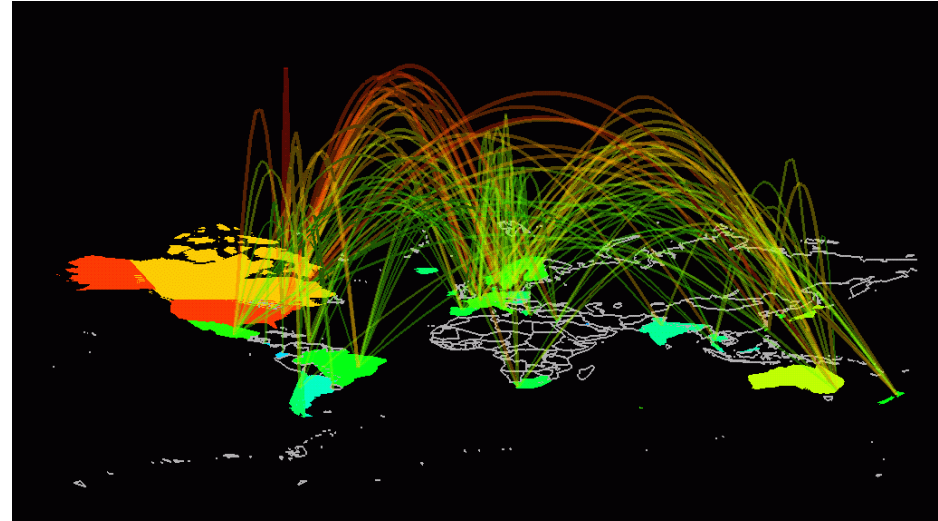
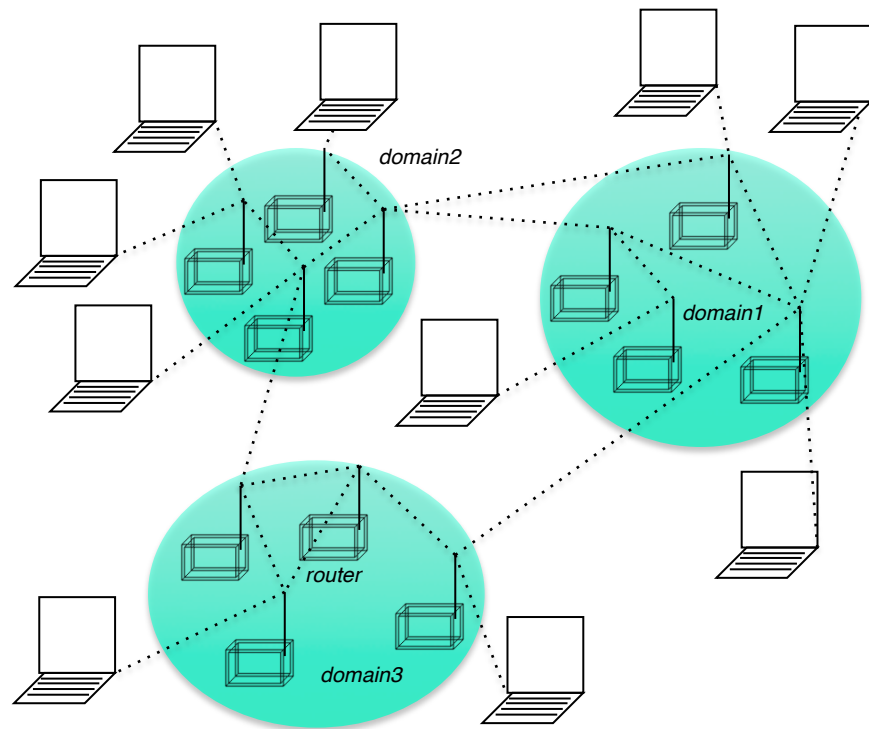


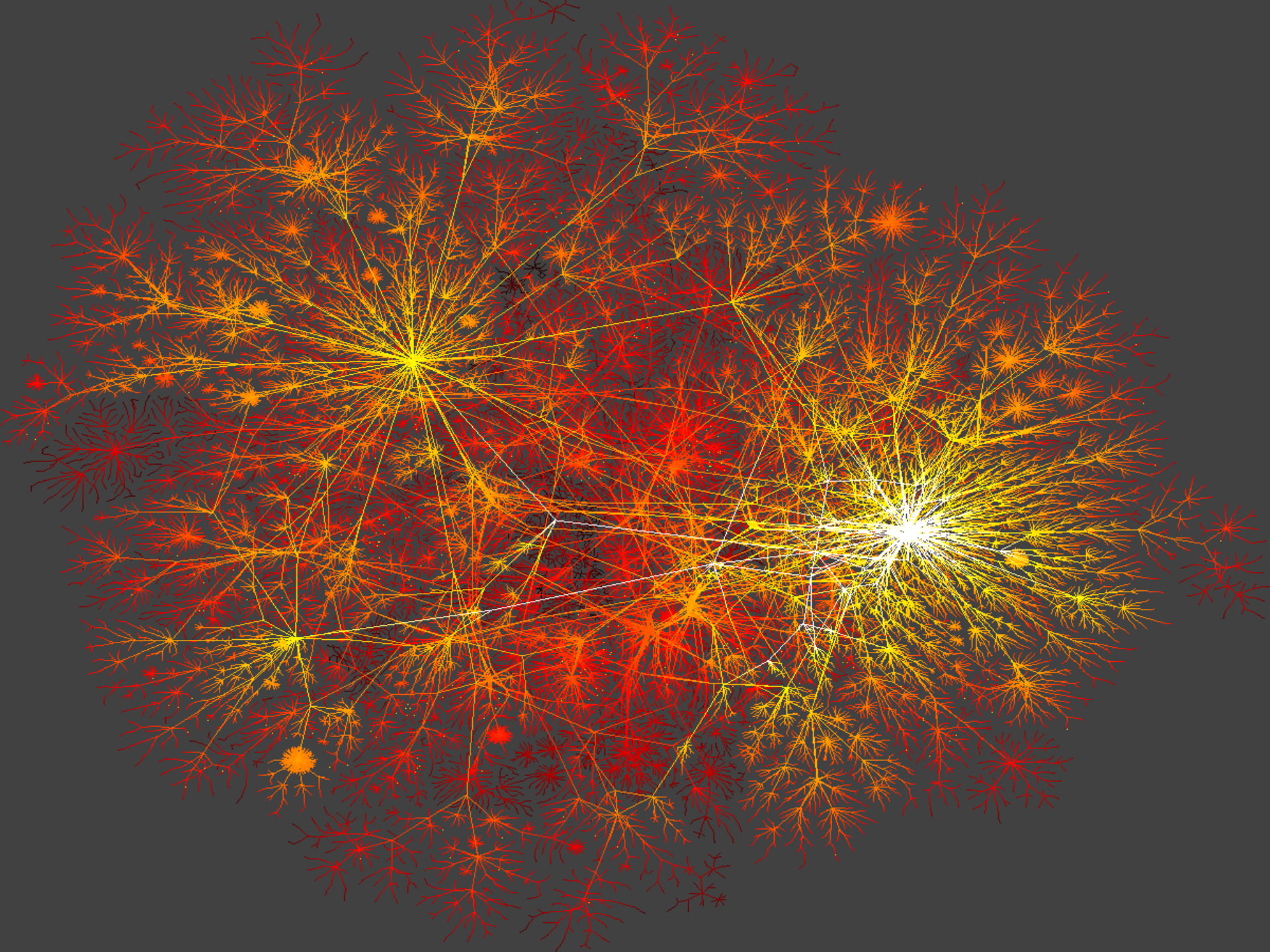
Waszczak	JV	AT&T (I)	USA, NJ	S
Shirane	G	Brookhaven (U)	USA, NY	S
Wiegmann	W	Brookhaven (U)	USA, NY	S
Vandover	RG	Gen Labs (I)	USA, NJ	M
Uchida*	S			
Hor	PE	Brookhaven (U)	USA, TX	S
Murphy	DW			A
Birgeneau	RJ	MIT (U)	USA, MA	S
Jorgensen	JD	Argonne (NL)	USA, IL	S
Hinks	DG	Argonne (NL)	USA, IL	S



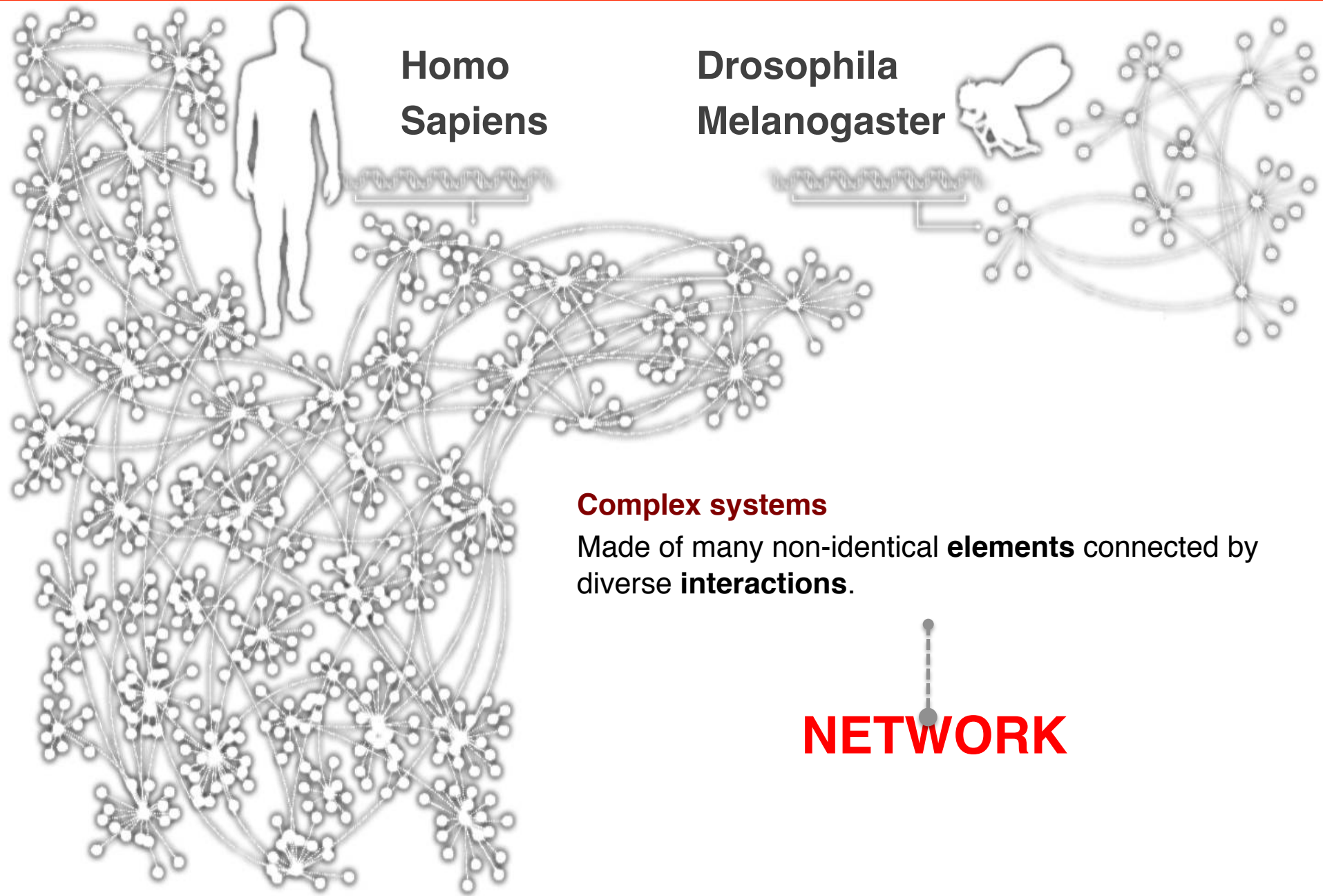
* citation total may be skewed because of multiple authors with the same name

INTERNET



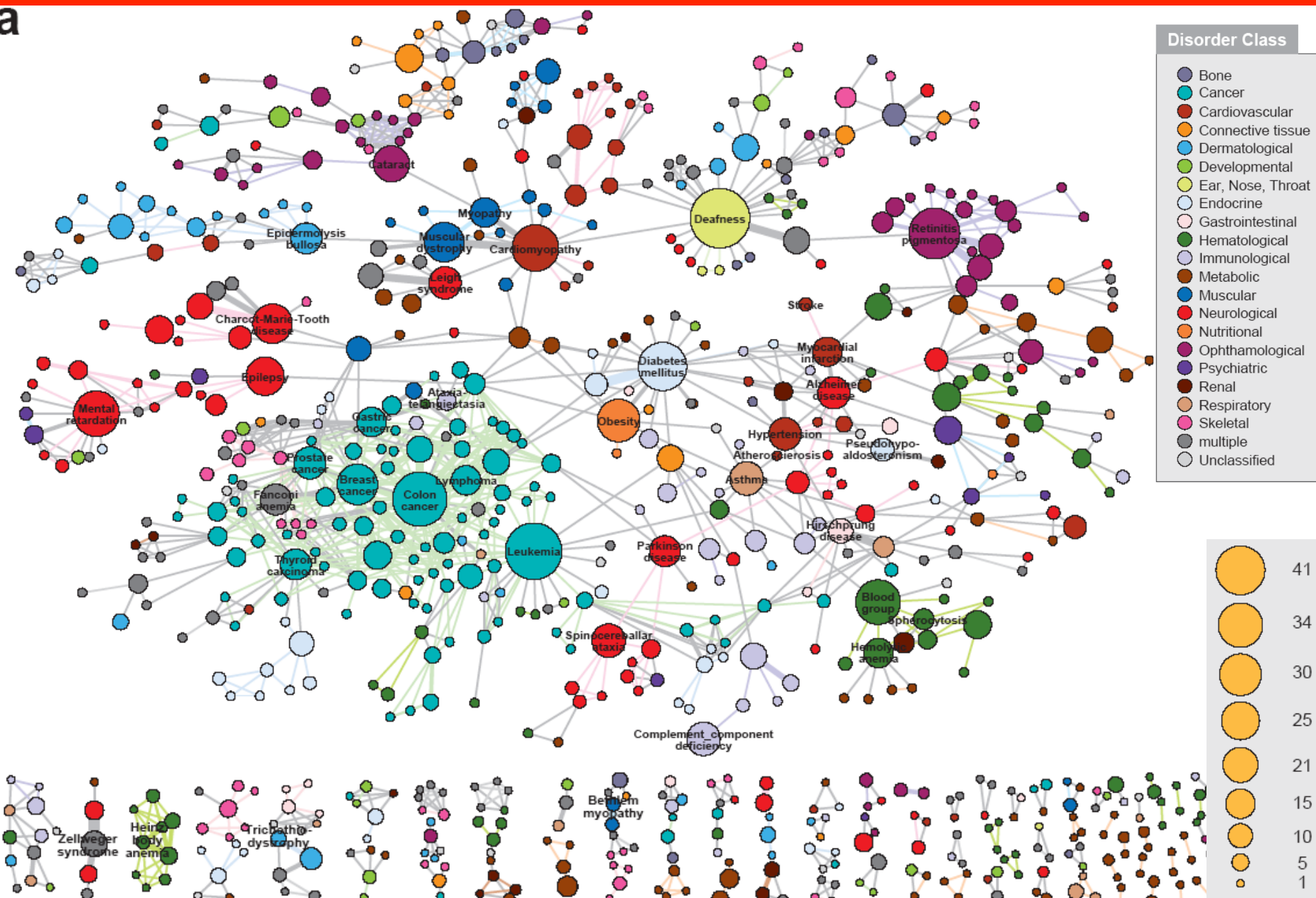


HUMANS GENES



HUMAN DISEASE NETWORK

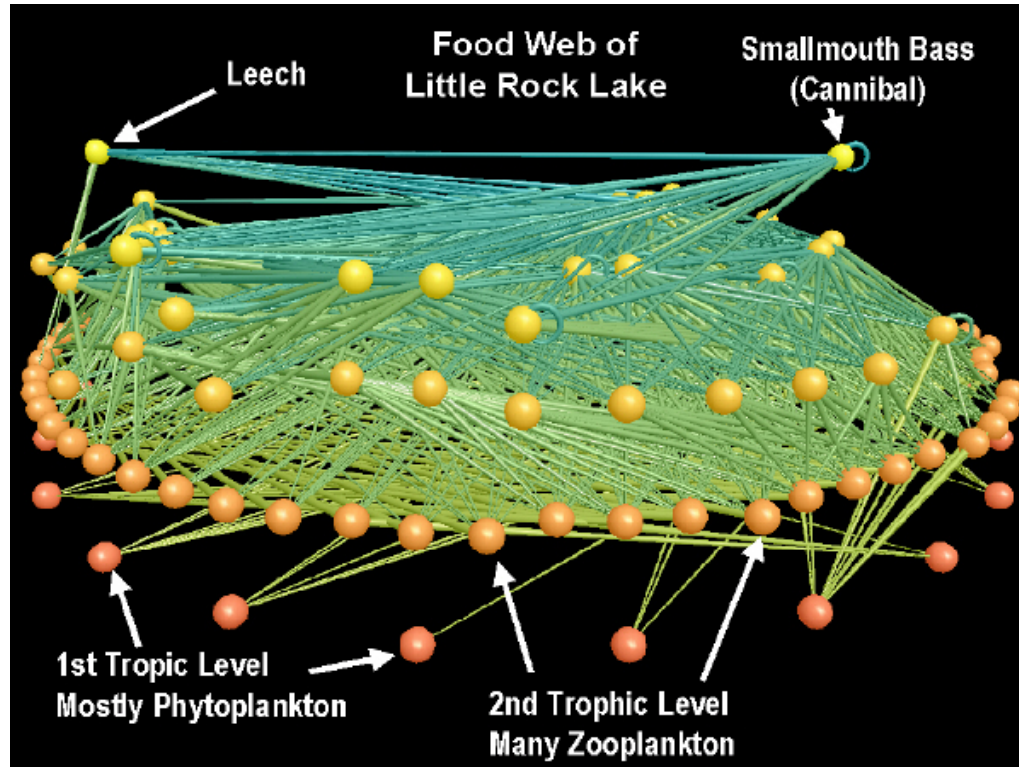
a



Biological networks: Food Web

Nodes: species

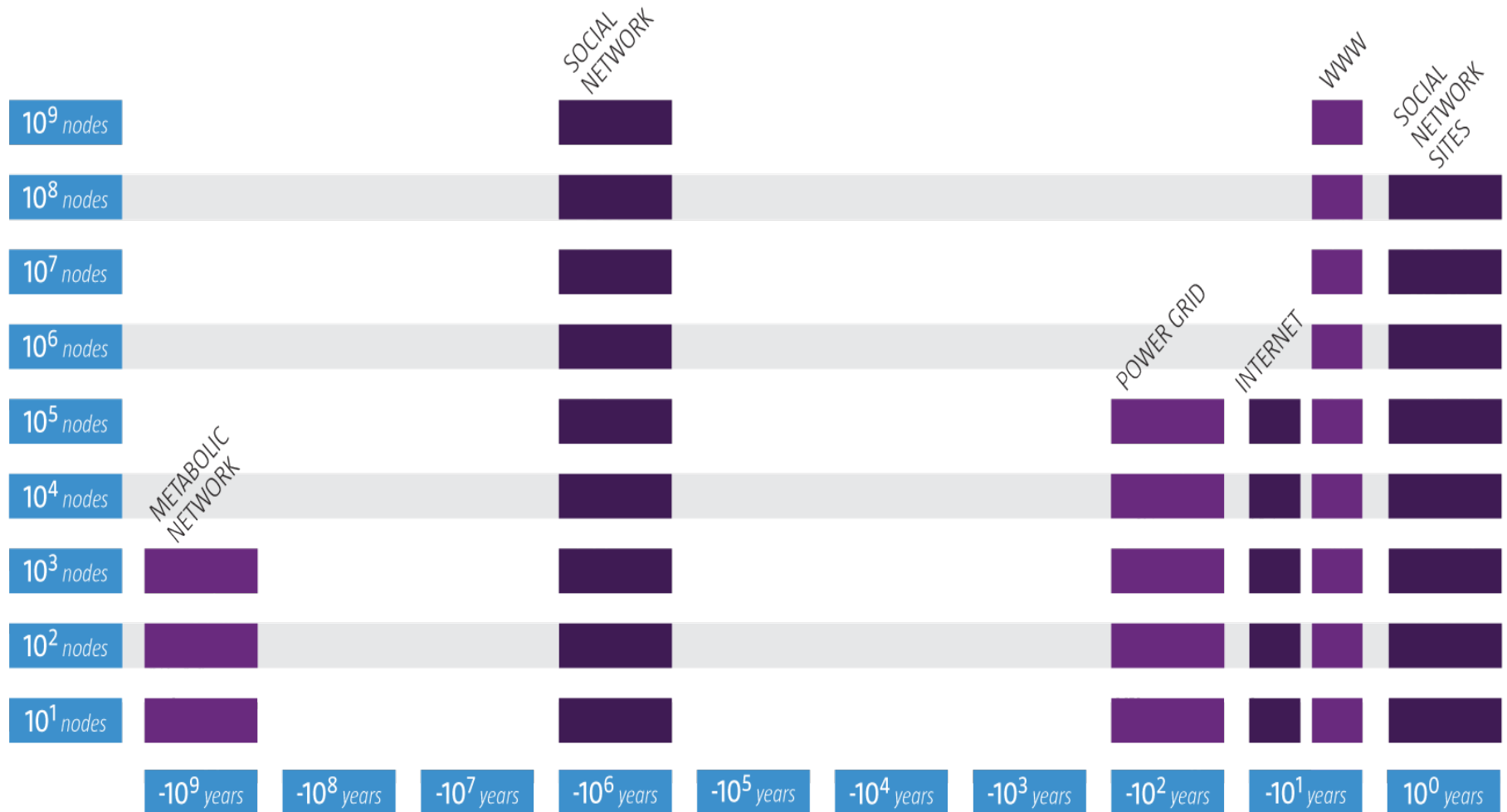
Links: trophic interactions



R. Sole (cond-mat/0011195)

R.J. Williams, N.D. Martinez *Nature* (2000)

THE LIFE OF NETWORKS



Data Availability: Movie Actor Network, 1998;
World Wide Web, 1999.
C elegans neural wiring diagram 1990
Citation Network, 1998
Metabolic Network, 2000;
PPI network, 2001

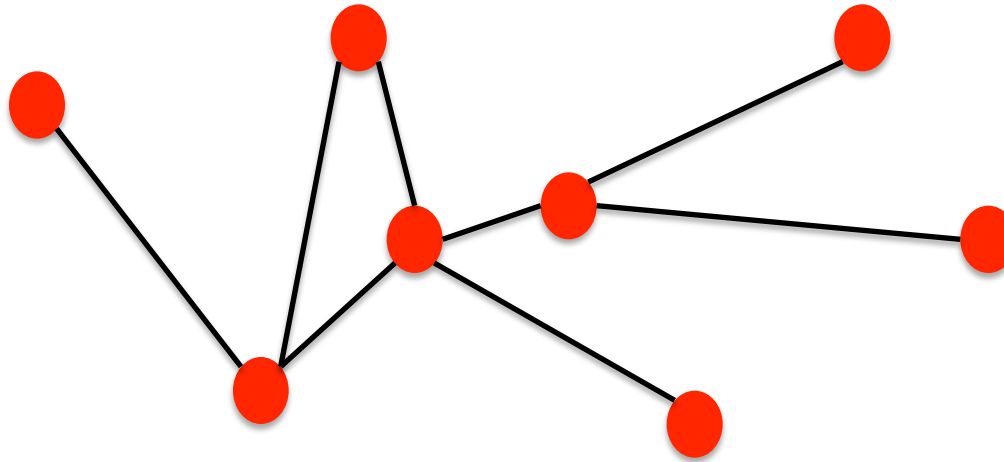
Universality: The architecture of networks emerging in various domains of science, nature, and technology are more similar to each other than one would have expected.

