

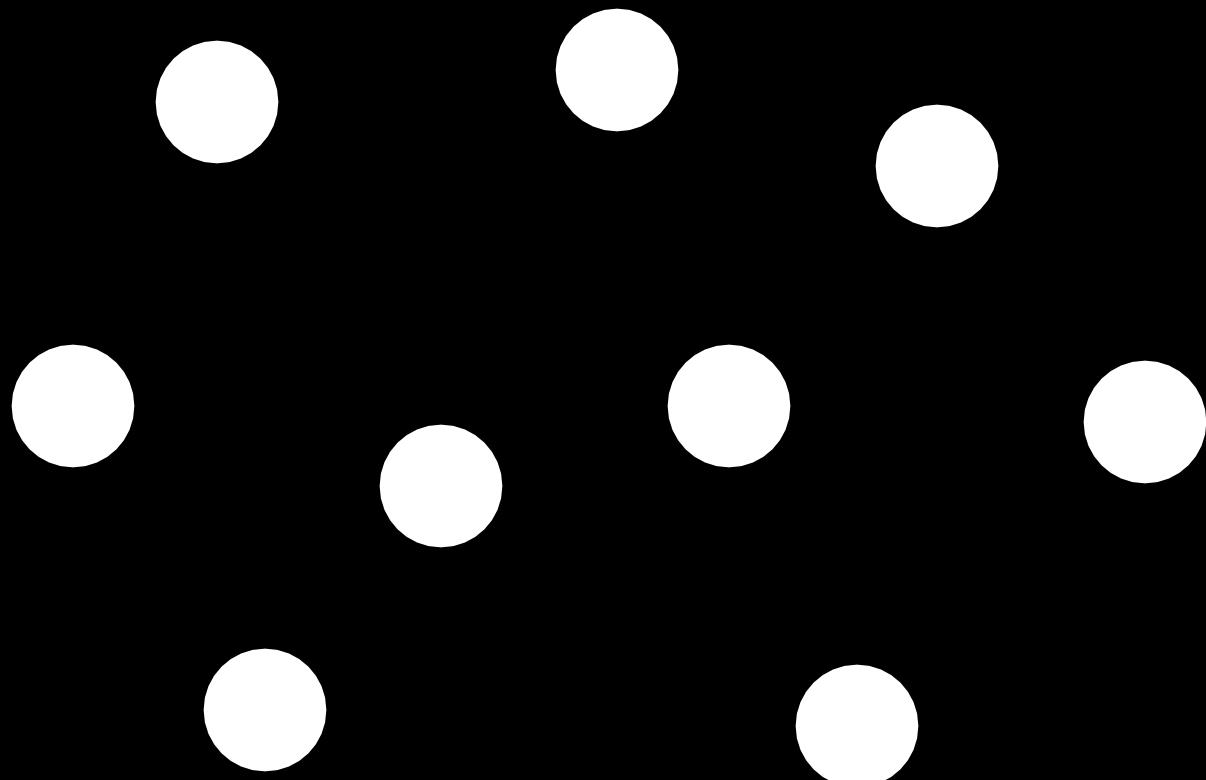
# CS224W: Analysis of Networks

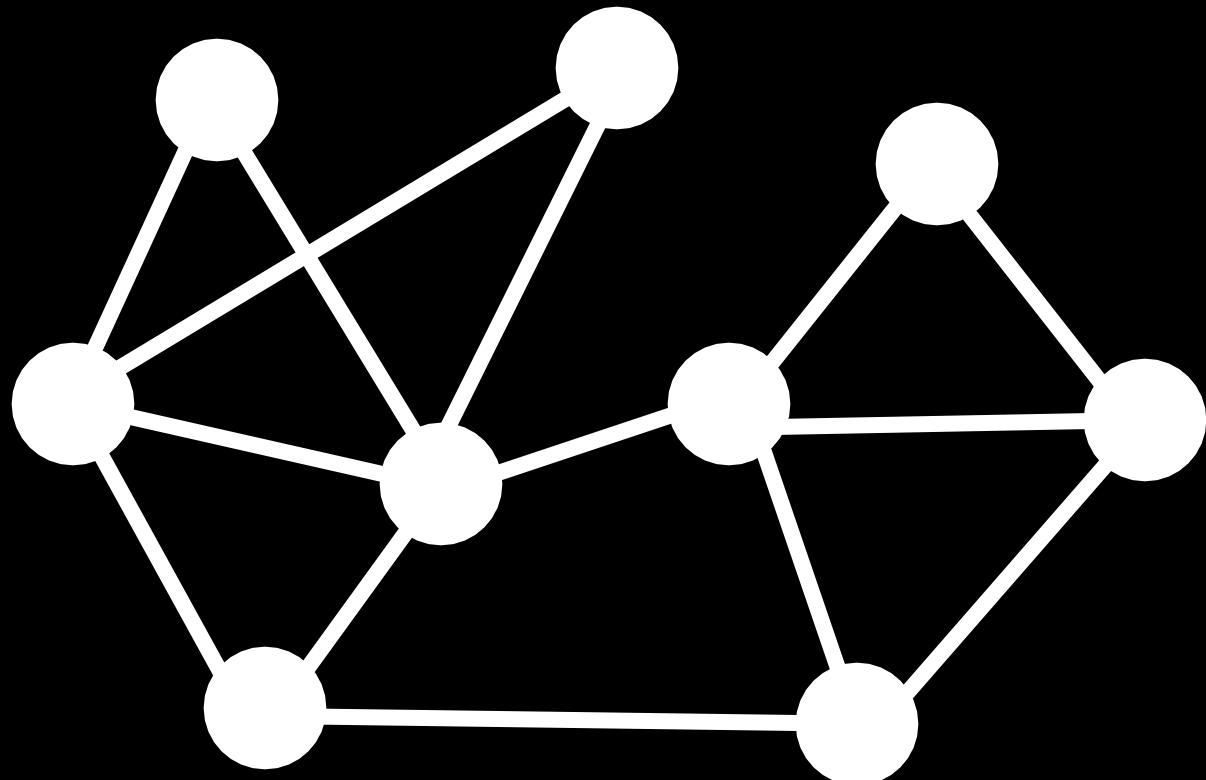
CS224W: Analysis of Network  
Jure Leskovec, Stanford University  
<http://cs224w.stanford.edu>