The (urgent) need to understand complexity: Despite the challenges complex systems offer us, we cannot afford to not address their behavior, a view increasingly shared both by scientists and policy makers. Networks are not only essential for this journey, but during the past decade some of the most important advances towards understanding complexity were provided in context of network theory.

A solid red horizontal bar at the top of the slide, divided into two segments by a thin white vertical line.

Networks and graphs

COMPONENTS OF A COMPLEX SYSTEM



▪ **components:** nodes, vertices

N

▪ **interactions:** links, edges

L

▪ **system:** network, graph

(N, L)

NETWORKS OR GRAPHS?

network often refers to real systems

- www,
- social network
- metabolic network.

Language: (Network, node, link)

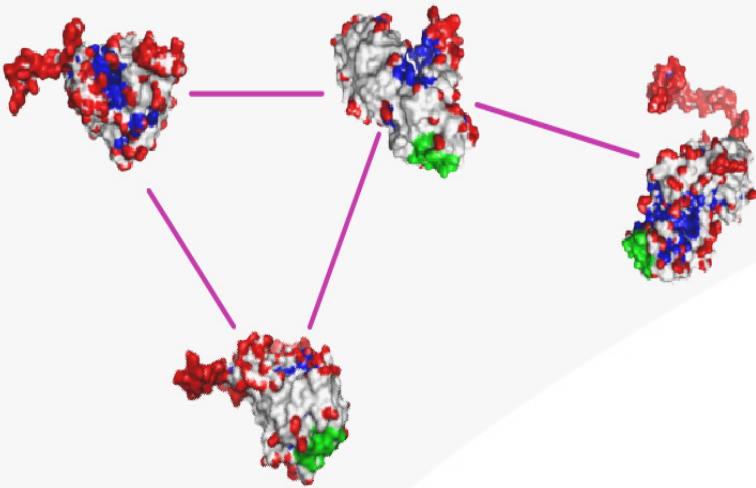
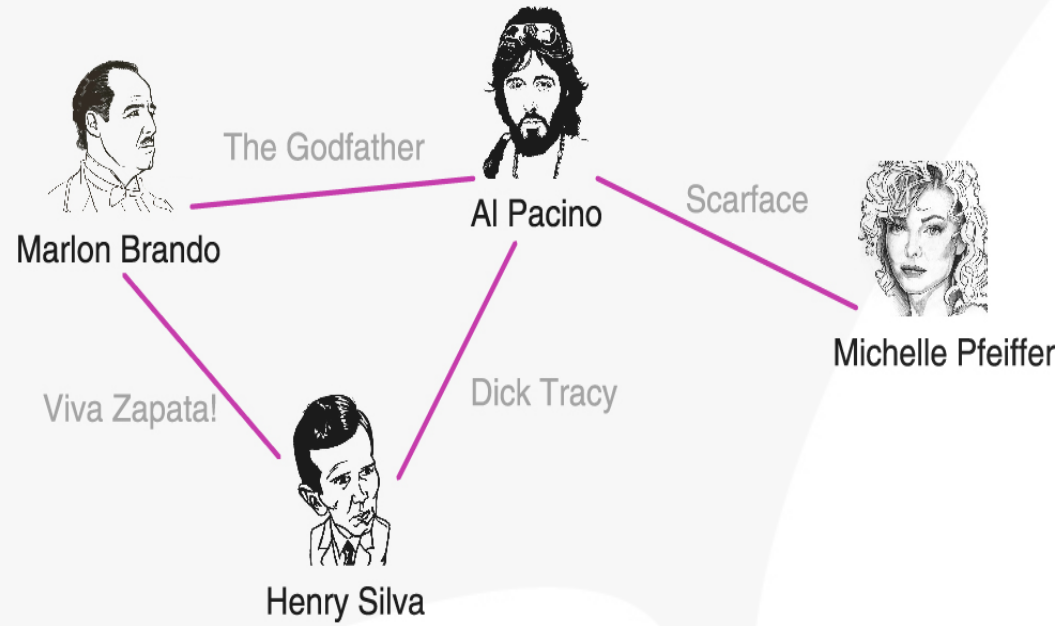
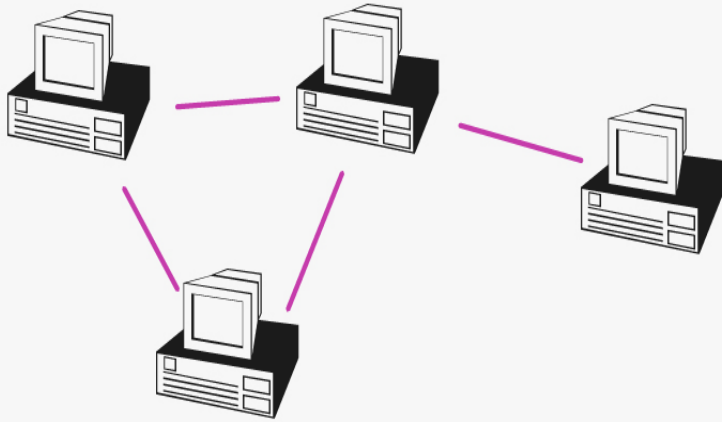
graph: mathematical representation of a network

- web graph,
- social graph (a Facebook term)

Language: (Graph, vertex, edge)

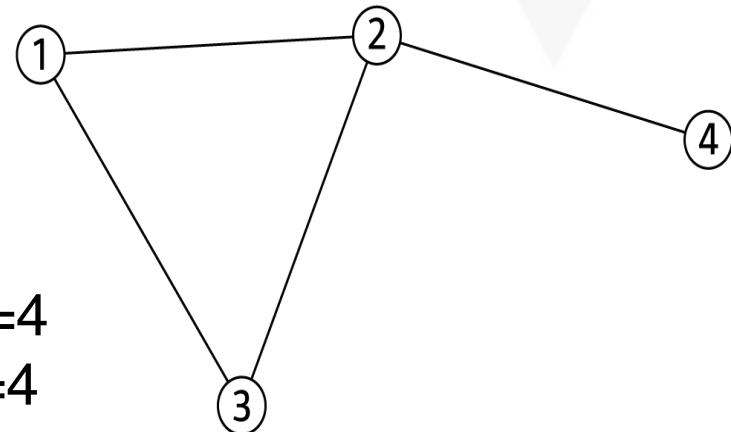
We will try to make this distinction whenever it is appropriate, but in most cases we will use the two terms interchangeably.

A COMMON LANGUAGE



$N=4$

$L=4$

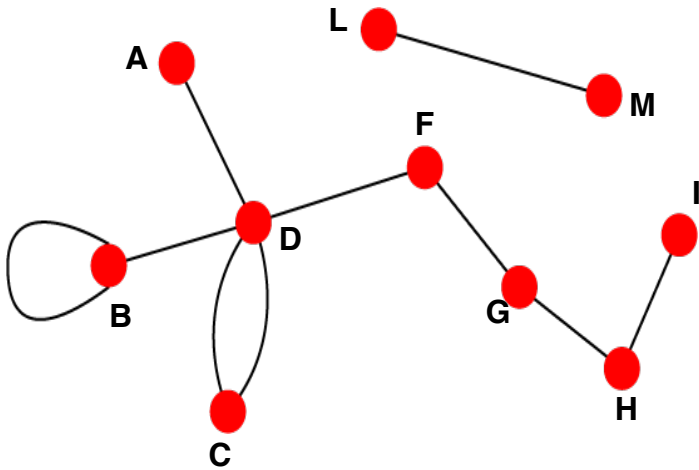


UNDIRECTED VS. DIRECTED NETWORKS

Undirected

Links: undirected (*symmetrical*)

Graph:



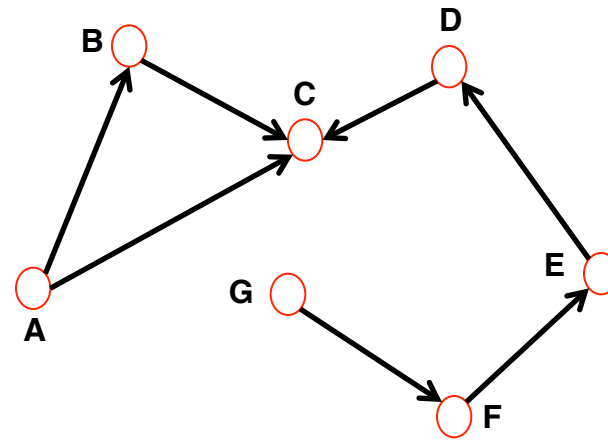
Undirected links :

coauthorship links
Actor network
protein interactions

Directed

Links: directed (*arcs*).

Digraph = directed graph:



An undirected link is the superposition of two opposite directed links.

Directed links :

URLs on the www
phone calls
metabolic reactions

Reference Networks

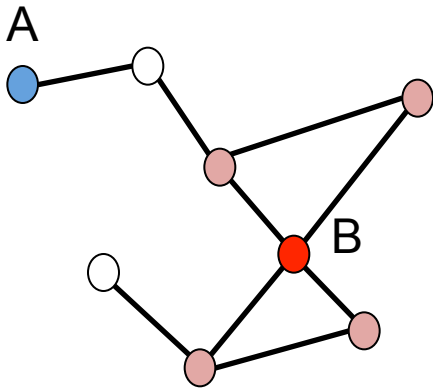
NETWORK	NODES	LINKS	DIRECTED UNDIRECTED	N	L
Internet	Routers	Internet connections	Undirected	192,244	609,066
WWW	Webpages	Links	Directed	325,729	1,497,134
Power Grid	Power plants, transformers	Cables	Undirected	4,941	6,594
Mobile Phone Calls	Subscribers	Calls	Directed	36,595	91,826
Email	Email addresses	Emails	Directed	57,194	103,731
Science Collaboration	Scientists	Co-authorship	Undirected	23,133	93,439
Actor Network	Actors	Co-acting	Undirected	702,388	29,397,908
Citation Network	Paper	Citations	Directed	449,673	4,689,479
E. Coli Metabolism	Metabolites	Chemical reactions	Directed	1,039	5,802
Protein Interactions	Proteins	Binding interactions	Undirected	2,018	2,930

A solid red horizontal bar at the top of the slide, divided into two segments by a thin white vertical line.

Degree, Average Degree and Degree Distribution

NODE DEGREES

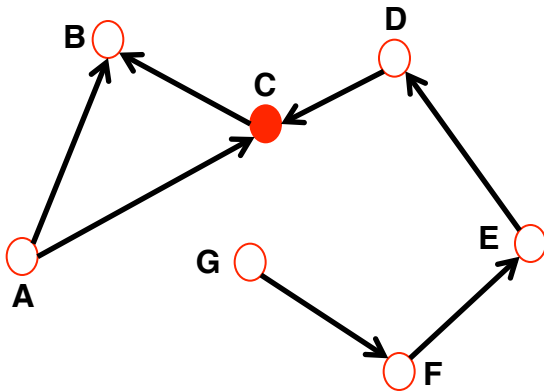
Undirected



Node degree: the number of links connected to the node.

$$k_A = 1 \quad k_B = 4$$

Directed



In *directed networks* we can define an **in-degree** and **out-degree**.

The (total) degree is the sum of in- and out-degree.

$$k_C^{in} = 2 \quad k_C^{out} = 1 \quad k_C = 3$$

Source: a node with $k^{in} = 0$; **Sink:** a node with $k^{out} = 0$.

BRIEF STATISTICS REVIEW

Four key quantities characterize a sample of N values x_1, \dots, x_N :

Average (mean):

$$\langle x \rangle = \frac{x_1 + x_2 + \dots + x_N}{N} = \frac{1}{N} \sum_{i=1}^N x_i$$

The n^{th} moment:

$$\langle x^n \rangle = \frac{x_1^n + x_2^n + \dots + x_N^n}{N} = \frac{1}{N} \sum_{i=1}^N x_i^n$$

Standard deviation:

$$\sigma_x = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \langle x \rangle)^2}$$

Distribution of x :

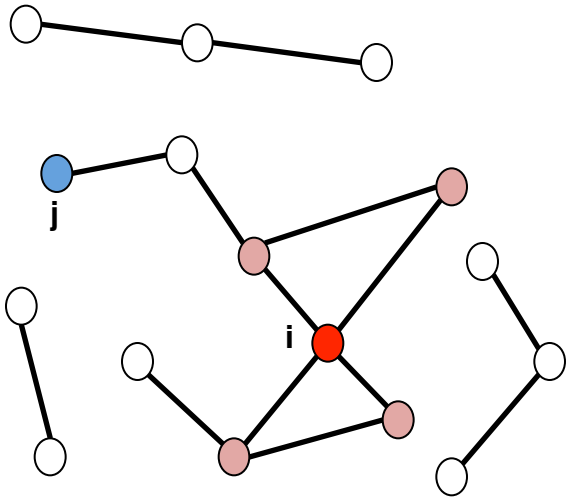
$$p_x = \frac{1}{N} \sum_i \delta_{x, x_i}$$

where p_x follows

$$\sum_i p_x = 1 \quad \left(\int p_x dx = 1 \right)$$

AVERAGE DEGREE

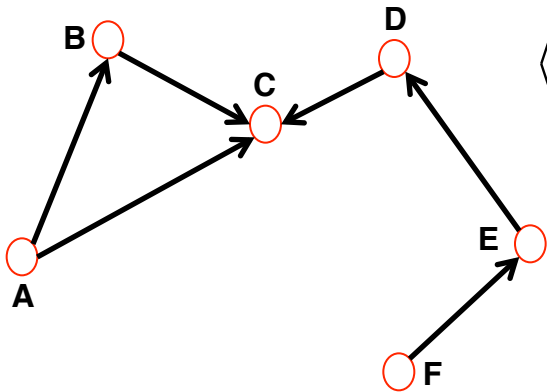
Undirected



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

N – the number of nodes in the graph

Directed



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

$$\langle k \rangle \equiv \frac{L}{N}$$

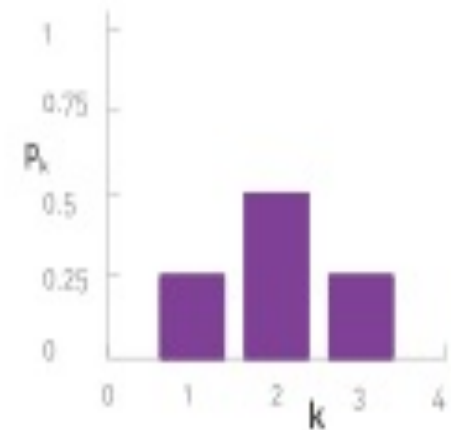
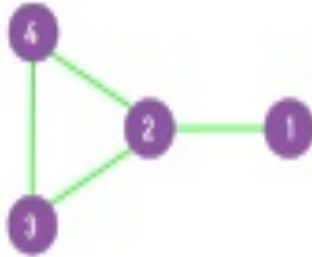
Average Degree

NETWORK	NODES	LINKS	DIRECTED UNDIRECTED	N	L	$\langle k \rangle$
Internet	Routers	Internet connections	Undirected	192,244	609,066	6.33
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-authorship	Undirected	23,133	93,439	8.08
Actor Network	Actors	Co-acting	Undirected	702,388	29,397,908	83.71
Citation Network	Paper	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

DEGREE DISTRIBUTION

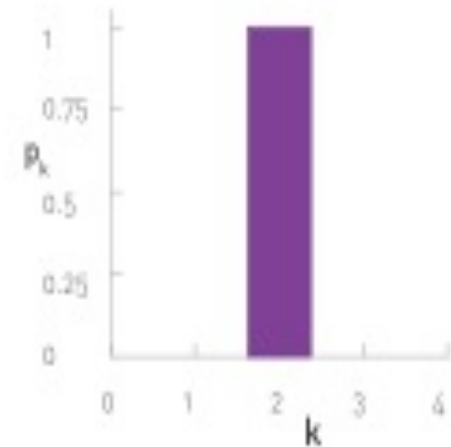
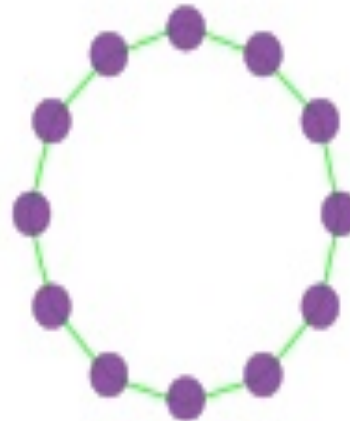
Degree distribution

$P(k)$: probability that a randomly chosen node has degree k



$N_k = \# \text{ nodes with degree } k$

$P(k) = N_k / N \rightarrow \text{plot}$



DEGREE DISTRIBUTION

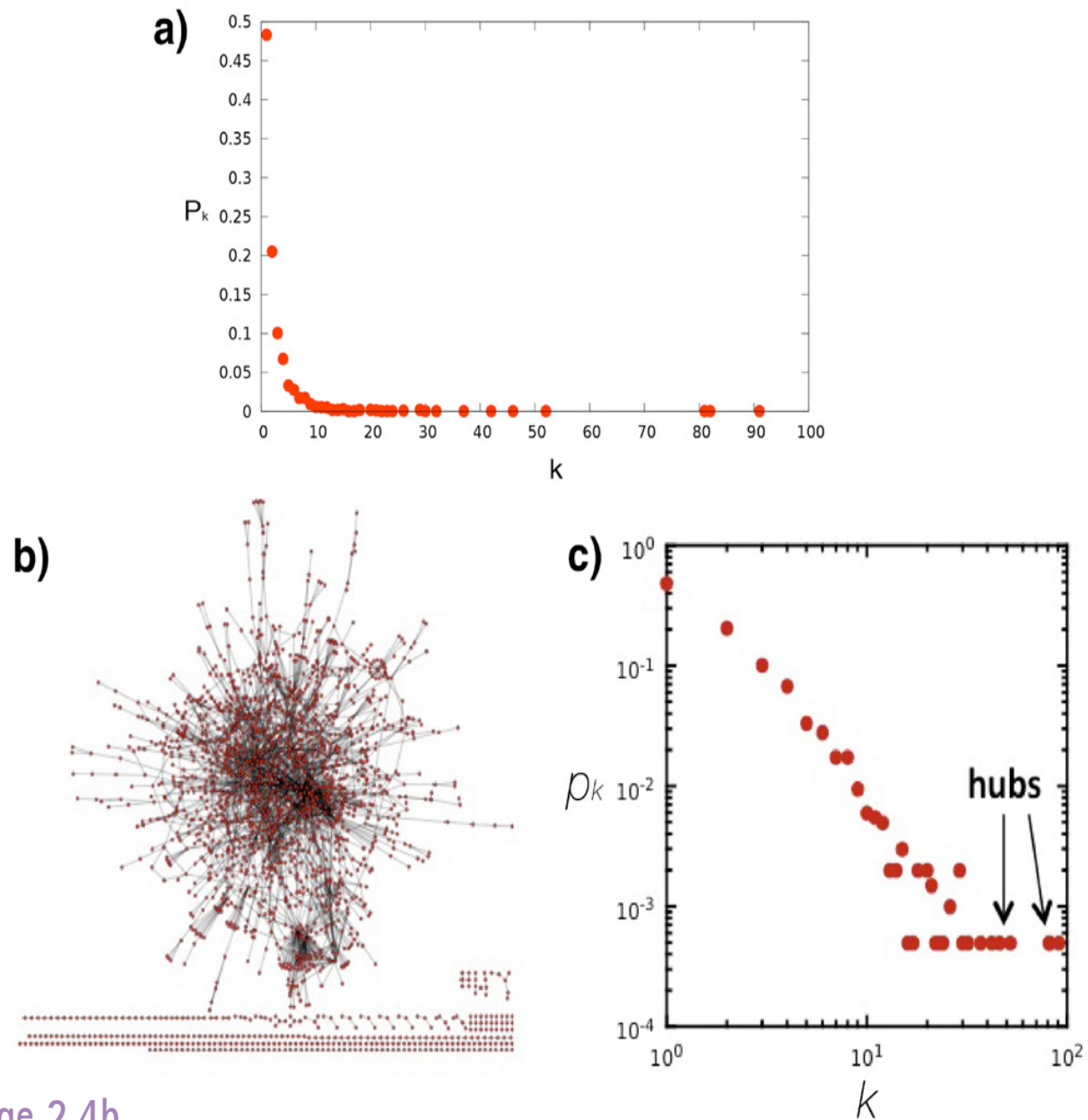


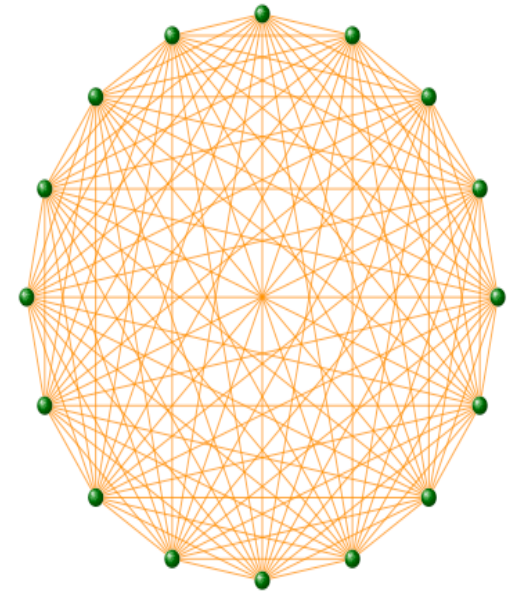
Image 2.4b

A solid red horizontal bar at the top of the slide, divided into two segments by a thin white vertical line.

Real networks are sparse

COMPLETE GRAPH

The maximum number of links a network of N nodes can have is:

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


A graph with degree $L=L_{\max}$ is called a **complete graph**, and its average degree is $\langle k \rangle = N-1$

Most networks observed in real systems are sparse:

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

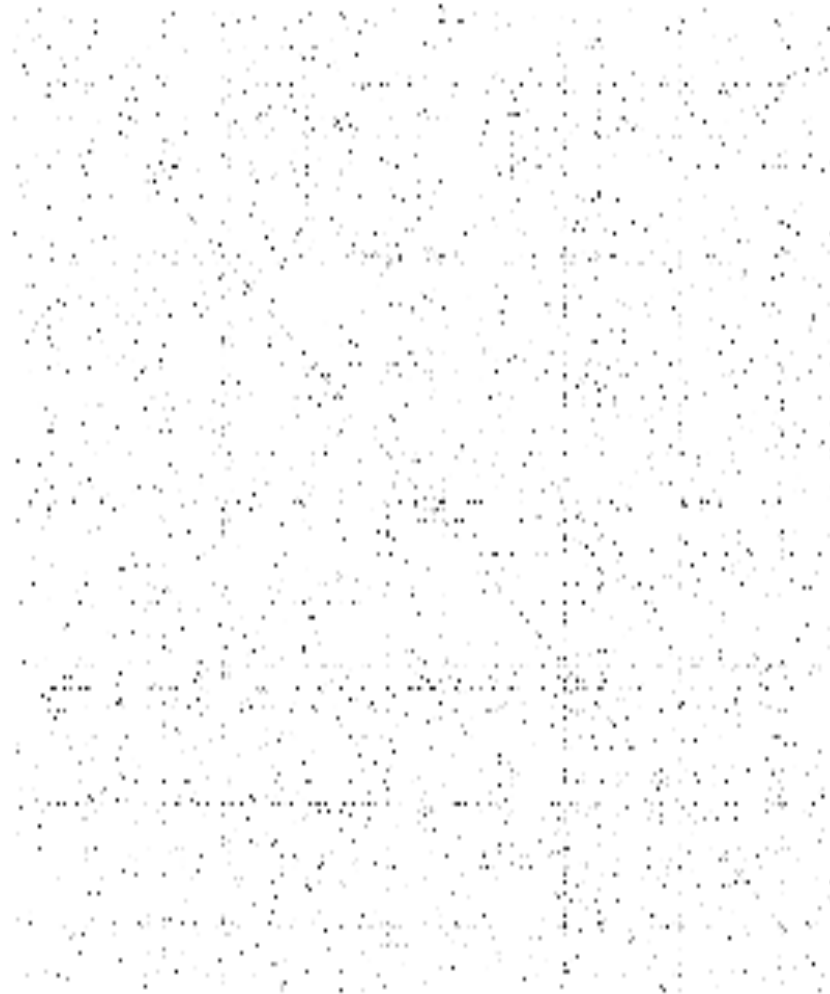
or

$$\langle k \rangle \ll N-1.$$

WWW (ND Sample):	$N=325,729;$	$L=1.4 \cdot 10^6$	$L_{\max}=10^{12}$	$\langle k \rangle=4.51$
Protein (<i>S. Cerevisiae</i>):	$N=1,870;$	$L=4,470$	$L_{\max}=10^7$	$\langle k \rangle=2.39$
Coauthorship (Math):	$N=70,975;$	$L=2 \cdot 10^5$	$L_{\max}=3 \cdot 10^{10}$	$\langle k \rangle=3.9$
Movie Actors:	$N=212,250;$	$L=6 \cdot 10^6$	$L_{\max}=1.8 \cdot 10^{13}$	$\langle k \rangle=28.78$

(Source: Albert, Barabasi, RMP2002)

ADJACENCY MATRICES ARE SPARSE



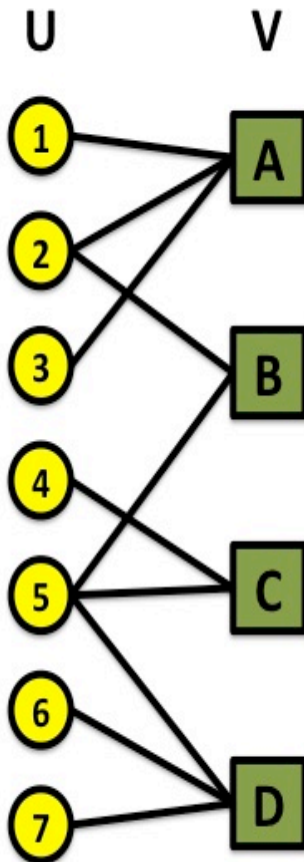
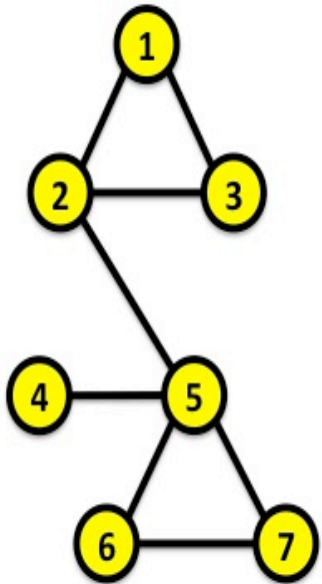
A solid red horizontal bar at the top of the slide, divided into two equal segments by a thin white vertical line.

BIPARTITE NETWORKS

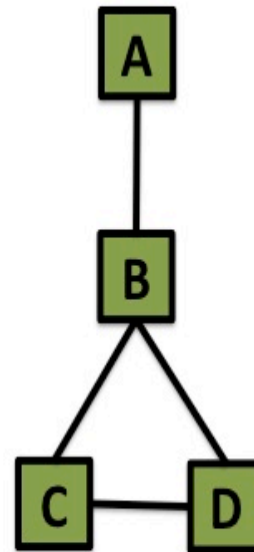
BIPARTITE GRAPHS

bipartite graph (or **bigraph**) is a [graph](#) whose nodes can be divided into two [disjoint sets](#) U and V such that every link connects a node in U to one in V ; that is, U and V are [independent sets](#).

Projection U



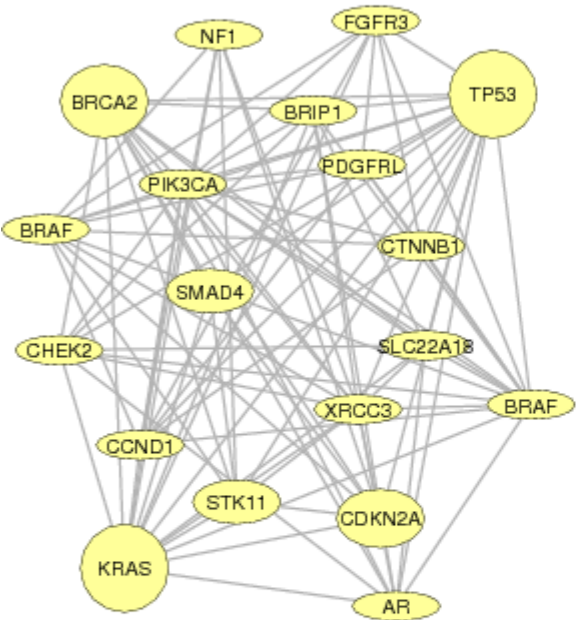
Projection V



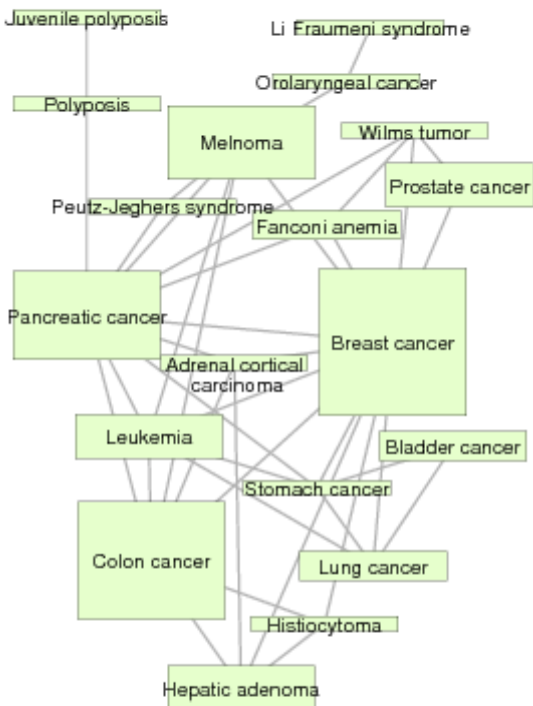
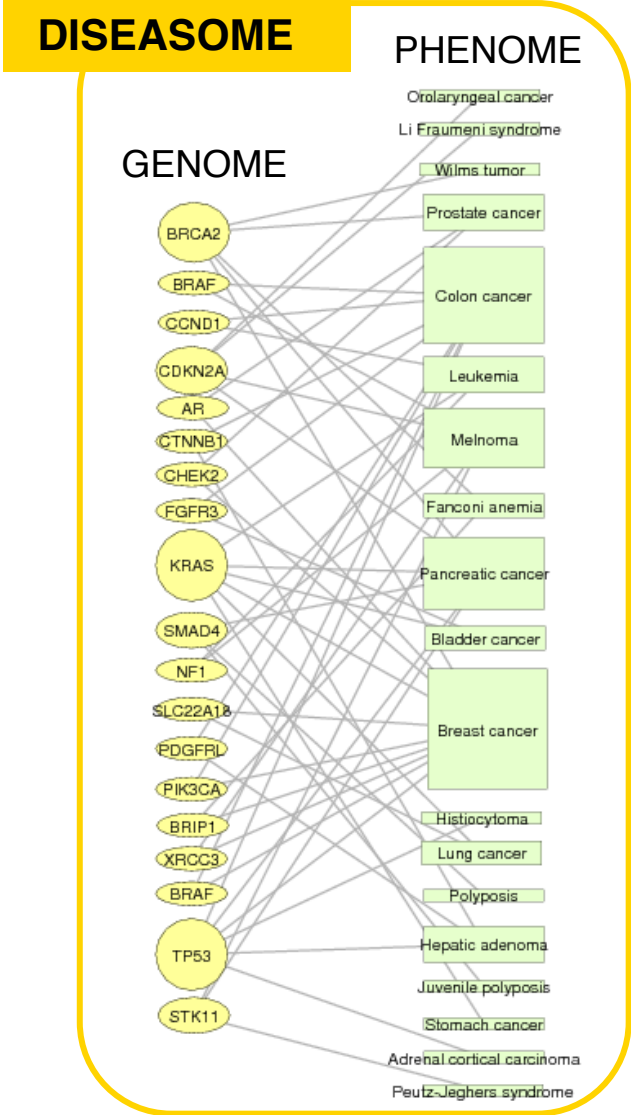
Examples:

Hollywood actor network
Collaboration networks
Disease network (diseasome)

GENE NETWORK – DISEASE NETWORK



Gene network

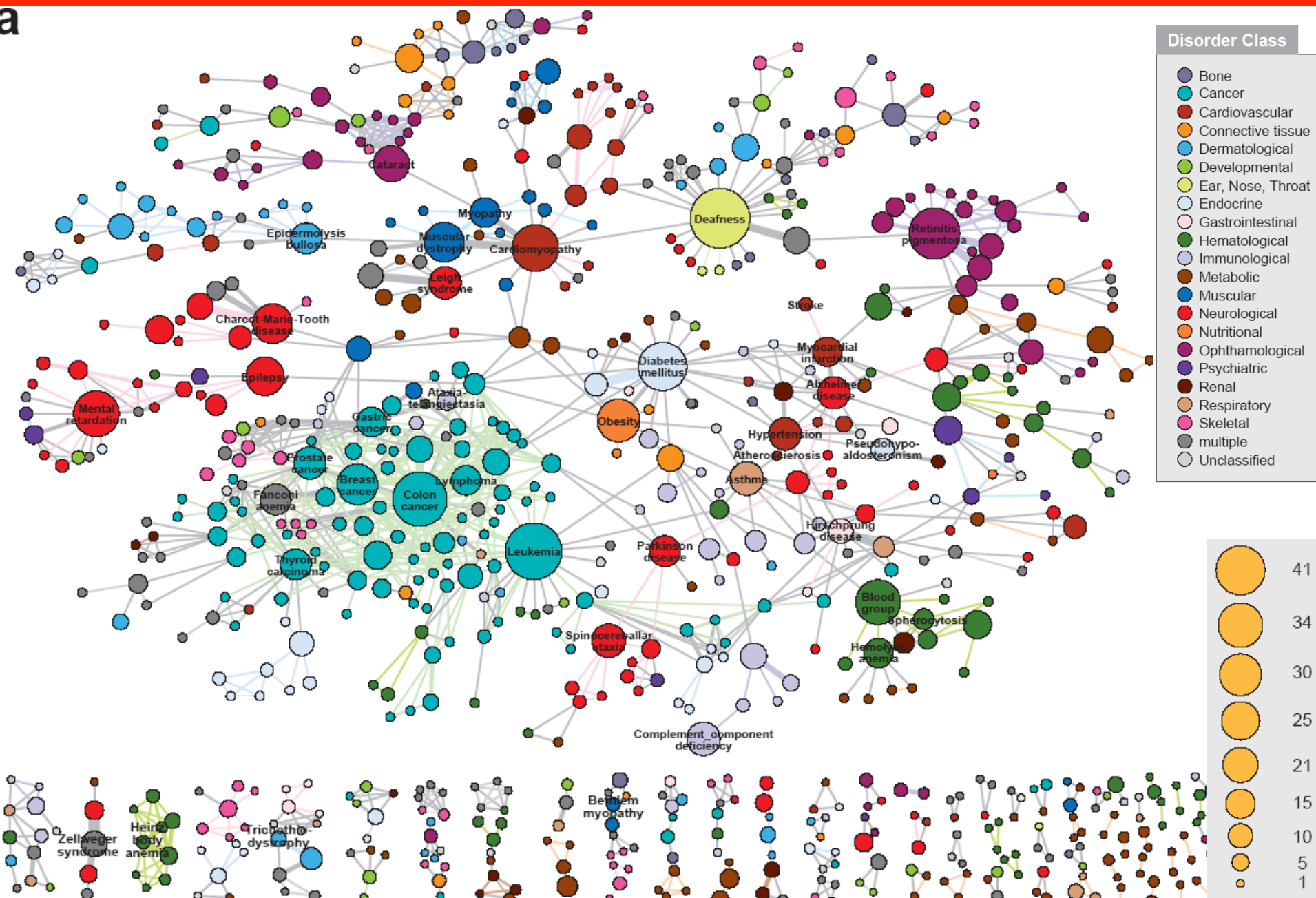


Disease network

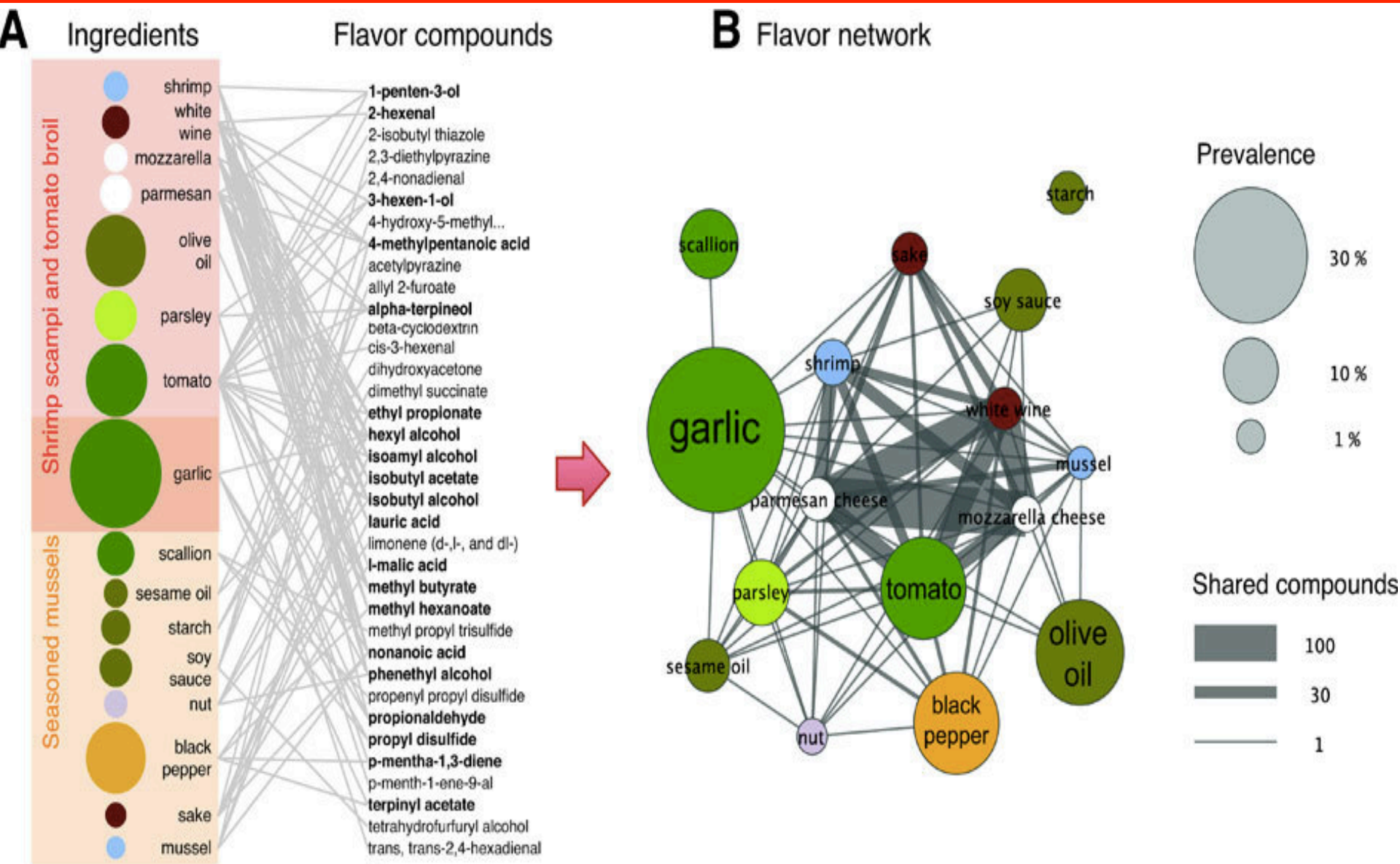
Goh, Cusick, Valle, Childs, Vidal & Barabási, PNAS (2007)

HUMAN DISEASE NETWORK

a



Ingredient-Flavor Bipartite Network



Y.-Y. Ahn, S. E. Ahnert, J. P. Bagrow, A.-L. Barabási

Flavor network and the principles of food pairing , Scientific Reports 196, (2011).

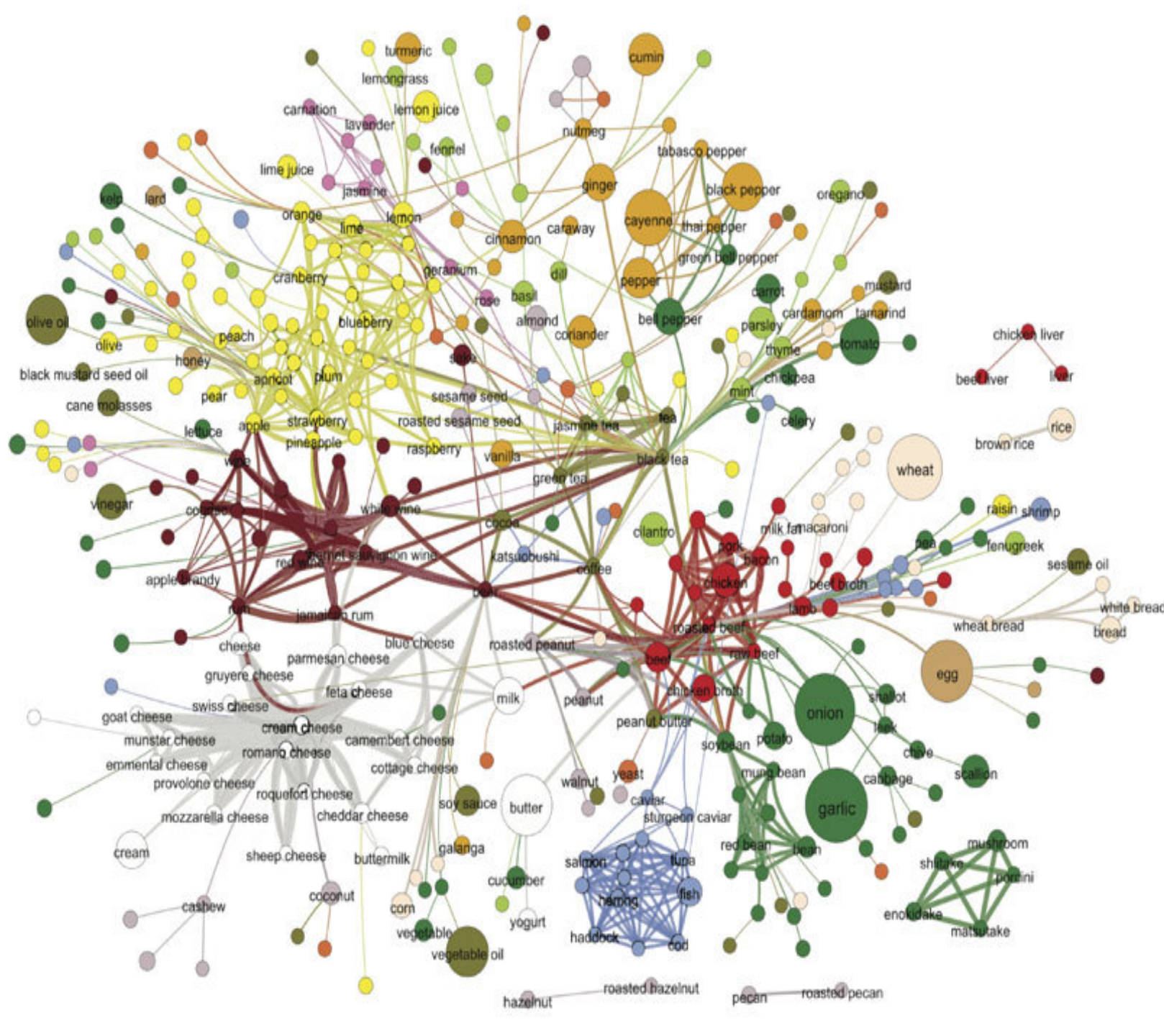
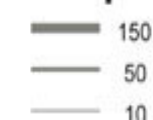
Categories

- fruits
- dairy
- spices
- alcoholic beverages
- nuts and seeds
- seafoods
- meats
- herbs
- plant derivatives
- vegetables
- flowers
- animal products
- plants
- cereal

Prevalence



Shared compounds



Basic network measures

Degree of a node

Distance between two nodes

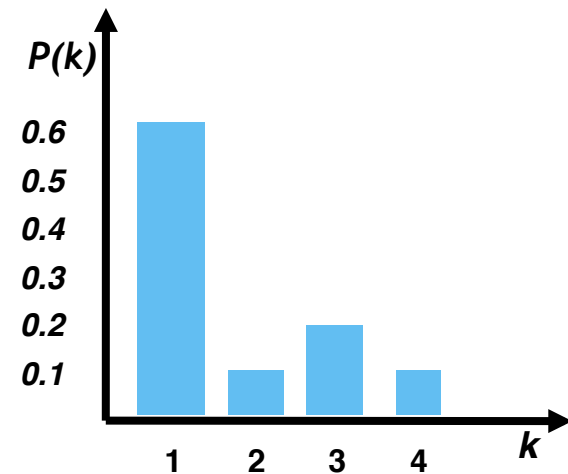
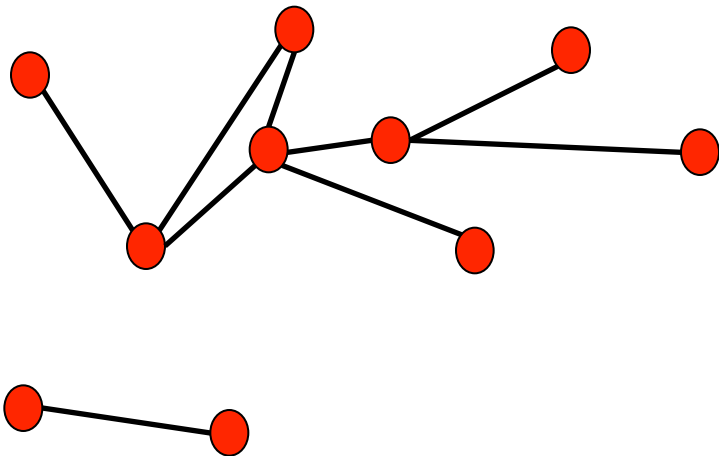
Clustering among three nodes

DEGREE DISTRIBUTION

Degree distribution $P(k)$: probability that a randomly chosen vertex has degree k

N_k = # nodes with degree k

$P(k) = N_k / N \rightarrow$ plot

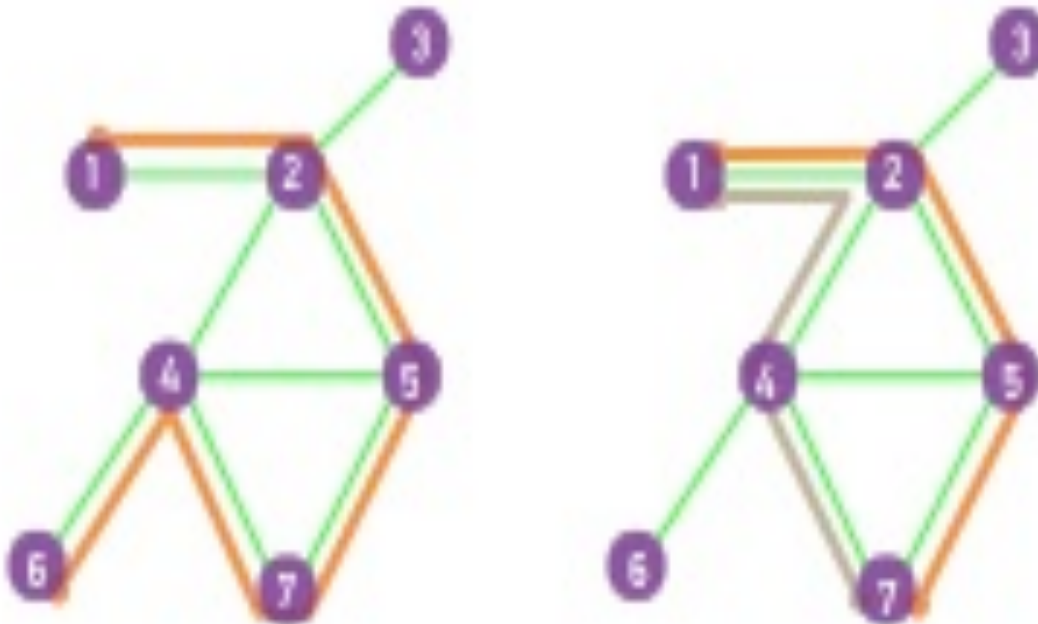


PATHS

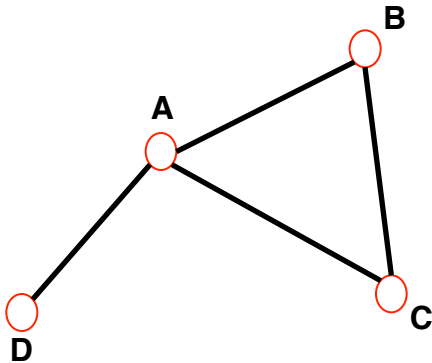
A *path* is a sequence of nodes in which each node is adjacent to the next one

P_{i_0, i_n} of length n between nodes i_0 and i_n is an ordered collection of $n+1$ nodes and n links

$$P_n = \{i_0, i_1, i_2, \dots, i_n\} \quad P_n = \{(i_0, i_1), (i_1, i_2), (i_2, i_3), \dots, (i_{n-1}, i_n)\}$$

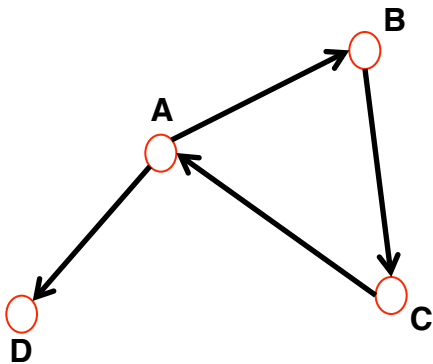


- In a directed network, the path can follow only the direction of an arrow.