# Why Networks?

Networks are a general  
language for describing  
complex systems



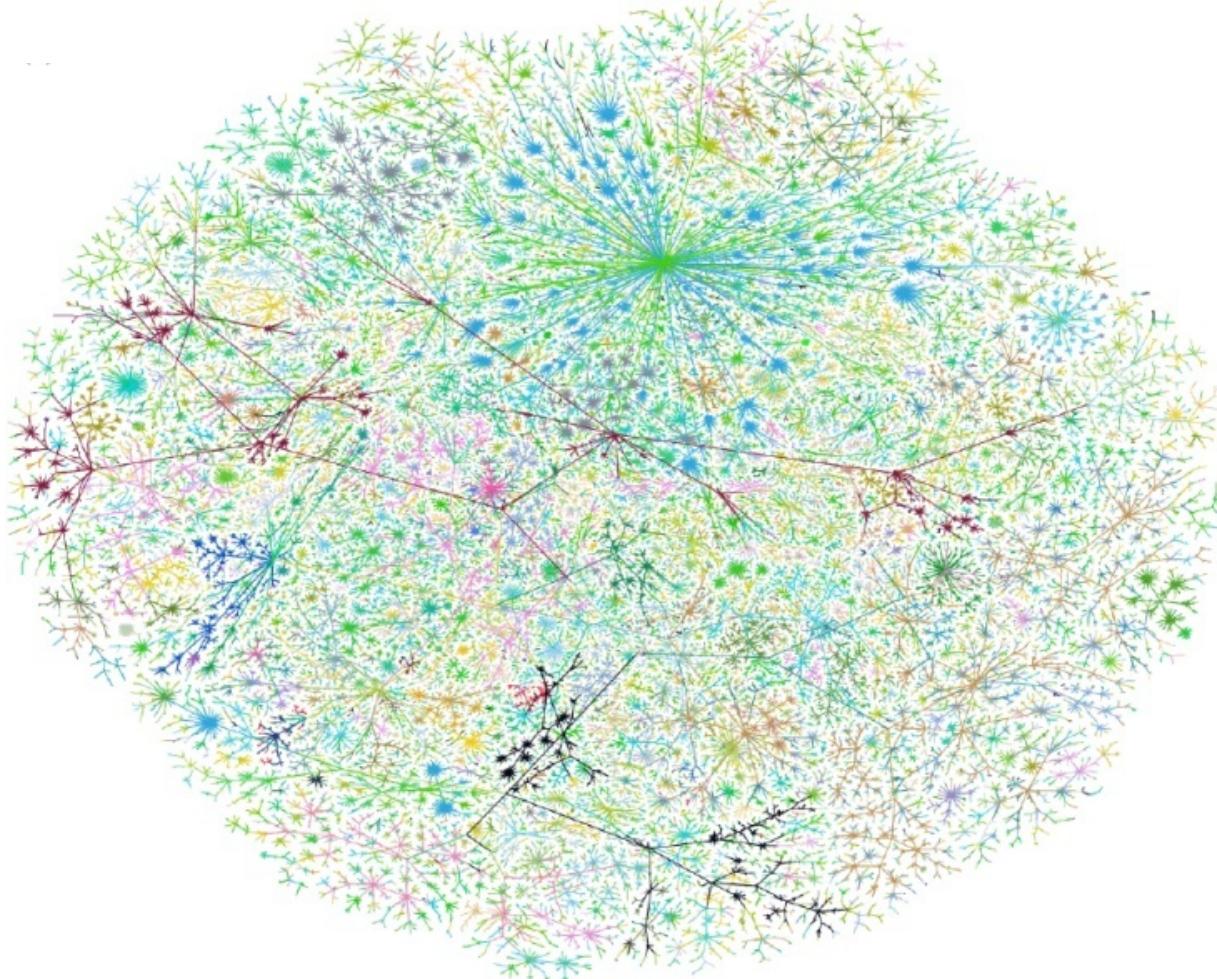


# Interactions!

# Networks & Complex Systems

- **Complex systems are all around us:**
  - **Society** is a collection of six billion individuals
  - **Communication systems** link electronic devices
  - **Information** and **knowledge** is organized and linked
  - Interactions between thousands of **genes/proteins** regulate life