The *distance (shortest path, geodesic path)* between two nodes is defined as the number of edges along the shortest path connecting them.

*If the two nodes are disconnected, the distance is infinity.



In **directed graphs** each path needs to follow the direction of the arrows.

Thus in a digraph the distance from node A to B (on an AB path) is generally different from the distance from node B to A (on a BCA path).

NETWORK DIAMETER AND AVERAGE DISTANCE

Diameter: d_{\max} the maximum distance between any pair of nodes in the graph.

Average path length/distance, $\langle d \rangle$, for a **connected graph**:

$$\langle d \rangle \equiv \frac{1}{2L_{\max}} \sum_{i,j \neq i} d_{ij} \quad \text{where } d_{ij} \text{ is the distance from node } i \text{ to node } j$$

In an *undirected graph* $d_{ij} = d_{ji}$, so we only need to count them once:

$$\langle d \rangle \equiv \frac{1}{L_{\max}} \sum_{i,j > i} d_{ij}$$

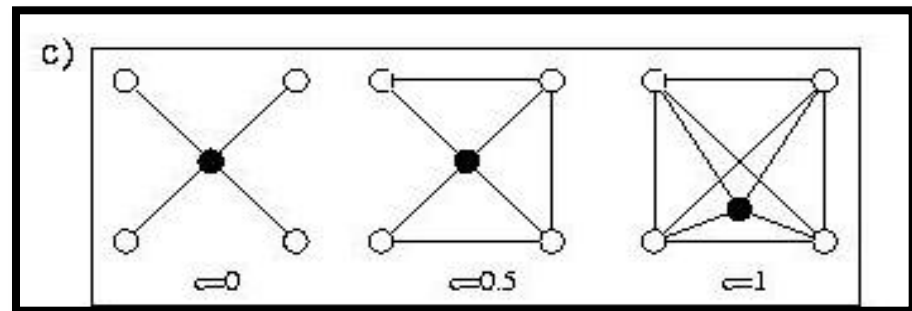
* Clustering coefficient:

what portion of your neighbors are connected?

* Node i with degree k_i

* C_i in $[0,1]$

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



KEY MEASURES

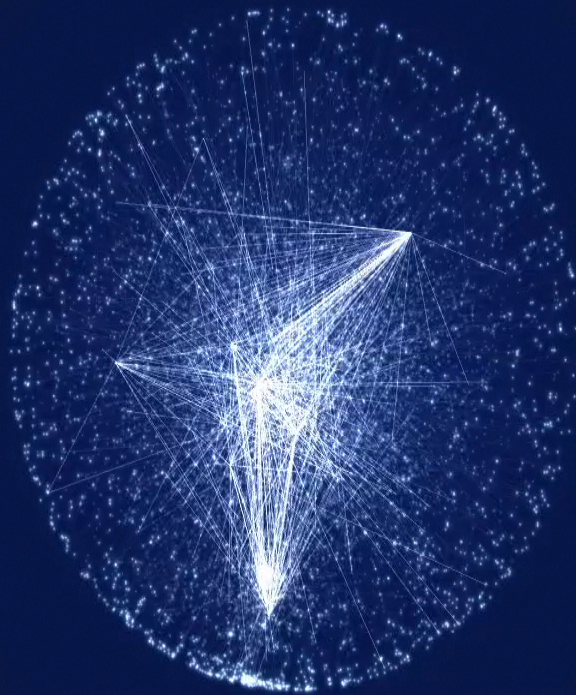
Degree distribution: $P(k)$

Path length: l

Clustering coefficient:

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

A CASE STUDY: PROTEIN-PROTEIN INTERACTION NETWORK



Undirected network

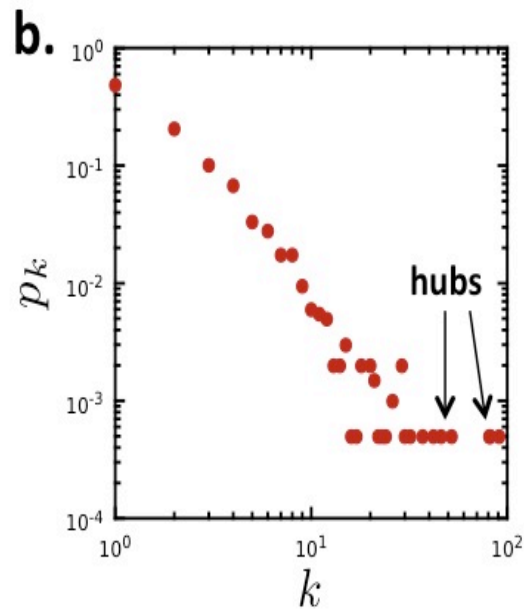
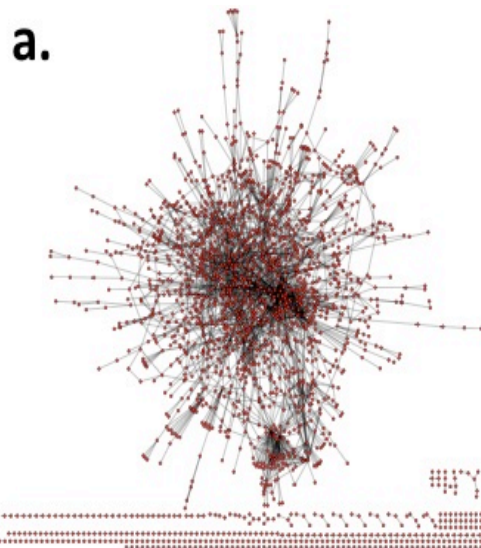
N=2,018 proteins as nodes

L=2,930 binding interactions as links.

Average degree $\langle k \rangle = 2.90$.

Not connected: 185 components
the largest (giant component)
1,647 nodes

A CASE STUDY: PROTEIN-PROTEIN INTERACTION NETWORK

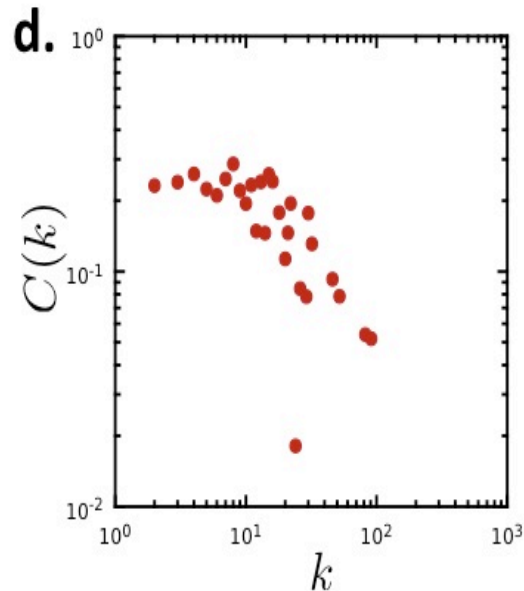
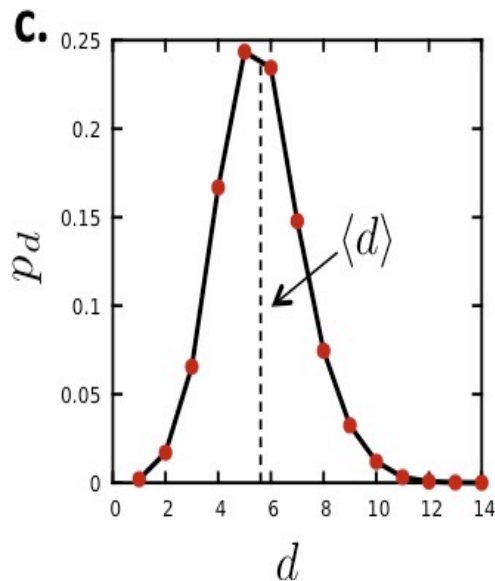


Undirected network

$N=2,018$ proteins as nodes

$L=2,930$ binding interactions as links.

Average degree $\langle k \rangle = 2.90$.

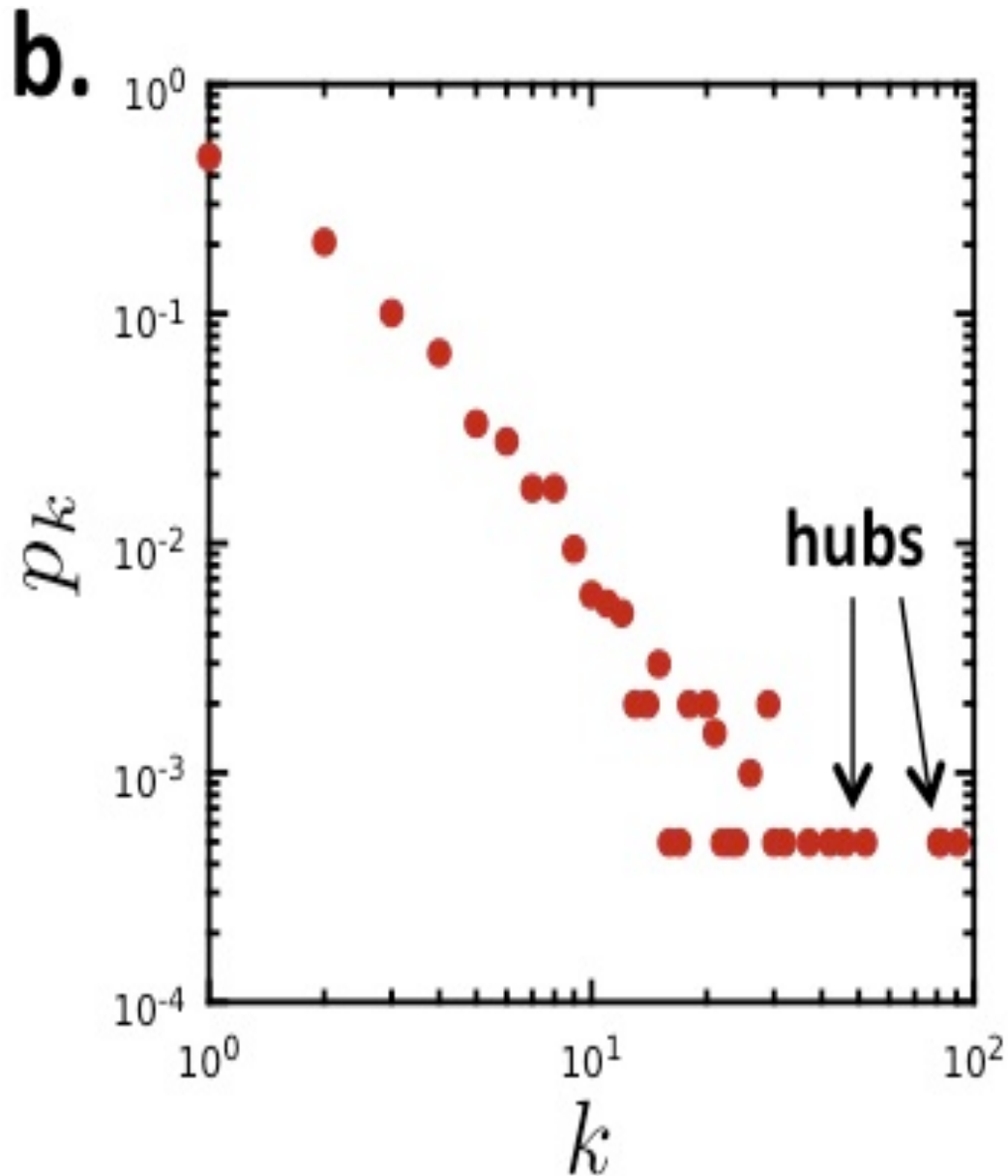


Not connected: 185 components

the largest (giant component)

1,647 nodes

A CASE STUDY: PROTEIN-PROTEIN INTERACTION NETWORK

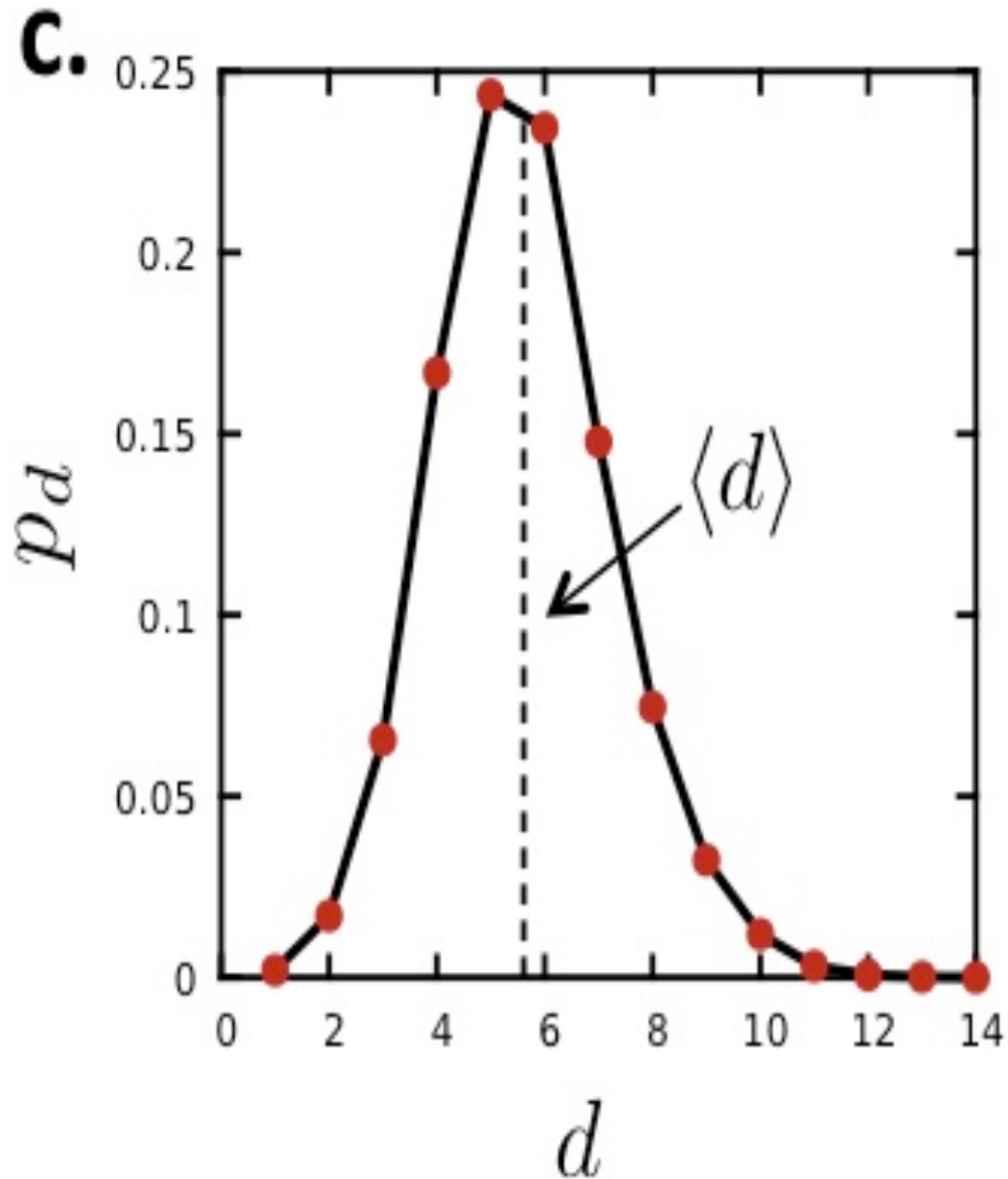


p_k is the probability that a node has degree k .

N_k = # nodes with degree k

$$p_k = N_k / N$$

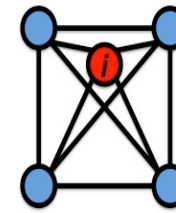
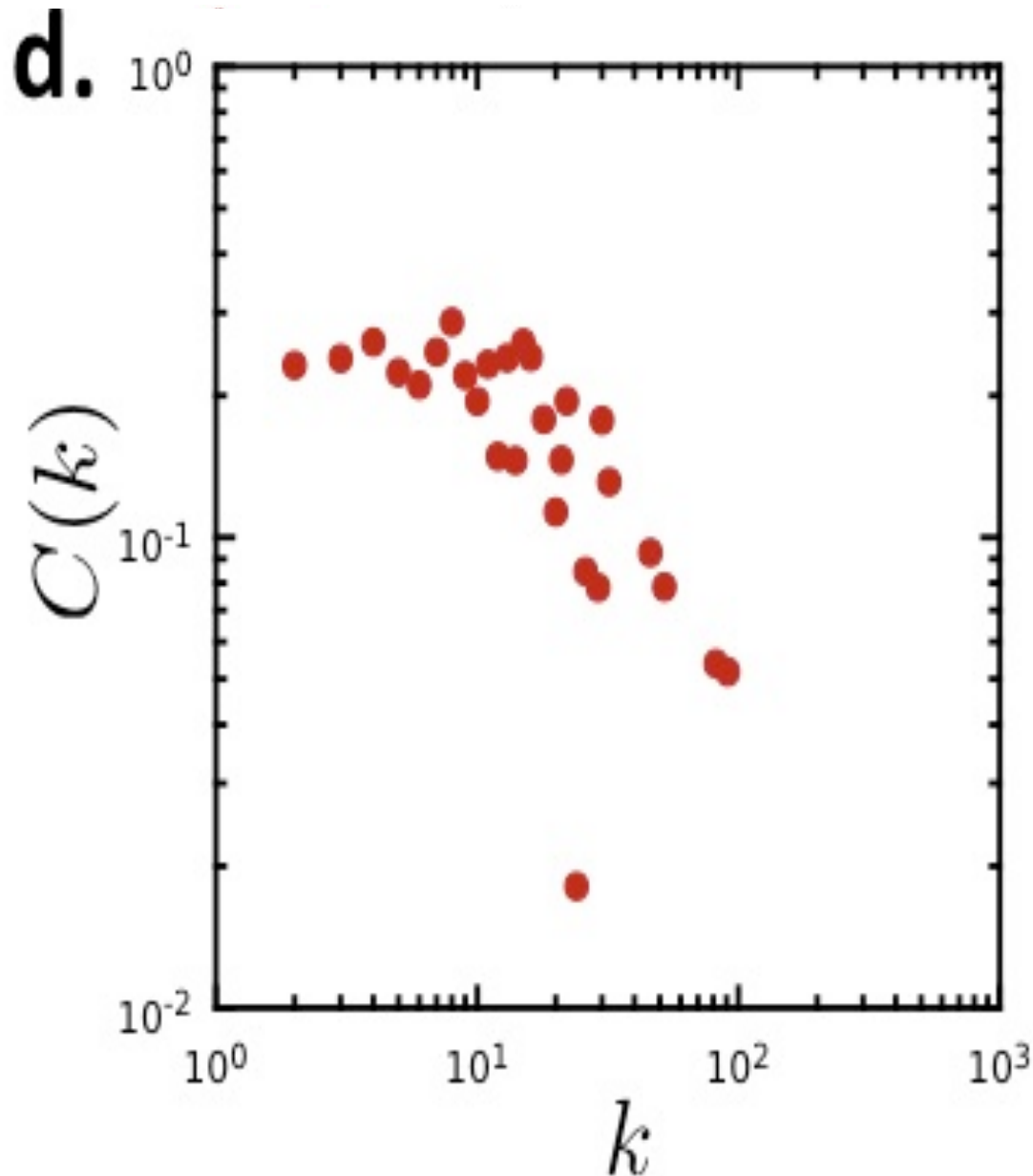
A CASE STUDY: PROTEIN-PROTEIN INTERACTION NETWORK



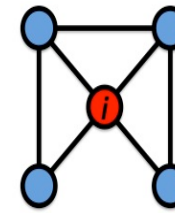
$$d_{\max}=14$$

$$\langle d \rangle = 5.61$$

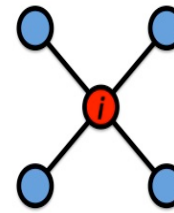
A CASE STUDY: PROTEIN-PROTEIN INTERACTION NETWORK



$$C_i = 1$$



$$C_i = 1/2$$



$$C_i = 0$$

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

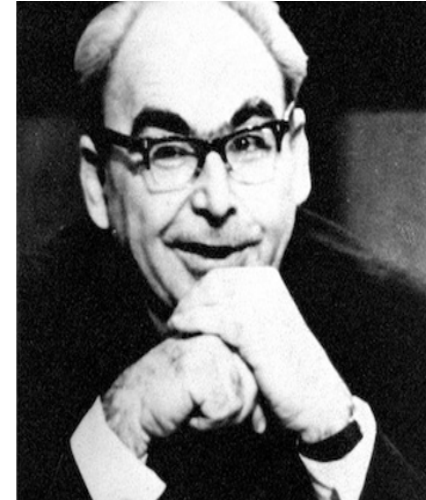
$$\langle C \rangle = 0.12$$

Random graphs

What are the expected basic measures emerging from random?