  - Our **thoughts** are hidden in the connections between billions of neurons in our brain

**What do these systems have in common?**  
**How can we represent them?**



# The Network!

# Networks!!

Behind many systems there is an intricate wiring diagram, **a network**, that defines the **interactions** between the components

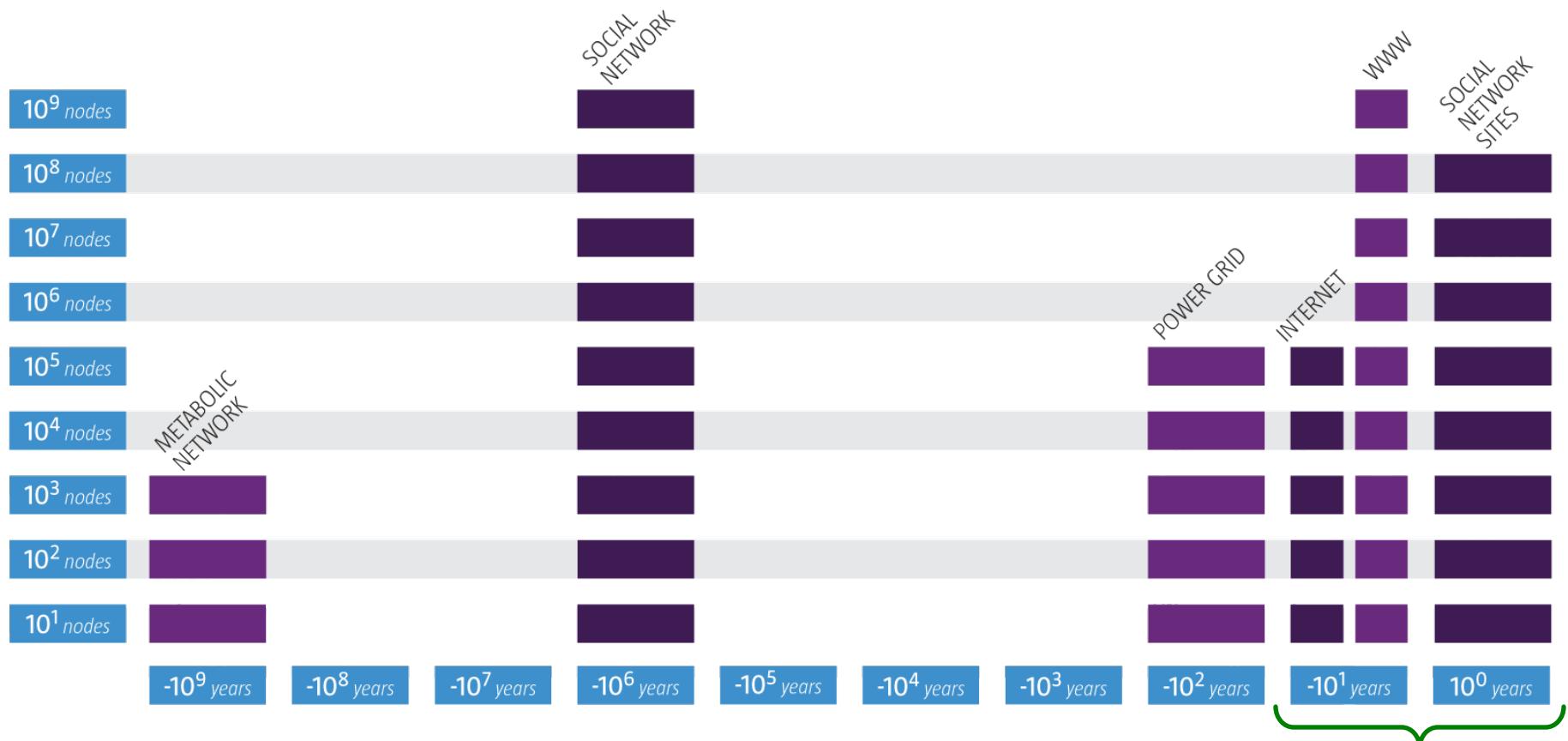
**We will never understand these systems unless we understand the networks behind them!**

**But Jure, why  
should I care about  
networks?**

# Why Networks? Why Now?

- **Universal language for describing complex data**
  - Networks from science, nature, and technology are more similar than one would expect
- **Shared vocabulary between fields**
  - Computer Science, Social science, Physics, Economics, Statistics, Biology
- **Data availability (/computational challenges)**
  - Web/mobile, bio, health, and medical
- **Impact!**
  - Social networking, Social media, Drug design

# Networks: Why Now?



Age and size of networks

CS!!

# Web – The Lab for Humanity



# **Networks and Applications**

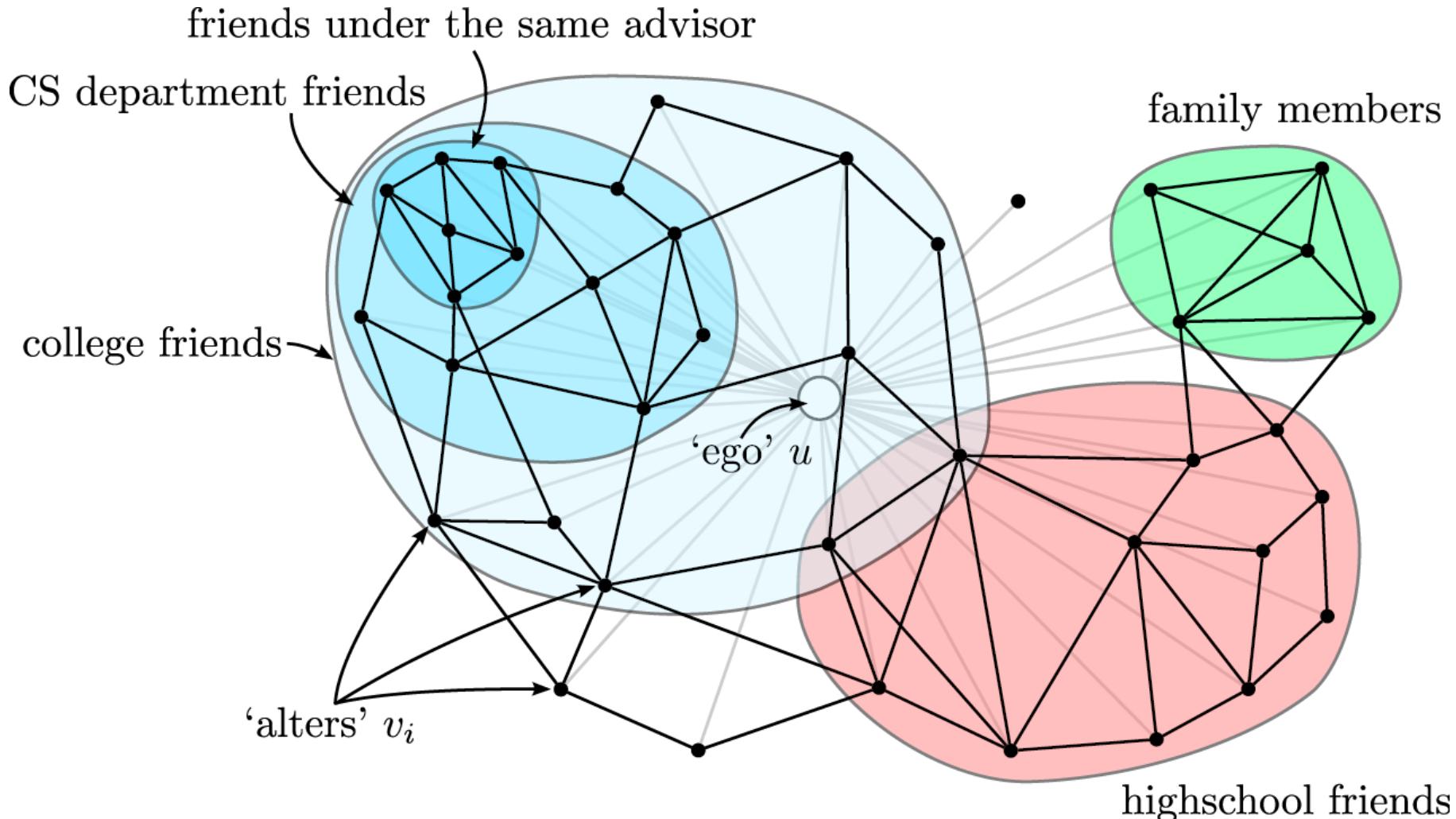
# (1) Networks: Social



Facebook social graph

4-degrees of separation [Backstrom-Boldi-Rosa-Ugander-Vigna, 2011]