RANDOM NETWORK MODEL

Pául Erdős
(1913-1996)

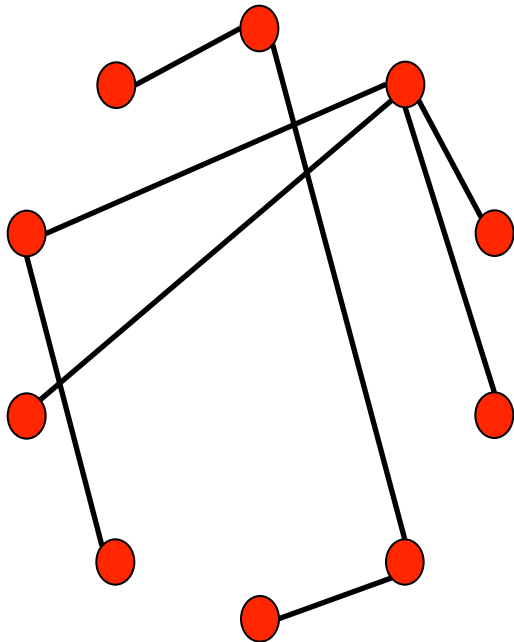


Erdős-Rényi model (1960)

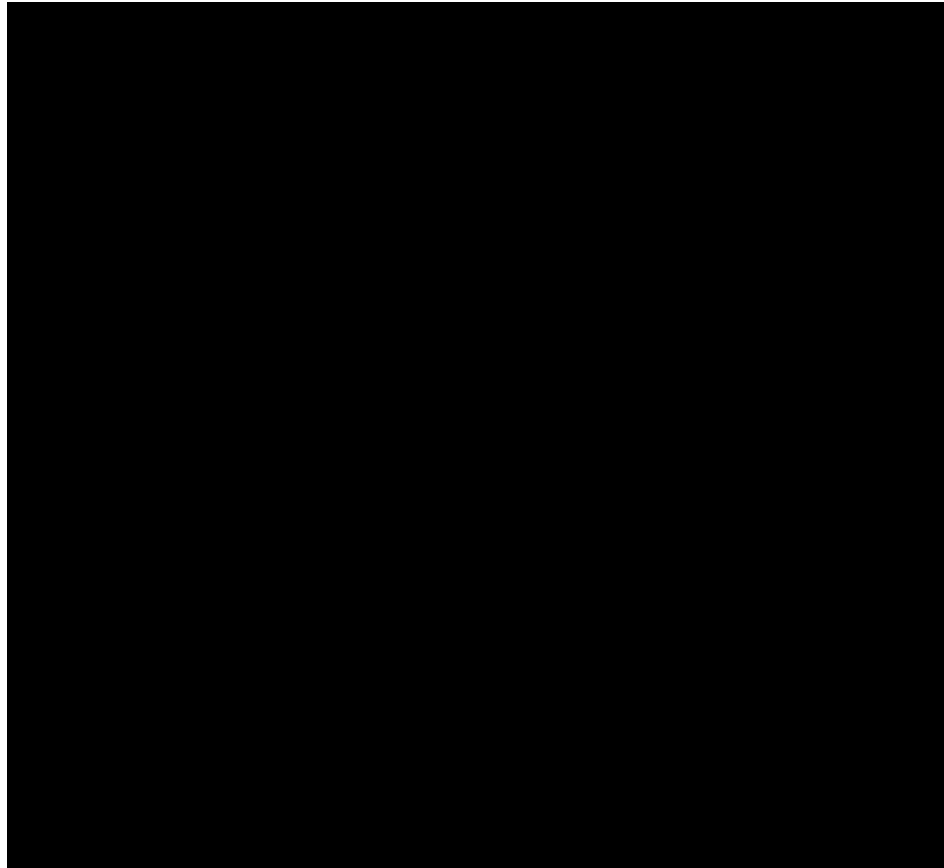
Connect with probability p

$p=1/6$ $N=10$

$\langle k \rangle \sim 1.5$



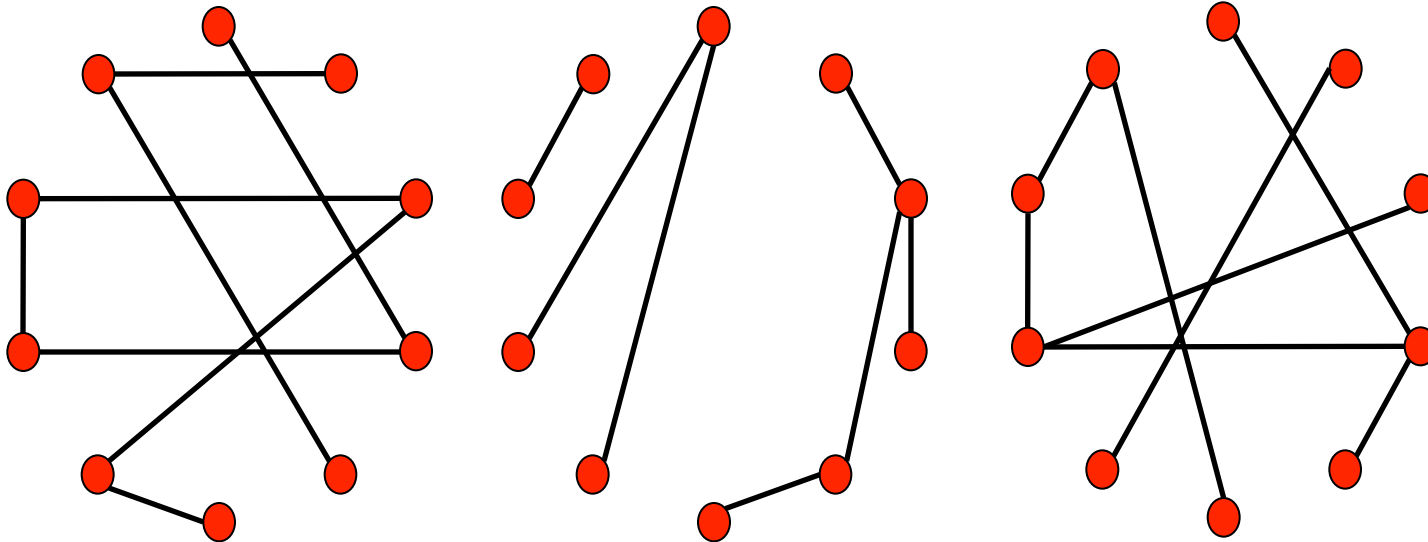
RANDOM NETWORK MODEL



Definition: A **random graph** is a graph of N labeled nodes where each pair of nodes is connected by a preset probability p .

RANDOM NETWORK MODEL

N and p do not uniquely define the network— we can have many different realizations of it. **How many?**



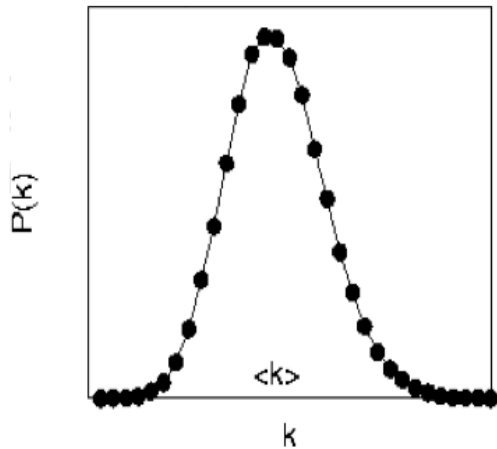
$N=10$
 $p=1/6$

The probability to form a *particular* graph $\mathbf{G}(N,L)$ is

$$P(\mathbf{G}(N,L)) = p^L (1-p)^{\frac{N(N-1)}{2} - L}$$

That is, each graph $\mathbf{G}(N,L)$ appears with probability $P(\mathbf{G}(N,L))$.

DEGREE DISTRIBUTION OF A RANDOM GRAPH



$$P(k) = \binom{N-1}{k} p^k (1-p)^{(N-1)-k}$$

Select k
nodes from $N-1$

probability of
having k edges

probability of
missing $N-1-k$
edges

$$\langle k \rangle = p(N-1)$$

$$\sigma_k^2 = p(1-p)(N-1)$$

$$\frac{\sigma_k}{\langle k \rangle} = \left[\frac{1-p}{p} \frac{1}{(N-1)} \right]^{1/2} \approx \frac{1}{(N-1)^{1/2}}$$

As the network size increases, the distribution becomes increasingly narrow—we are increasingly confident that the degree of a node is in the vicinity of $\langle k \rangle$.

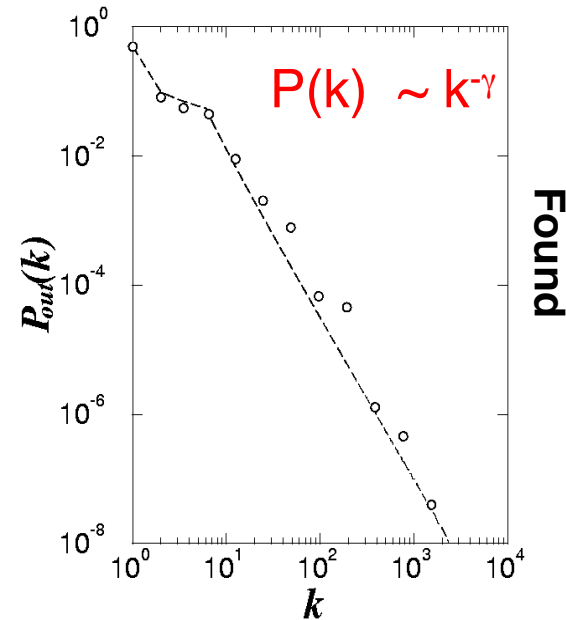
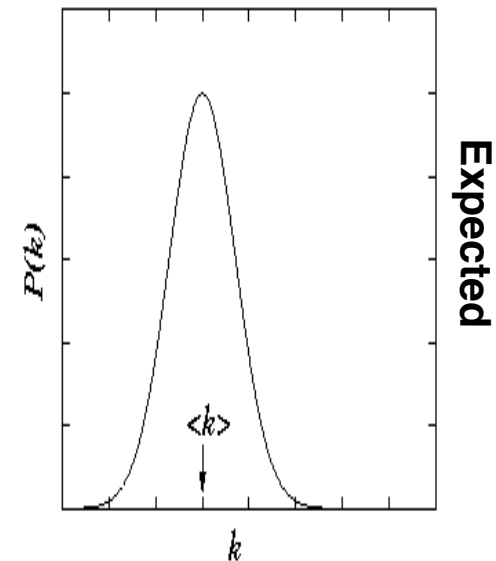
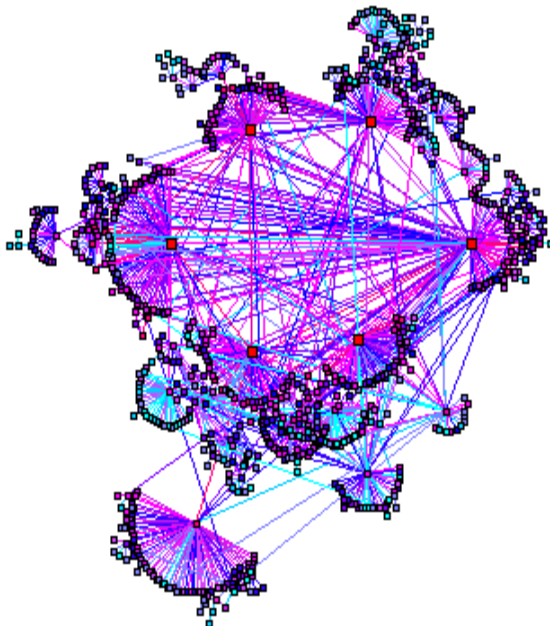
WORLD WIDE WEB

Nodes: **WWW documents**

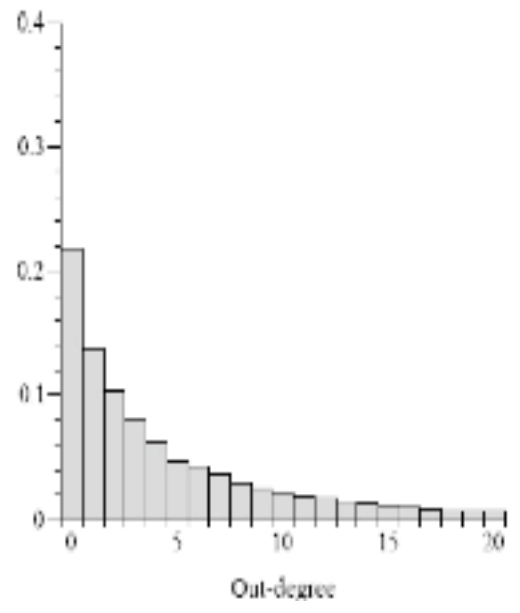
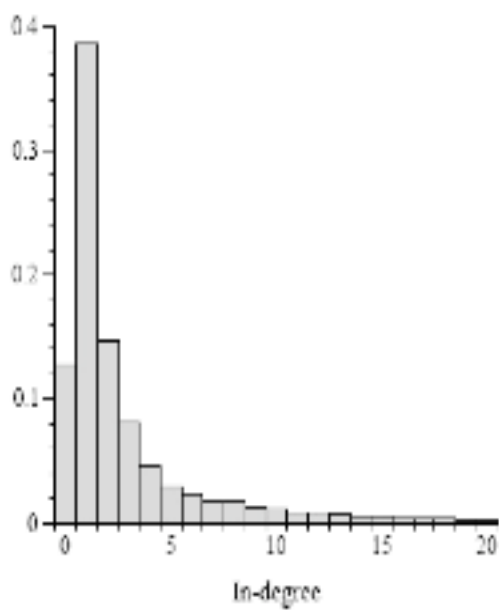
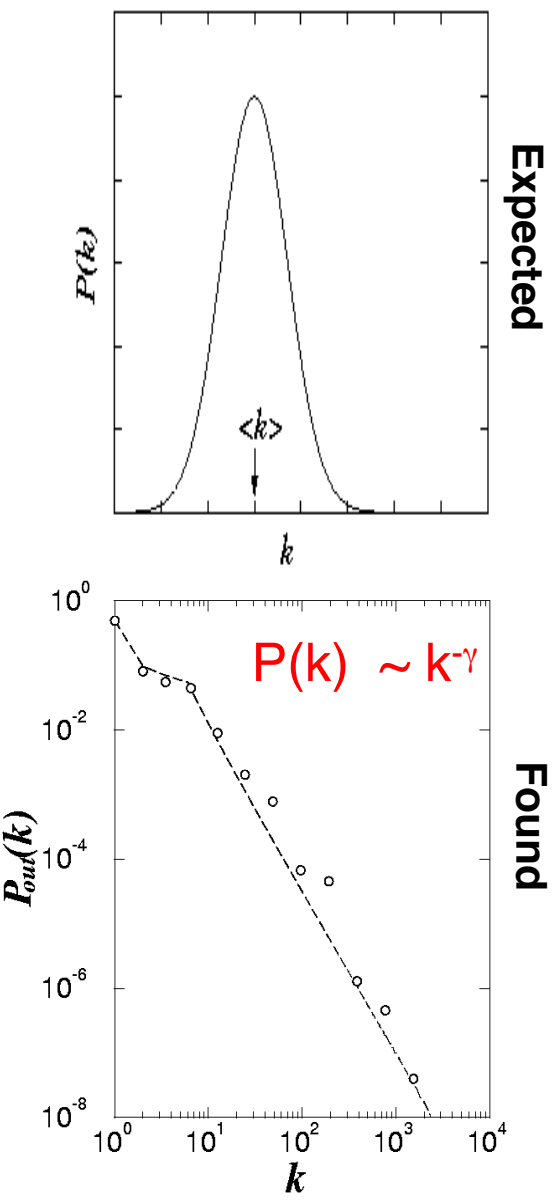
Links: **URL links**

Over 3 billion documents

ROBOT: collects all URL's
found in a document and
follows them recursively

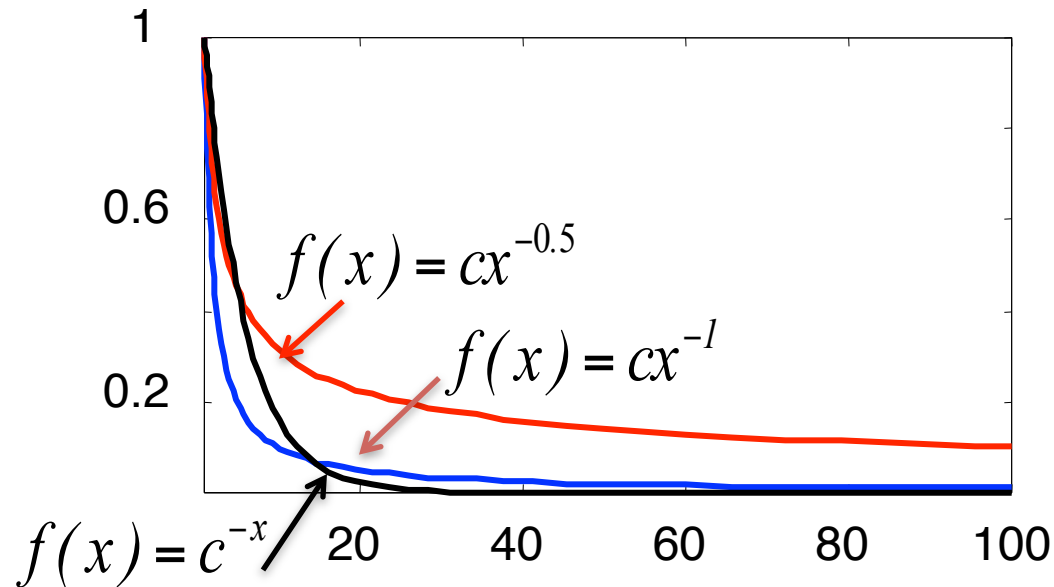


Degree distribution of the WWW



R. Albert, H. Jeong, A-L Barabasi, *Nature*, 401 130 (1999).

The difference between a power law and an exponential distribution



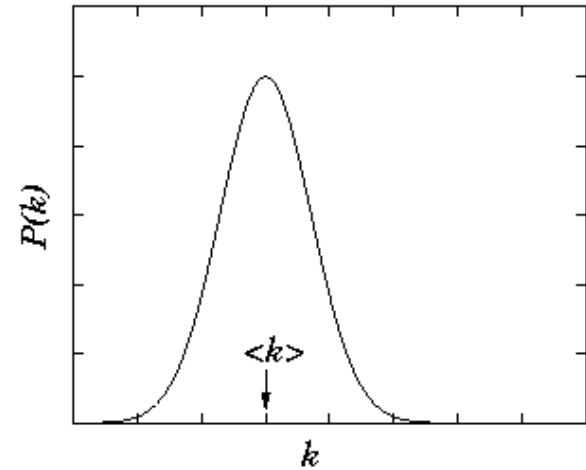
Above a certain x value, the power law is always higher than the exponential.

What does the difference mean? Visual representation.

Exponential
Network

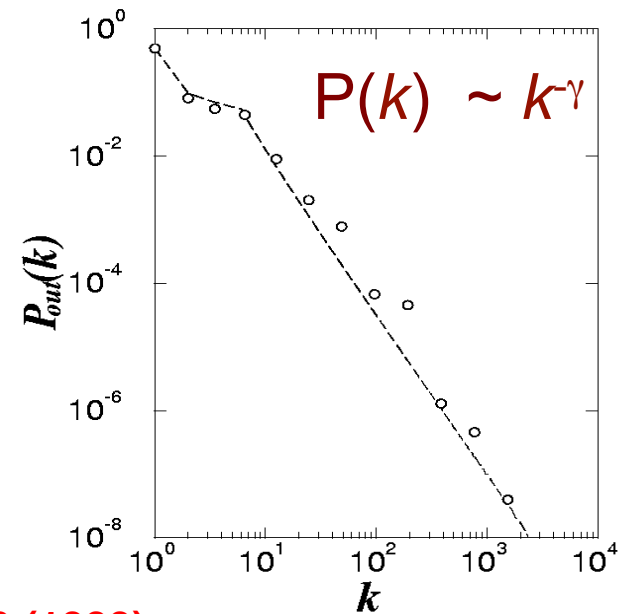
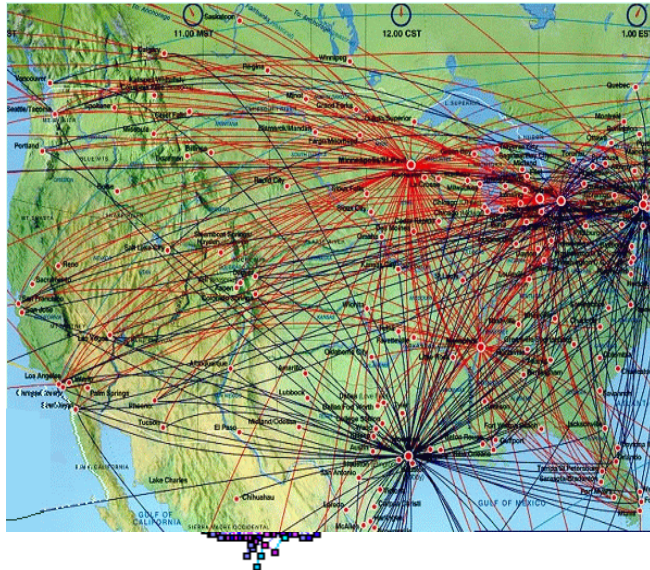


S



Expected

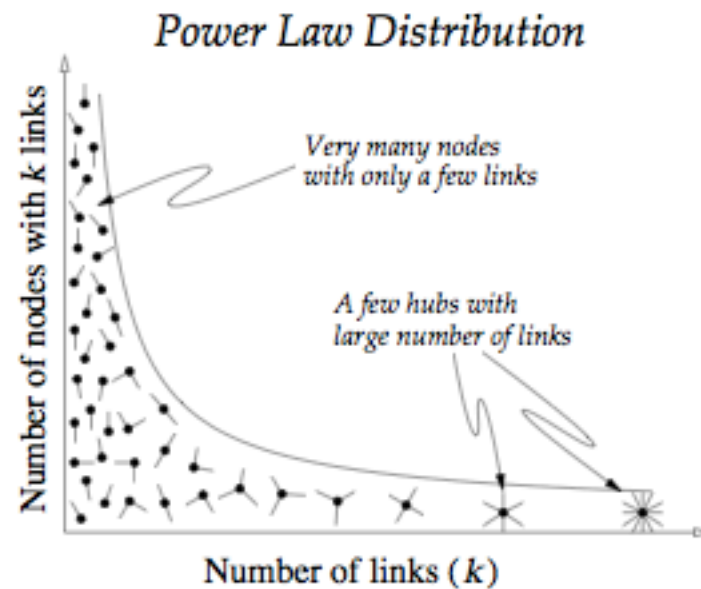
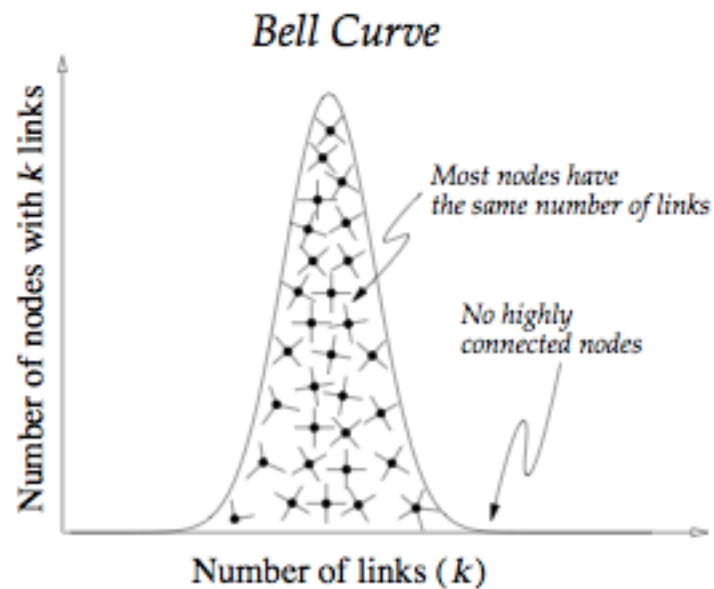
Scale-free
Network



Found

R. Albert, H. Jeong, A-L Barabasi, *Nature*, 401 130 (1999).

WORLD WIDE WEB

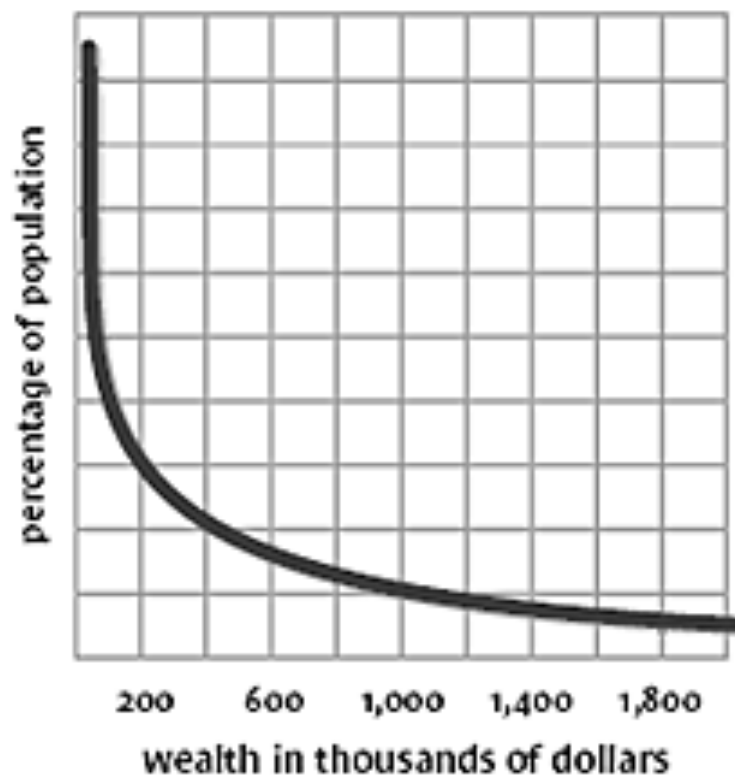


PARETO DISTRIBUTION OF WEALTH

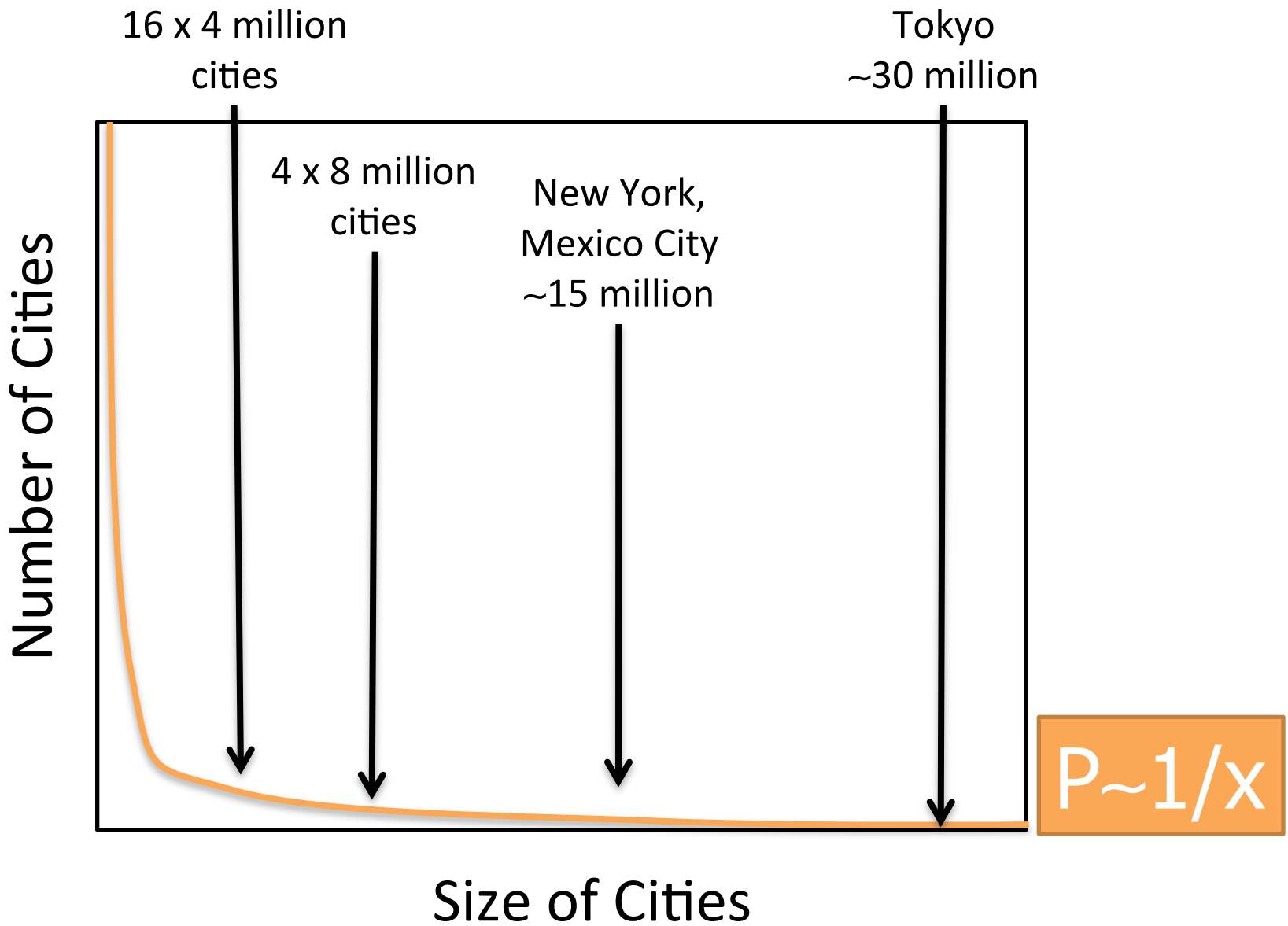


Vilfredo Pareto (1848-1923)

Rich and Poor in America



This plot of household wealth in the United States, taken from 1998 census figures, clearly shows a distribution of rich and poor forming a Pareto curve. The highest percentage of households fall at the lower levels of wealth, but at the higher end, the curve drops off relatively slowly, displaying Pareto's "fat-tailed" pattern.



NO OUTLIERS IN A RANDOM SOCIETY

$$P(k) = e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!}$$

- The most connected individual has degree $k_{\max} \sim 1,185$
- The least connected individual has degree $k_{\min} \sim 816$

The probability to find an individual with degree $k > 2,000$ is 10^{-27} . Hence the chance of finding an individual with 2,000 acquaintances is so tiny that such nodes are virtually inexistent in a random society.

- a random society would consist of mainly average individuals, with everyone with roughly the same number of friends.
- It would lack outliers, individuals that are either highly popular or recluse.

After Bill enters the arena the average wealth of the public ~ \$1,000,000

~ \$100 billion

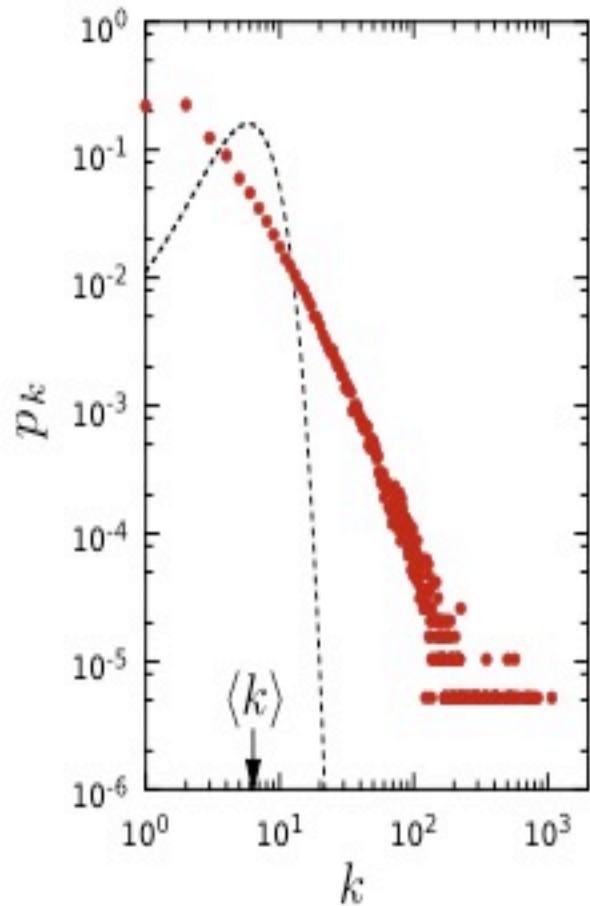


10^5 people, 10^5 \$ average wealth per capita

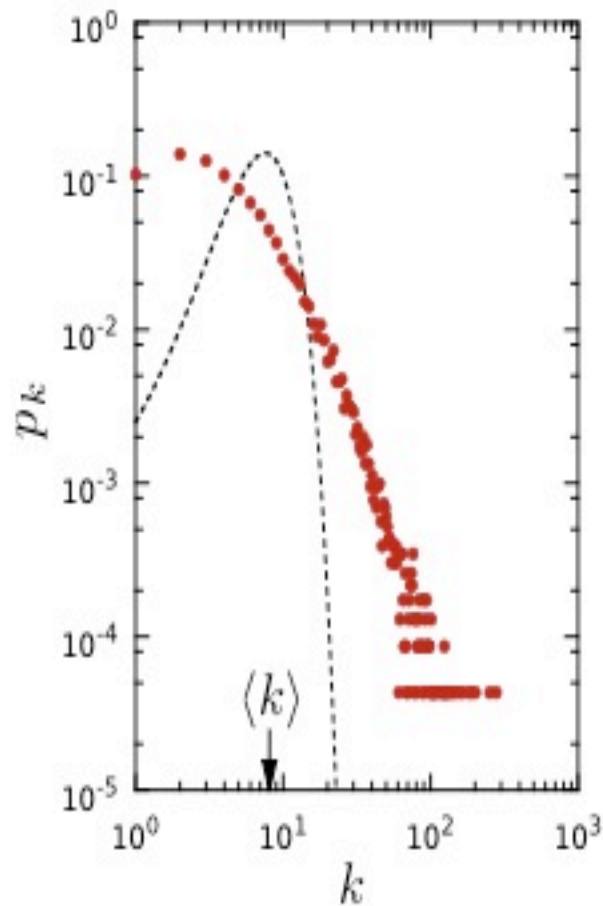
FACING REALITY: Degree distribution of real networks

$$P(k) = e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!}$$

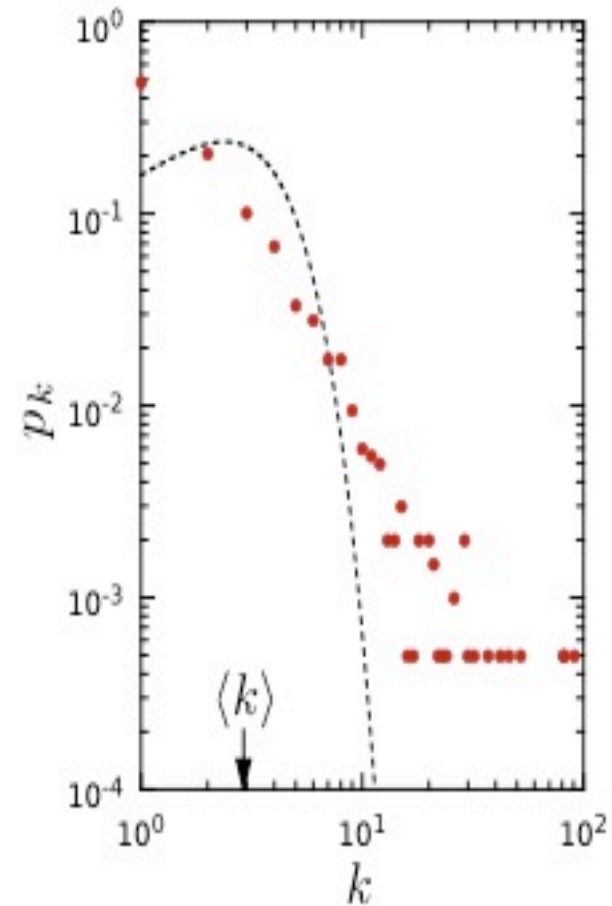
Internet



Science Collaboration



Protein Interactions



UNIVERSALITY

Network	Size	$\langle k \rangle$	κ	γ_{out}	γ_{in}
WWW	325 729	4.51	900	2.45	2.1
WWW	4×10^7	7		2.38	2.1
WWW	2×10^8	7.5	4000	2.72	2.1
WWW, site	260 000				1.94
Internet, domain*	3015–4389	3.42–3.76	30–40	2.1–2.2	2.1–2.2
Internet, router*	3888	2.57	30	2.48	2.48
Internet, router*	150 000	2.66	60	2.4	2.4
Movie actors*	212 250	28.78	900	2.3	2.3
Co-authors, SPIRES*	56 627	173	1100	1.2	1.2
Co-authors, neuro.*	209 293	11.54	400	2.1	2.1
Co-authors, math.*	70 975	3.9	120	2.5	2.5
Sexual contacts*	2810			3.4	3.4
Metabolic, <i>E. coli</i>	778	7.4	110	2.2	2.2
Protein, <i>S. cerev.</i> *	1870	2.39		2.4	2.4
Ythan estuary*	134	8.7	35	1.05	1.05
Silwood Park*	154	4.75	27	1.13	1.13
Citation	783 339	8.57			3
Phone call	53×10^6	3.16		2.1	2.1
Words, co-occurrence*	460 902	70.13		2.7	2.7
Words, synonyms*	22 311	13.48		2.8	2.8

Networks:

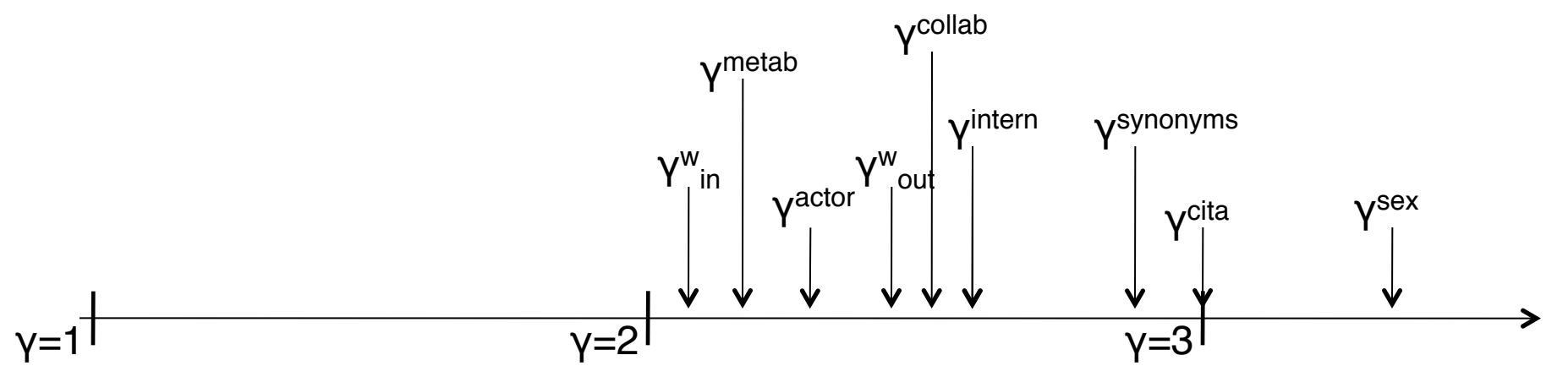
The exponents vary from system to system.

Most are between 2 and 3

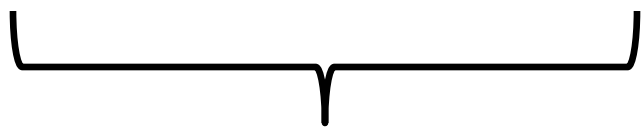
Universality:

the emergence of common features across different networks. Like the scale-free property.

VARIANCE DIVERGES!



$\langle k^2 \rangle$ diverges		$\langle k^2 \rangle$ finite
Regime full of anomalies...	The scale-free behavior is relevant	Behaves like a random network



Why are most exponents in this regime?

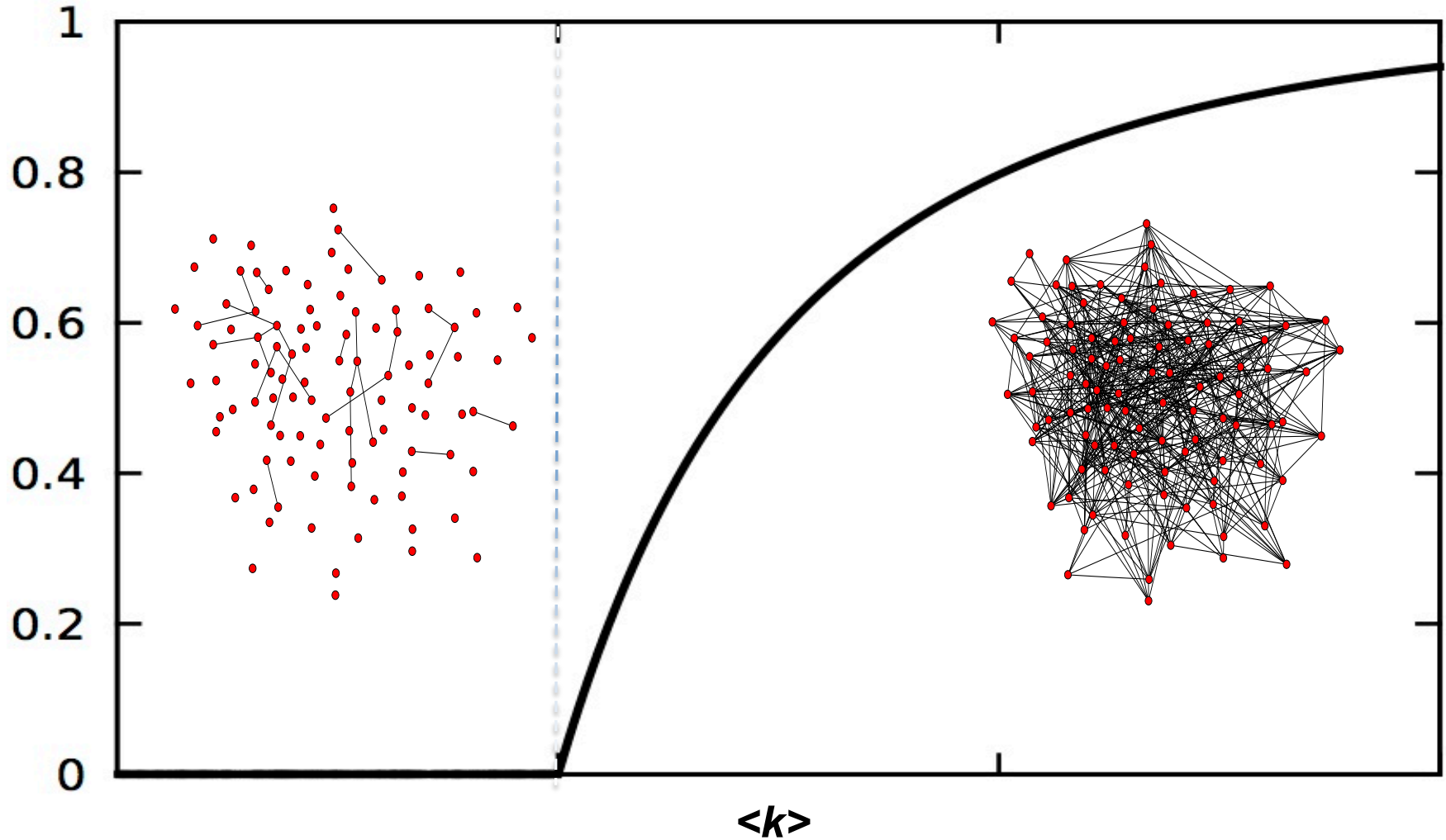
The evolution of a random network

EVOLUTION OF A RANDOM NETWORK

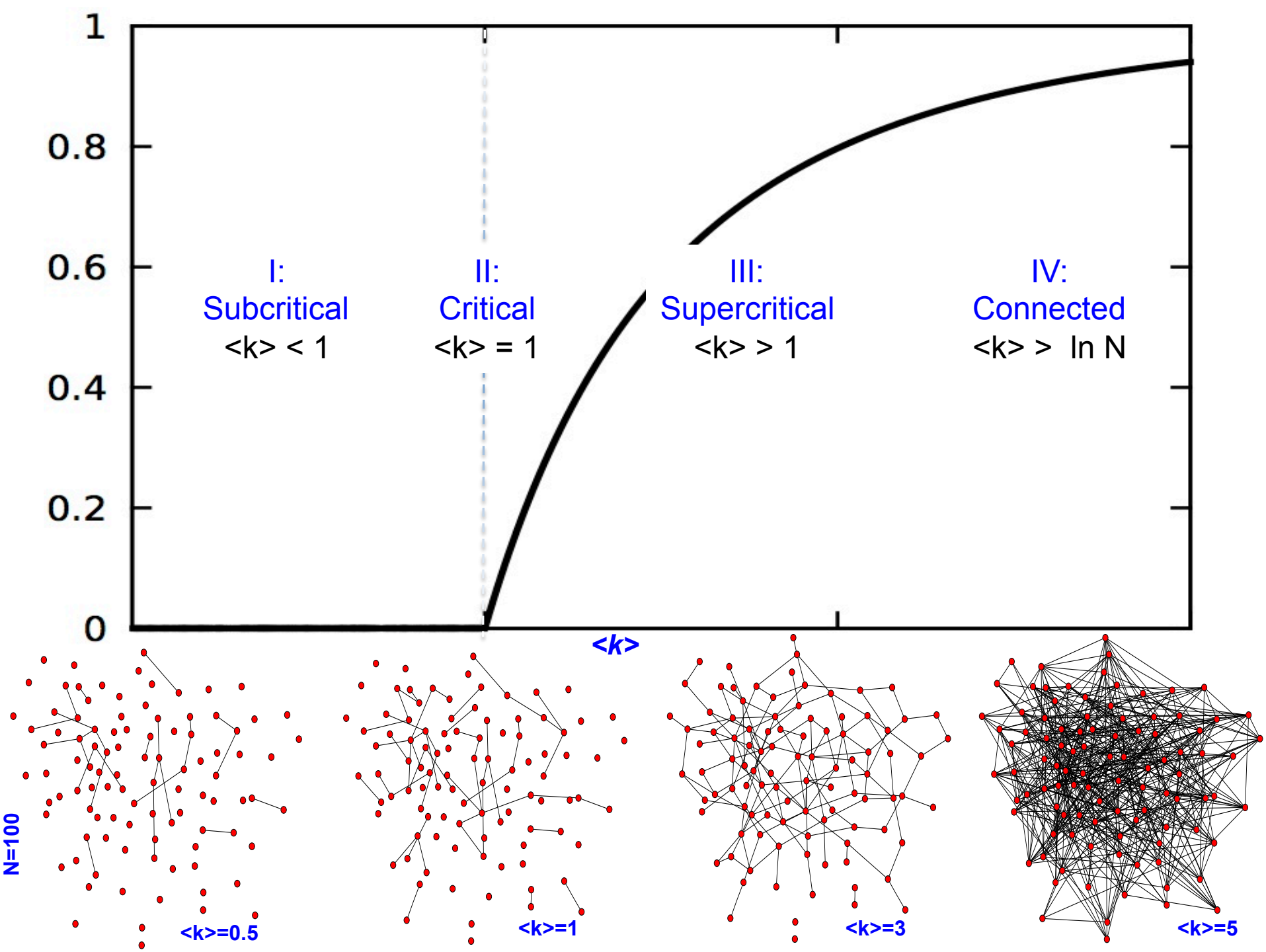
disconnected nodes



NETWORK.