# Application: Social Circle Detection



## Discover circles and why they exist

# (2) Networks: Infrastructure



Water supply distribution networks



Airline networks

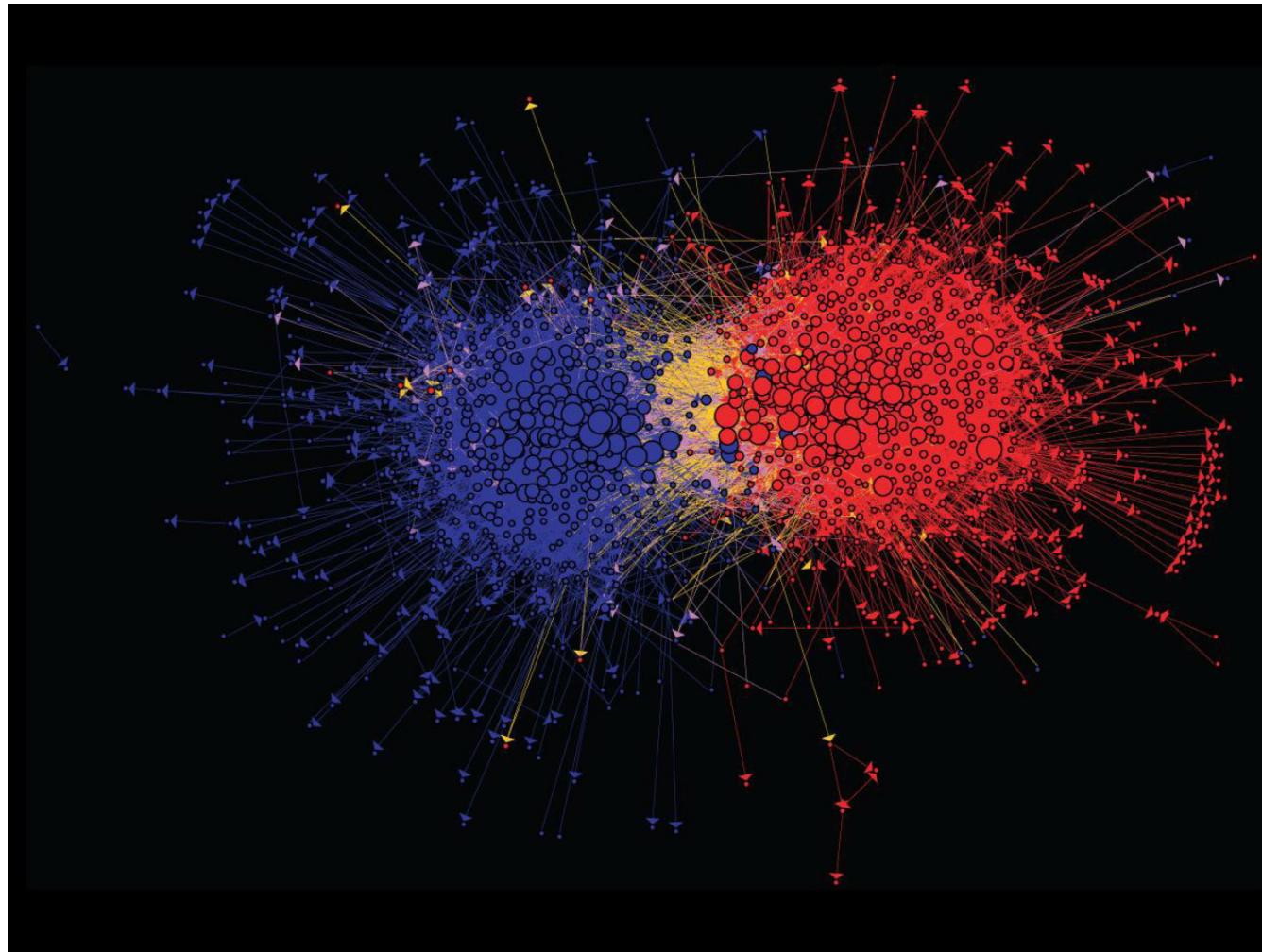