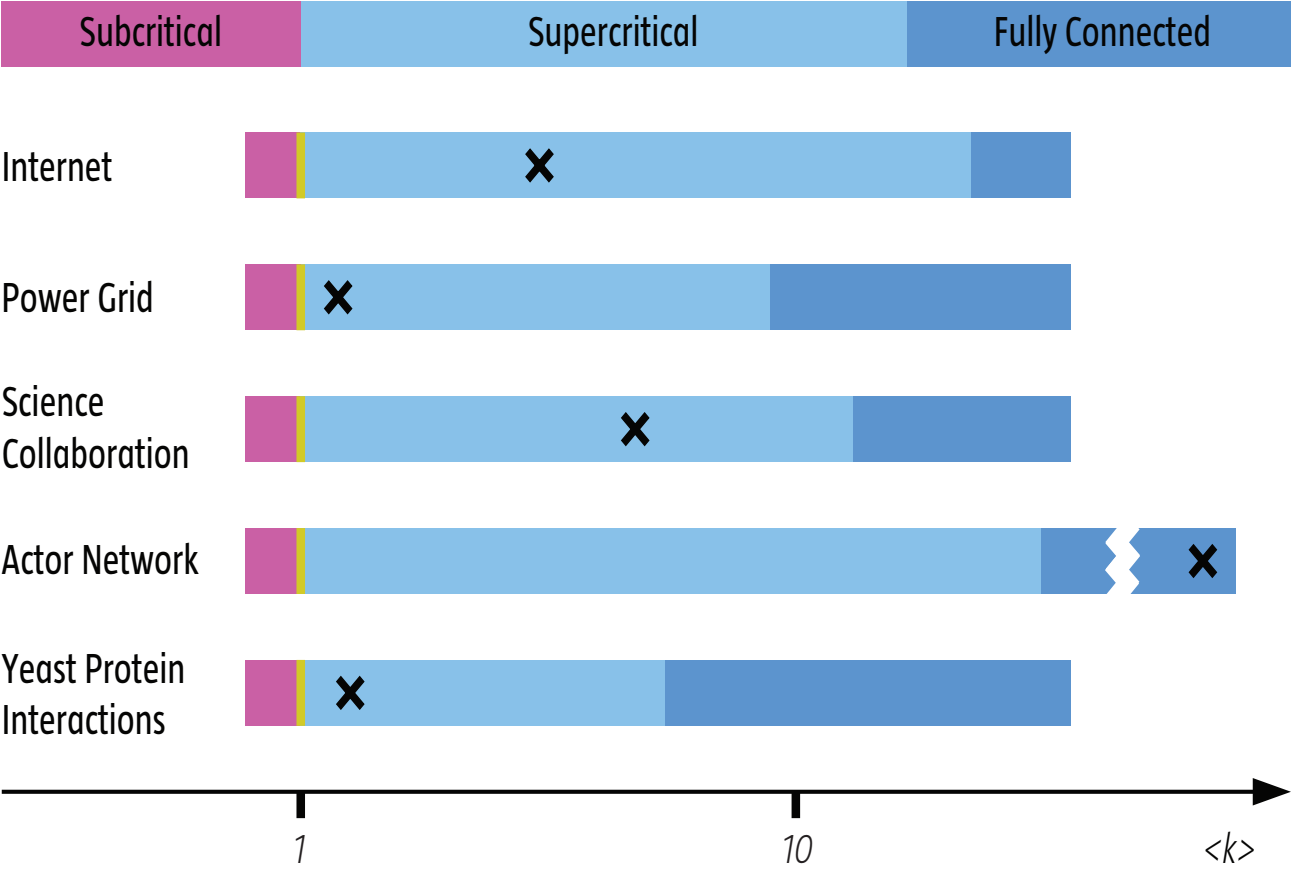


How does this transition happen?



Real networks are supercritical

Section 7

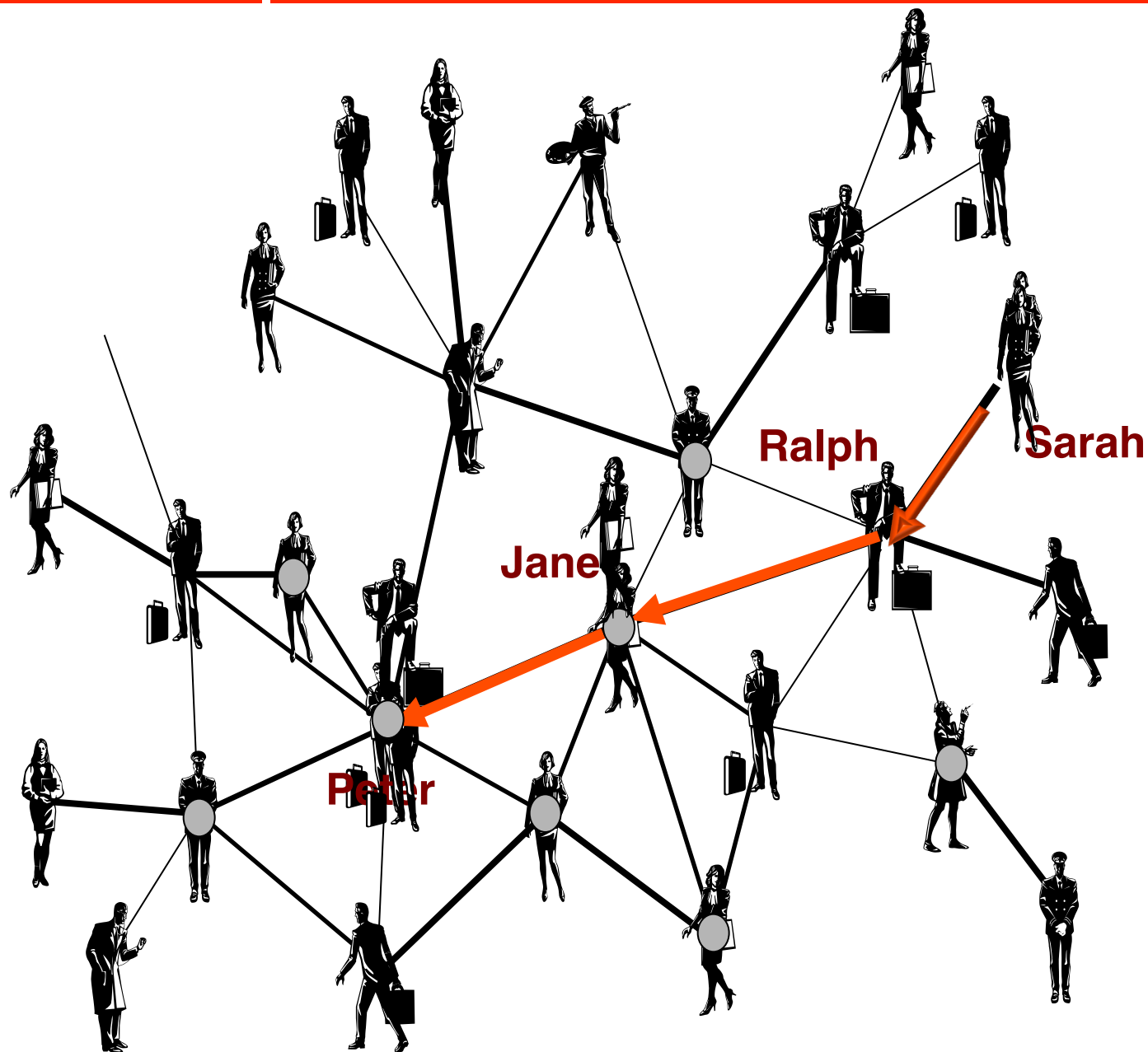


Network	N	L	$\langle k \rangle$	$\ln N$
Internet	192,244	609,066	6.34	12.17
Power Grid	4,941	6,594	2.67	8.51
Science Collaboration	23,133	186,936	8.08	10.04
Actor Network	212,250	3,054,278	28.78	12.27
Yeast Protein Interactions	2,018	2,930	2.90	7.61

Small world property

SIX DEGREES

small worlds



*Frigyes Karinthy, 1929
Stanley Milgram, 1967*



Frigyes Karinthy (1887-1938)
Hungarian Writer

1929: *Minden más*  *leppen van* (Everything is Different)
Láncszemek (Chains)

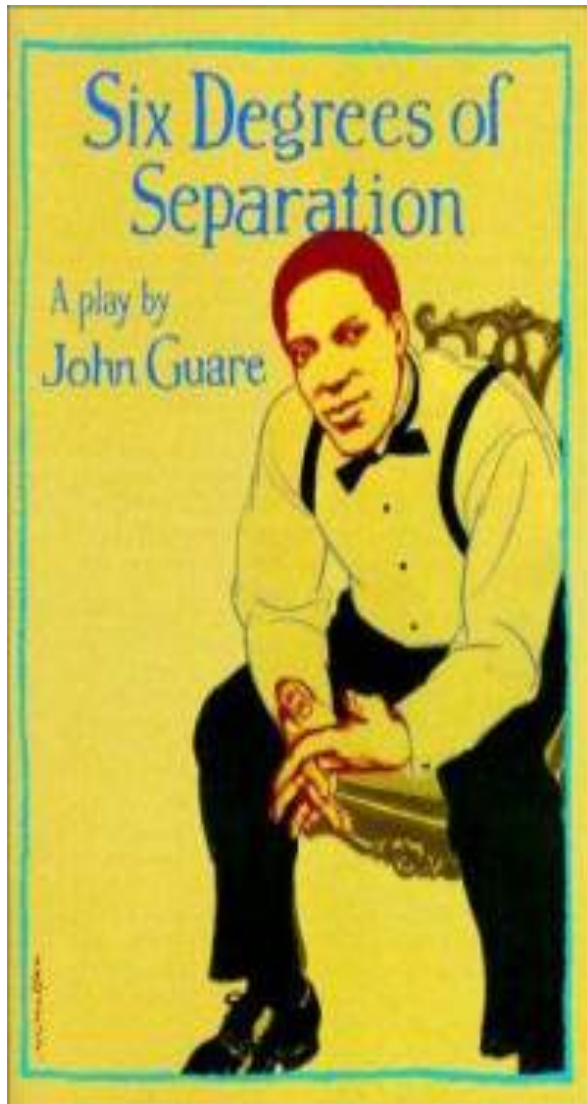
“Look, Selma Lagerlöf just won the Nobel Prize for Literature, thus she is bound to know King Gustav of Sweden, after all he is the one who handed her the Prize, as required by tradition. King Gustav, to be sure, is a passionate tennis player, who always participates in international tournaments. He is known to have played Mr. Kehrling, whom he must therefore know for sure, and as it happens I myself know Mr. Kehrling quite well.”

"The worker knows the manager in the shop, who knows Ford; Ford is on friendly terms with the general director of Hearst Publications, who last year became good friends with Arpad Pasztor, someone I not only know, but to the best of my knowledge a good friend of mine. So I could easily ask him to send a telegram via the general director telling Ford that he should talk to the manager and have the worker in the shop quickly hammer together a car for me, as I happen to need one."



HOW TO TAKE PART IN THIS STUDY

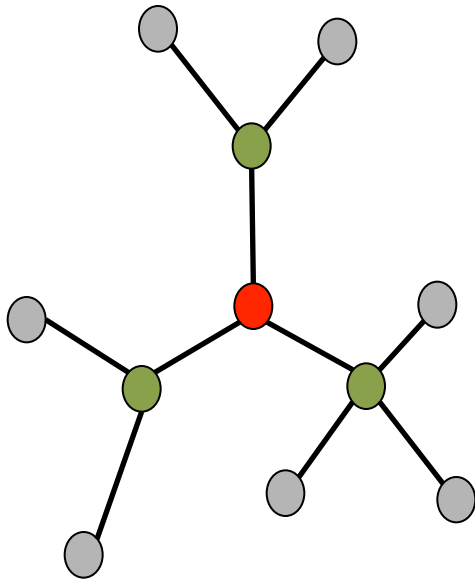
1. ADD YOUR NAME TO THE ROSTER AT THE BOTTOM OF THIS SHEET, so that the next person who receives this letter will know who it came from.
2. DETACH ONE POSTCARD. FILL IT AND RETURN IT TO HARVARD UNIVERSITY. No stamp is needed. The postcard is very important. It allows us to keep track of the progress of the folder as it moves toward the target person.
3. IF YOU KNOW THE TARGET PERSON ON A PERSONAL BASIS, MAIL THIS FOLDER DIRECTLY TO HIM (HER). Do this only if you have previously met the target person and know each other on a first name basis.
4. IF YOU DO NOT KNOW THE TARGET PERSON ON A PERSONAL BASIS, DO NOT TRY TO CONTACT HIM DIRECTLY. INSTEAD, MAIL THIS FOLDER (POST CARDS AND ALL) TO A PERSONAL ACQUAINTANCE WHO IS MORE LIKELY THAN YOU TO KNOW THE TARGET PERSON. You may send the folder to a friend, relative or acquaintance, but it must be someone you know on a first name basis.



"Everybody on this planet is separated by only six other people. Six degrees of separation. Between us and everybody else on this planet. The president of the United States. A gondolier in Venice.... It's not just the big names. It's anyone. A native in a rain forest. A Tierra del Fuegan. An Eskimo. I am bound to everyone on this planet by a trail of six people. It's a profound thought. How every person is a new door, opening up into other worlds."

DISTANCES IN RANDOM GRAPHS

Random graphs tend to have a tree-like topology with almost constant node degrees.



- nr. of first neighbors:

$$N_1 \cong \langle k \rangle$$

- nr. of second neighbors:

$$N_2 \cong \langle k \rangle^2$$

- nr. of neighbours at distance d:

$$N_d \cong \langle k \rangle^d$$

- estimate maximum distance:

$$1 + \sum_{l=1}^{l_{\max}} \langle k \rangle^l = N \quad \Rightarrow \quad l_{\max} = \frac{\log N}{\log \langle k \rangle}$$

DISTANCES IN RANDOM GRAPHS

compare with real data

$$l_{\max} = \frac{\log N}{\log \langle k \rangle}$$

Network	Size	(k)	l	l _{rand}	C	C _{rand}	Reference	Nr
www, site level, undir	153127	35.21	3.1	3.35	0.1078	0.00023	Adamic, 1999	1
Internet, domain level	3015-6209	3.52-4.11	3.7-3.76	6.36-6.18	0.18-0.3	0.001	Yook e al., 2001a, Pastor-Satorras et al., 2001	2
Movie actors	225226	61	3.65	2.99	0.79	0.00027	Watts and Strogatz,1998	3
LANL co-authorship	52909	9.7	5.9	4.79	0.43	1.8 x 10 ⁻⁴	Newman, 2001a, 2001b, 2001c	4
MEDLINE eo-authorship	1520251	18.1	4.6	4.91	0.066	1.1 x 10 ⁻⁵	Newman, 2001a, 2001b, 2001c	5
SPIRES co-authorship	56627	173	4.0	2.12	0.726	0.003	Newman, 2001a, 2001b, 2001c	6
NCSTRL co-authorship	11994	3.59	9.7	7.34	0.496	3 x 10 ⁻⁴	Newman, 2001a, 2001b, 2001c	7
Math. co-authorship	70975	3.9	9.5	8.2	0.59	5.4 x 10 ⁻⁵	Barabasi et al, 2001	8
Neurosci. co-authorship	209293	11.5	6	5.01	0.76	5.5 x 10 ⁻⁵	Barabasi et al, 2001	9
E. coli, sustrate graph	282	7.35	2.9	3.04	0.32	0.026	Wagner and Fell, 2000	10
E. coli, reaction graph	315	28.3	2.62	1.98	0.59	0.09	Wagner and Fell, 2000	11
Ythan estuary food web	134	8.7	2.43	2.26	0.22	0.06	Montoya and Sole, 2000	12
Silwood Park food web	154	4.75	3.40	3.23	0.15	0.03	Montoya and Sole, 2000	13
Words, co-occurrence	460902	70.13	2.67	3.03	0.437	0.0001	Ferrer i Cancho and Sole, 2001	14
Words, synonyms	22311	13.48	4.5	3.84	0.7	0.0006	Yook et al. 2001b	15
Power grid	4941	2.67	18.7	12.4	0.08	0.005	Watts and Strogatz, 1998	16
C.Elegans	282	14	2.65	2.25	0.28	0.05	Watts and Strogatz, 1998	17

Given the huge differences in scope, size, and average degree, the agreement is excellent.

CLUSTERING COEFFICIENT

$$C_i \equiv \frac{2n_i}{k_i(k_i - 1)}$$

Since edges are independent and have the same probability p ,

$$n_i \cong p \frac{k_i(k_i - 1)}{2} \quad \Rightarrow \quad C \cong p = \frac{\langle k \rangle}{N}$$

The clustering coefficient of random graphs is small.

For fixed degree C decreases with the system size N .

CLUSTERING IN RANDOM GRAPHS

compare with real data

Network	Size	(k)	\bar{l}	\bar{l}_{rand}	C	C_{rand}	Reference	Nr
www, site level, undir	153127	35.21	3.1	3.35	0.1078	0.00023	Adamic, 1999	1
Internet, domain level	3015-6209	3.52-4.11	3.7-3.76	6.36-6.18	0.18-0.3	0.001	Yook et al., 2001a, Pastor-Satorras et al., 2001	2
Movie actors	225226	61	3.65	2.99	0.79	0.00027	Watts and Strogatz, 1998	3
LANL co-authorship	52909	9.7	5.9	4.79	0.43	1.8×10^{-4}	Newman, 2001a, 2001b, 2001c	4
MEDLINE eo-authorship	1520251	18.1	4.6	4.91	0.066	1.1×10^{-5}	Newman, 2001a, 2001b, 2001c	5
SPIRES co-authorship	56627	173	4.0	2.12	0.726	0.003	Newman, 2001a, 2001b, 2001c	6
NCSTRL co-authorship	11994	3.59	9.7	7.34	0.496	3×10^{-4}	Newman, 2001a, 2001b, 2001c	7
Math. co-authorship	70975	3.9	9.5	8.2	0.59	5.4×10^{-5}	Barabasi et al, 2001	8
Neurosci. co-authorship	209293	11.5	6	5.01	0.76	5.5×10^{-5}	Barabasi et al, 2001	9
E. coli, sustrate graph	282	7.35	2.9	3.04	0.32	0.026	Wagner and Fell, 2000	10
E. coli, reaction graph	315	28.3	2.62	1.98	0.59	0.09	Wagner and Fell, 2000	11
Ythan estuary food web	134	8.7	2.43	2.26	0.22	0.06	Montoya and Sole, 2000	12
Silwood Park food web	154	4.75	3.40	3.23	0.15	0.03	Montoya and Sole, 2000	13
Words, co-occurrence	460902	70.13	2.67	3.03	0.437	0.0001	Ferrer i Cancho and Sole, 2001	14
Words, synonyms	22311	13.48	4.5	3.84	0.7	0.0006	Yook et al. 2001b	15
Power grid	4941	2.67	18.7	12.4	0.08	0.005	Watts and Strogatz, 1998	16
C.Elegans	282	14	2.65	2.25	0.28	0.05	Watts and Strogatz, 1998	17

- **Degree distribution**

Binomial, Poisson (exponential tails)

- **Clustering coefficient**

Vanishing for large network sizes

- **Average distance among nodes**

Logarithmically small

**Are real networks like
random graphs?
NO!**

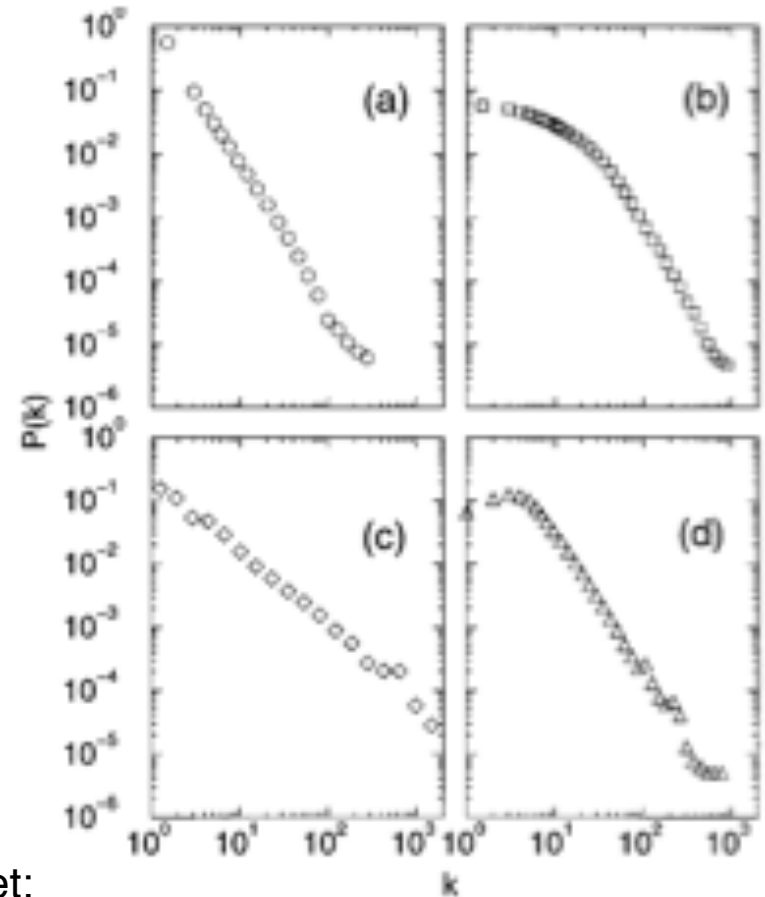
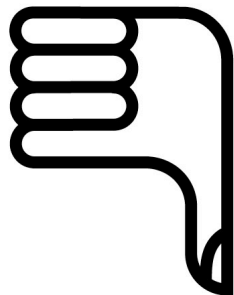
THE DEGREE DISTRIBUTION

Prediction:

$$P_{rand}(k) \cong C_{N-1}^k p^k (1-p)^{N-1-k}$$

Data:

$$P(k) \approx k^{-\gamma}$$



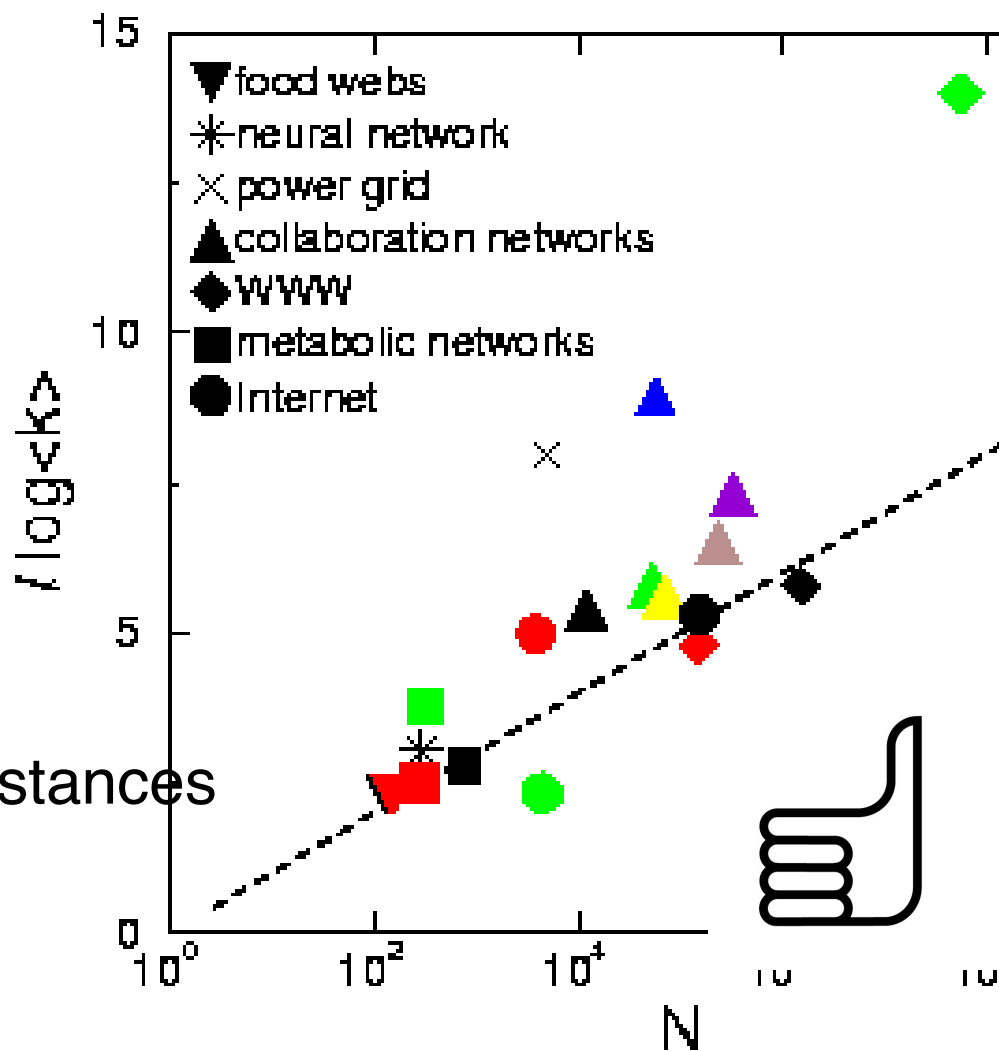
- (a) Internet;
- (b) Movie Actors;
- (c) Coauthorship, high energy physics;
- (d) Coauthorship, neuroscience

PATH LENGTHS IN REAL NETWORKS

Prediction:

$$l_{rand} = \frac{\log N}{\log \langle k \rangle}$$

Data:



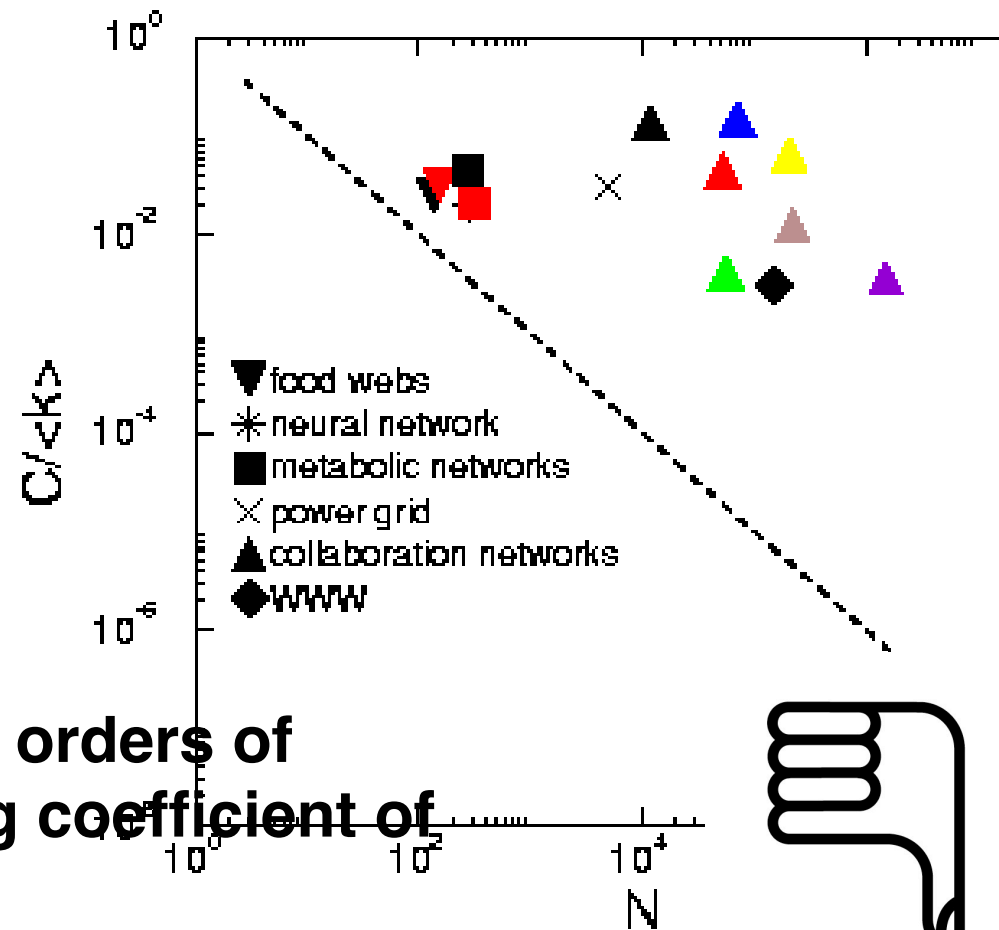
Real networks have short distances like random graphs.

CLUSTERING COEFFICIENT

Prediction:

$$C_{rand} = \frac{\langle k \rangle}{N}$$

Data:

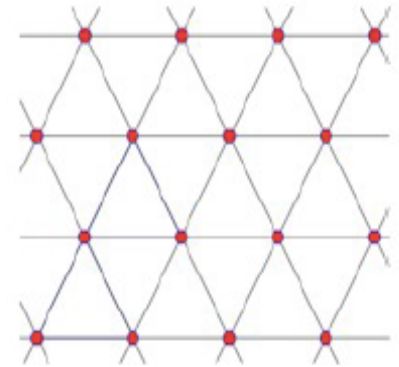
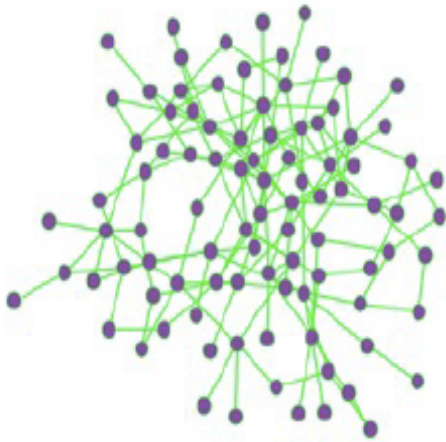


C_{rand} underestimates with orders of magnitudes the clustering coefficient of real networks.

A solid red horizontal bar at the top of the slide, divided into two segments by a thin white vertical line.

The small-world model

Real networks are between random networks and lattices



Real networks are
somewhere here

Watts-Strogatz model



Duncan Watts



Steve Strogatz

NATURE | VOL 393 | 4 JUNE 1998

Collective dynamics of 'small-world' networks

Duncan J. Watts* & Steven H. Strogatz

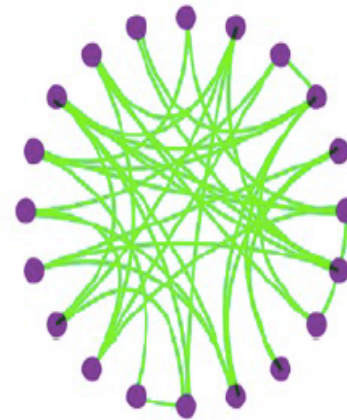
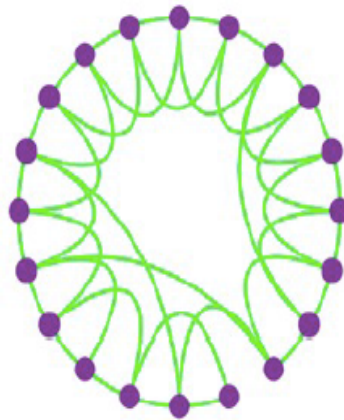
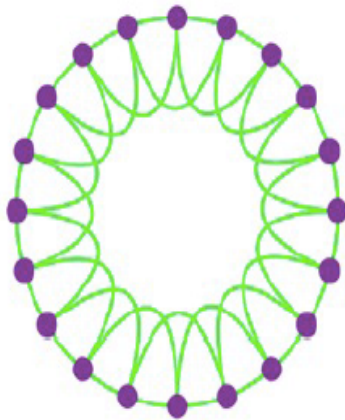
*Department of Theoretical and Applied Mechanics, Kimball Hall,
Cornell University, Ithaca, New York 14853, USA*

Networks of coupled dynamical systems have been used to model biological oscillators¹⁻⁴, Josephson junction arrays^{5,6}, excitable media⁷, neural networks⁸⁻¹⁰, spatial games¹¹, genetic control networks¹² and many other self-organizing systems. Ordinarily, the connection topology is assumed to be either completely regular or completely random. But many biological, technological and social networks lie somewhere between these two extremes.

REGULAR

SMALL-WORLD

RANDOM



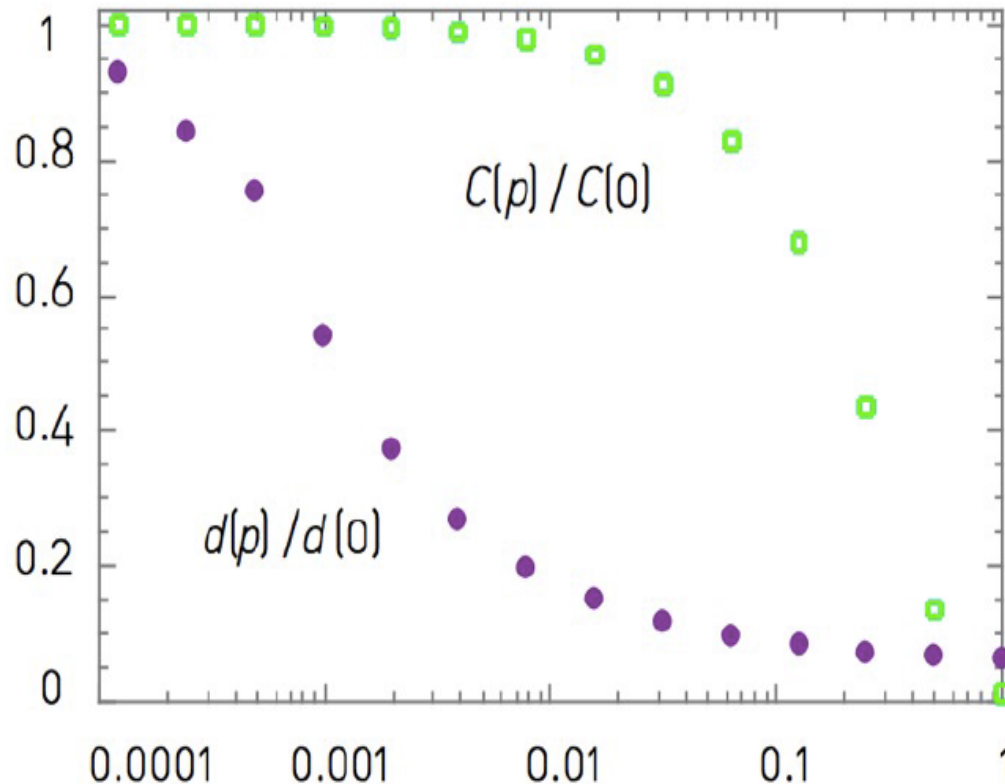
$p=0$



$p=1$

Increasing randomness

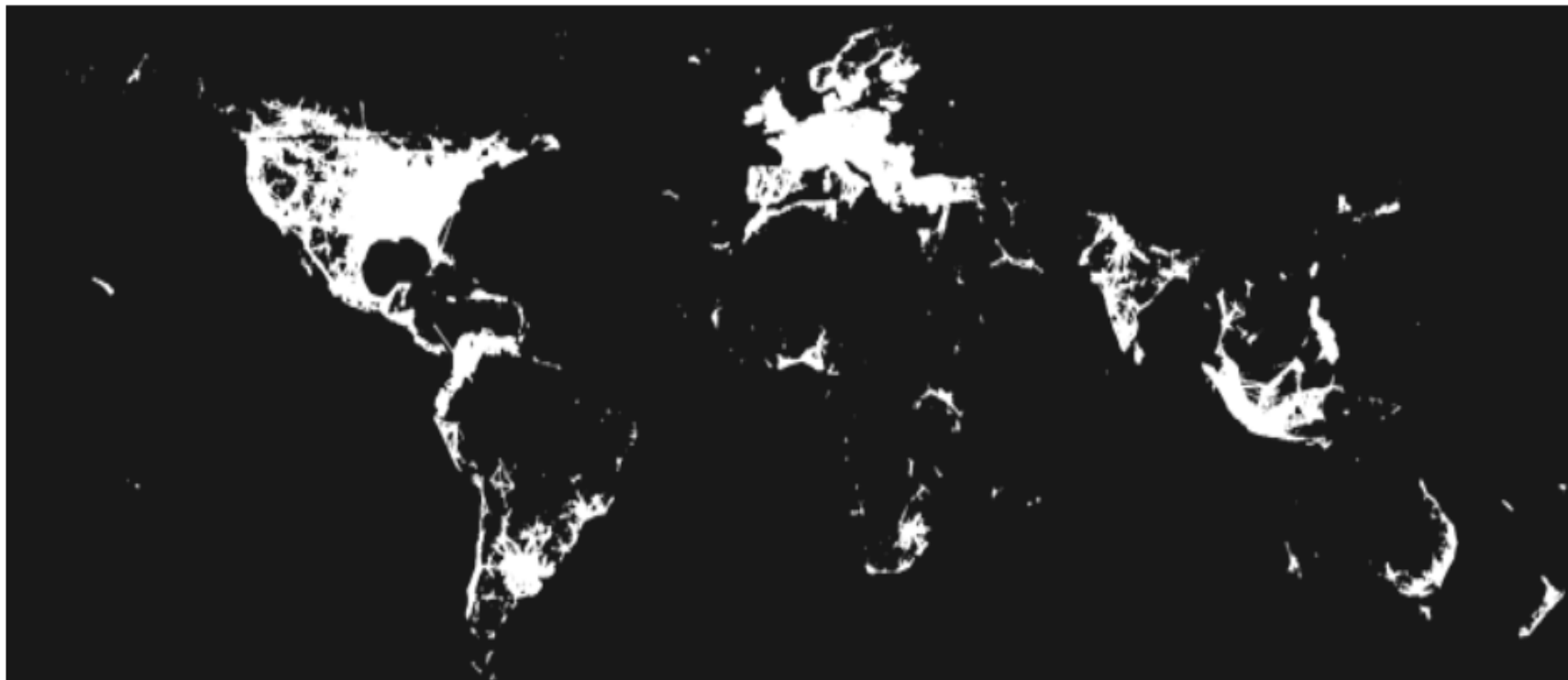
Average path length vs. clustering coefficient



p

The Watts Strogatz Model:

It takes a lot of randomness to ruin the clustering, but a very small amount to overcome locality





Hubs represent the most striking difference between a random and a scale-free network. Their emergence in many real systems raises several fundamental questions:

- Why does the random network model of Erdős and Rényi fail to reproduce the hubs and the power laws observed in many real networks?
- Why do so different systems as the WWW or the cell converge to a similar scale-free architecture?

A solid red horizontal bar at the top of the slide, divided into two segments by a thin white vertical line.

Growth and preferential attachment

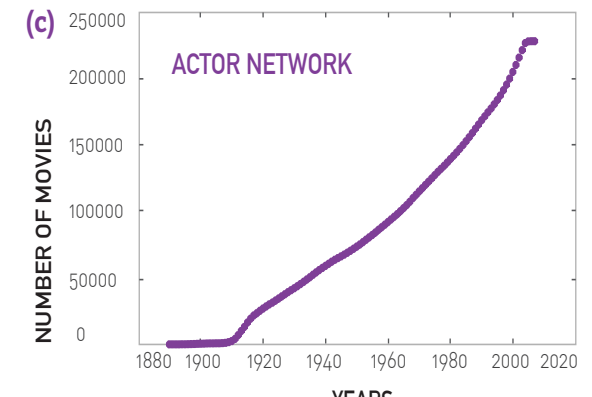
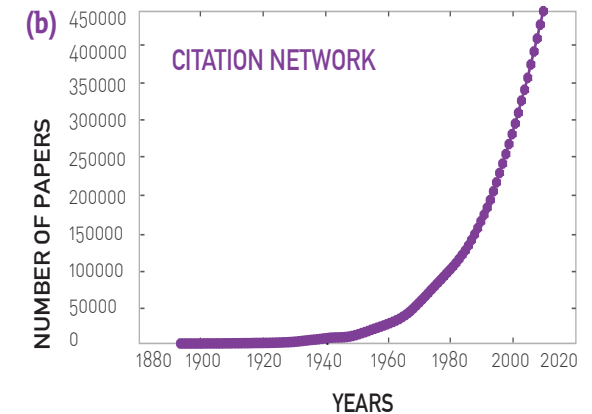
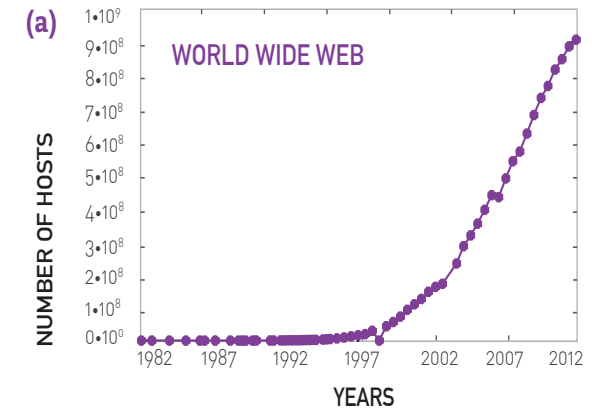
BA MODEL: Growth

ER model:

the number of nodes, N , is fixed (static models)

**networks expand through the addition
of new nodes**

Barabási & Albert, *Science* **286**, 509 (1999)



ER model: links are added randomly to the network

New nodes prefer to connect to the more connected nodes

Growth and Preferential Attachment

The random network model differs from real networks in two important characteristics:

Growth: While the random network model assumes that the number of nodes is fixed (time invariant), real networks are the result of a growth process that continuously increases.

Preferential Attachment: While nodes in random networks randomly choose their interaction partner, in real networks new nodes prefer to link to the more connected nodes.

A solid red horizontal bar at the top of the slide, divided into two segments by a thin white vertical line.

The Barabási-Albert model

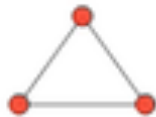
Origin of SF networks: Growth and preferential attachment

(1) Networks continuously expand by the addition of new nodes

WWW : addition of new documents

(2) New nodes prefer to link to highly connected nodes.

WWW : linking to well known sites

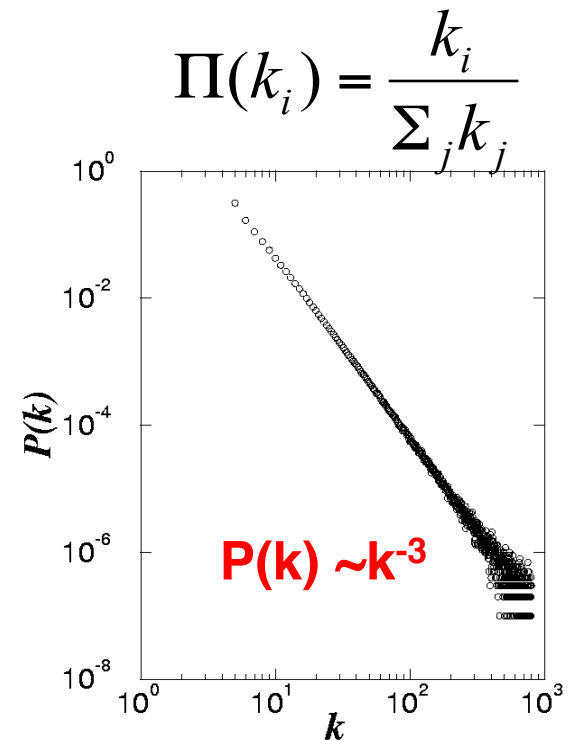


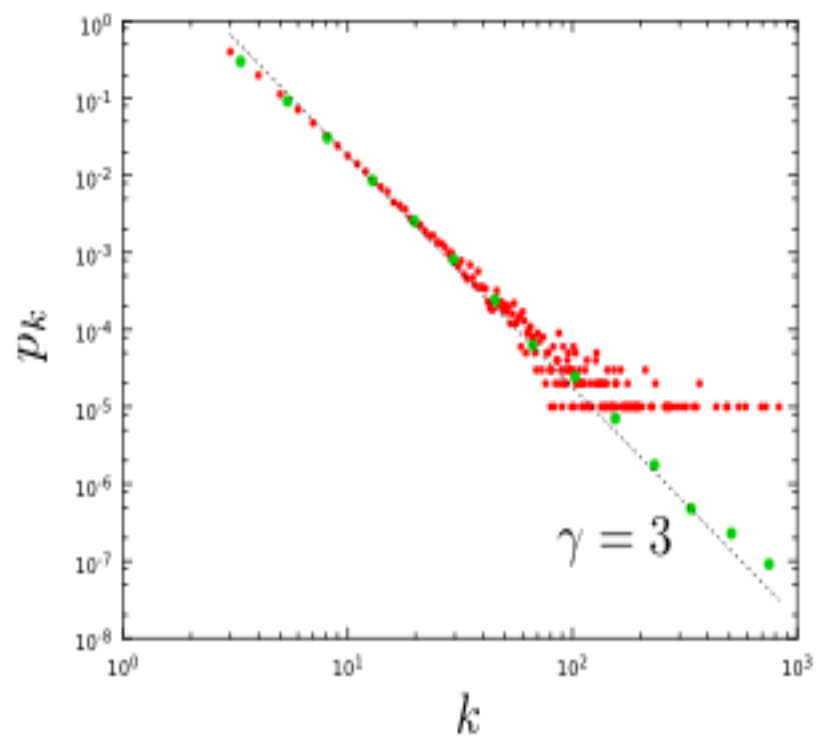
GROWTH:

add a new node with m links

PREFERENTIAL ATTACHMENT:

the probability that a node connects to a node with k links is proportional to k .







György Pólya
PÓLYA PROCESS
MATHEMATICIAN



George Kinsley Zipf
WEALTH DISTRIBUTION
ECONOMIST



Herbert Alexander Simon
MASTER EQUATION
POLITICAL SCIENTIST



Robert Merton
MATTHEW EFFECT
SOCIOLOGIST



Albert-László Barabási & Réka Albert
PREFERENTIAL ATTACHMENT
NETWORK SCIENTISTS



George Udny Yule
YULE PROCESS
STATISTICIAN



Robert Gibrat
PROPORTIONAL GROWTH
ECONOMIST



Derek de Solla Price
CUMULATIVE ADVANTAGE
PHYSICIST

MILESTONES

PUBLICATION
DATE

1923 1925 1931 1935 1941 1945 1950 1955 1960 1968 1970 1976 1980 1985 1990 1995 1999 2000 2005 2010

György Pólya (1887-1985)
Preferential attachment made its first appearance in 1923 in the celebrated urn model of the Hungarian mathematician György Pólya [2]. Hence, in mathematics preferential attachment is often called a **Pólya process**.

George Udny Yule (1871-1951)
used preferential attachment to explain the power-law distribution of the number of species per genus of flowering plants [3]. Hence, in statistics preferential attachment is often called a **Yule process**.

Robert Gibrat (1904-1980)
proposed that the size and the growth rate of a firm are independent. Hence, larger firms grow faster [4]. Called **proportional growth**, this is a form of preferential attachment.

George Kinsley Zipf (1902-1950)
used preferential attachment to explain the fat tailed distribution of wealth in the society [5].

Herbert Alexander Simon (1916-2001)
used preferential attachment to explain the fat-tailed nature of the distributions describing city sizes, word frequencies, or the number of papers published by scientists [6].

Derek de Solla Price (1922-1983)
used preferential attachment to explain the citation statistics of scientific publications, referring to it as **cumulative advantage** [7].

Robert Merton (1910-2003)
In sociology preferential attachment is often called the **Matthew effect**, named by Merton [8] after a passage in the Gospel of Matthew.

Barabási (1967) & **Albert** (1972)
introduce the term **preferential attachment** in the context of networks [1] to explain the origin of their power-law degree distribution.