# Application: Modeling Epidemics

- Infrastructure networks are crucial for modeling epidemics



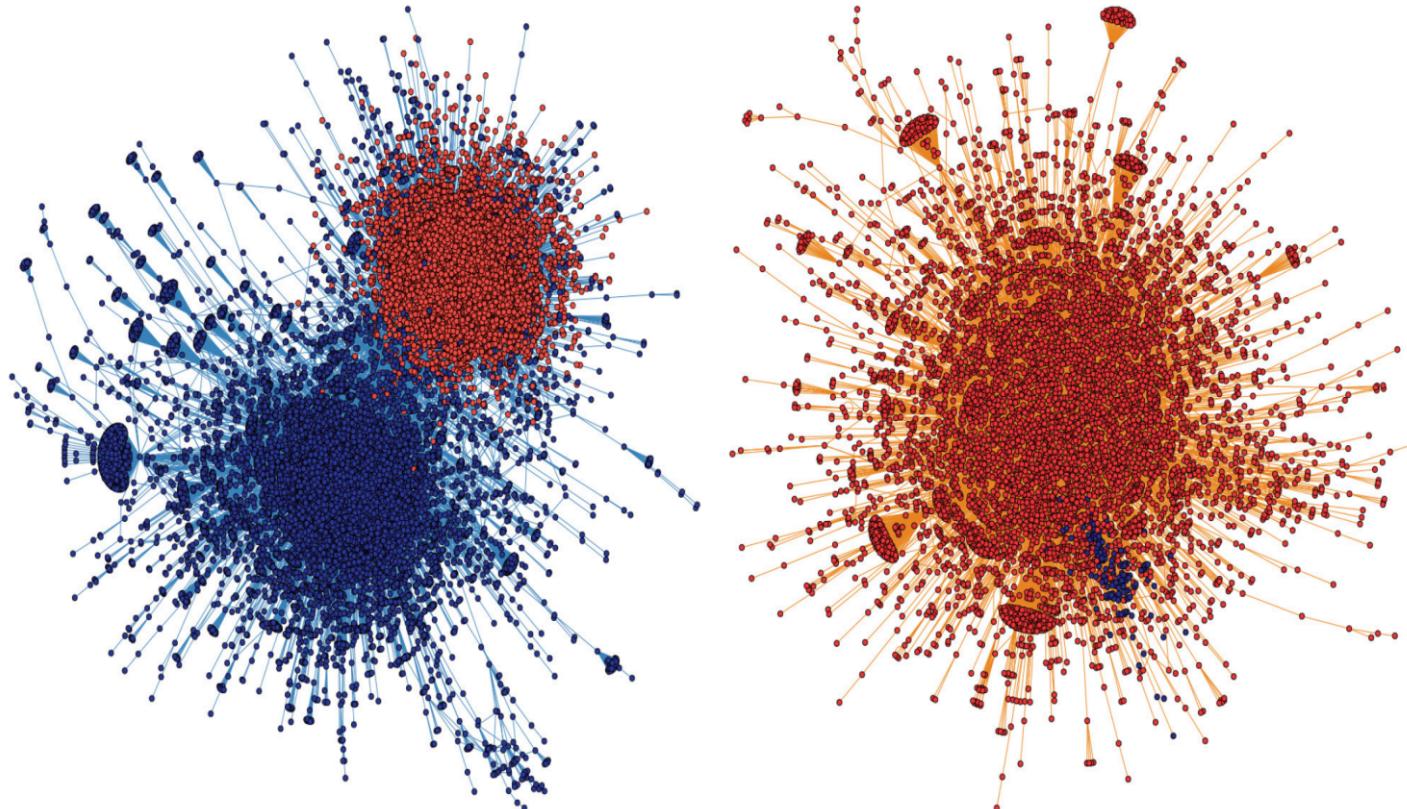
<http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0040961>

# (3) Networks: Online Media



**Connections between political blogs**  
Polarization of the network [Adamic-Glance, 2005]

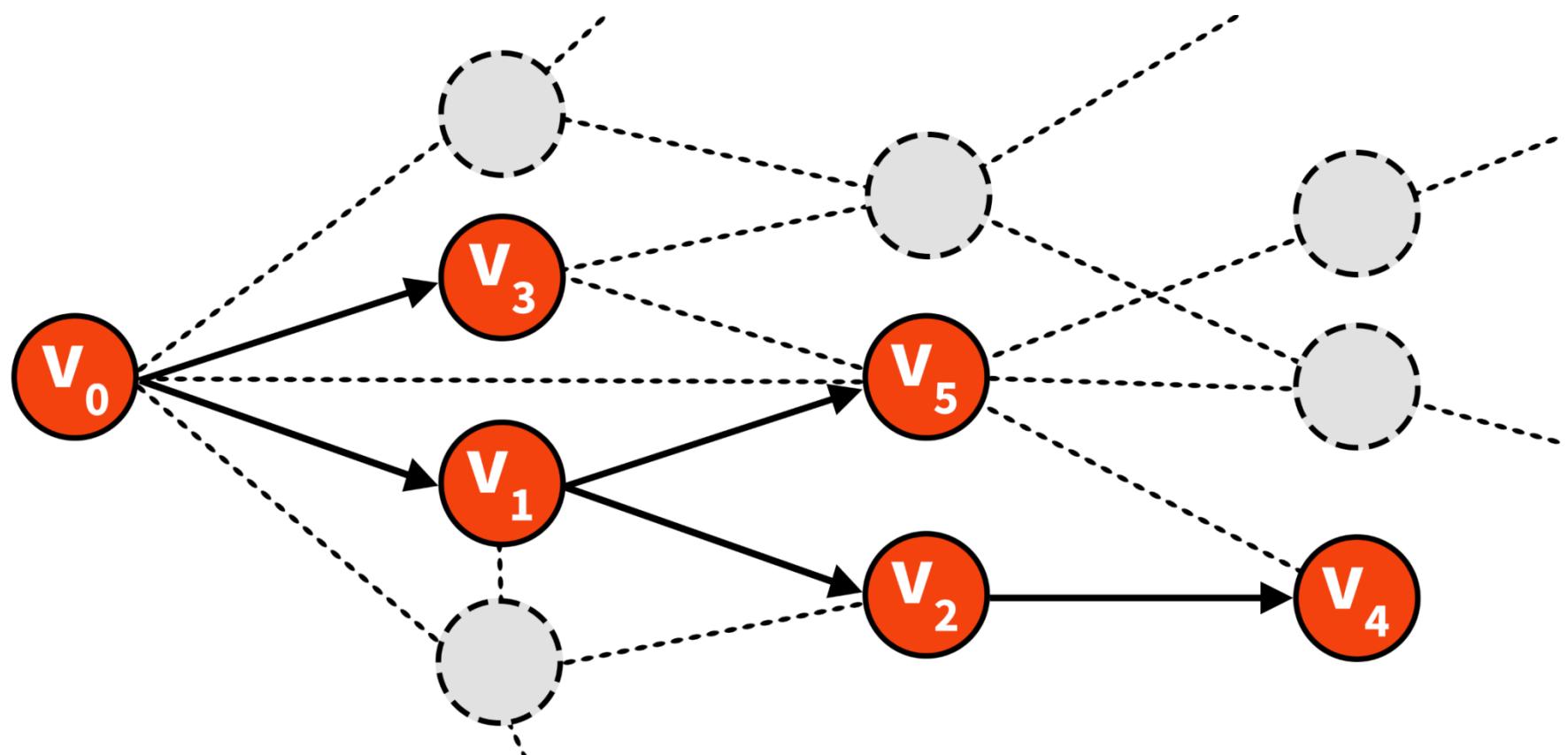
# Application: Polarization on Twitter



- **Retweet networks:**  
Polarized (left), Unpolarized (right)

Conover, M., Ratkiewicz, J., Francisco, M. R., Gonçalves, B., Menczer, F., & Flammini, A. "Political Polarization on Twitter." (2011)

# Application: Information Diffusion



**Information cascade in a social network**

# Facebook Information Cascades

## Timeline Photos

[Back to Album](#) · I fucking love science's Photos · I fucking love science's Page

[Previous](#) · [Next](#)

Thickness  $a$  ————— Radius  $z$



$$V = \pi z^2 a$$

$$V = \text{Pi}(z*z)a$$



I fucking love science

Seriously. If you have a pizza with radius "z" and thickness "a", its volume is  $\text{Pi}(z*z)a$ .

3,311 likes · Lina-von Der Stein, Iman Khallaf, 周明佳 and 73,191 others like this.

27,761 shares

1,470 comments

46 of 1,470

[archive](#)

Album: Timeline Photos

Shared with: Public

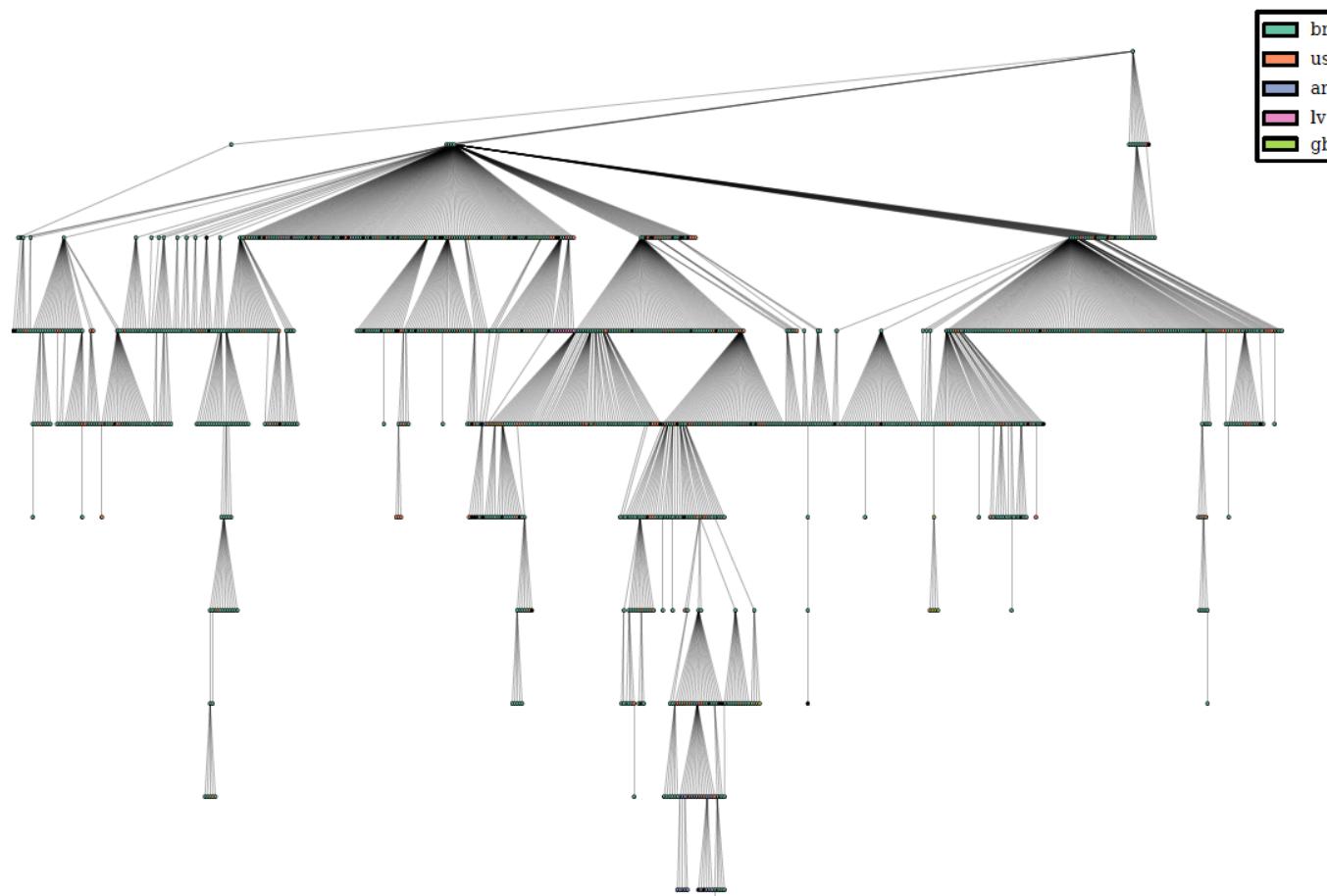
[Open Photo Viewer](#)

[Download](#)

[Embed Post](#)

Can cascades be predicted? Cheng et al., WWW '14.

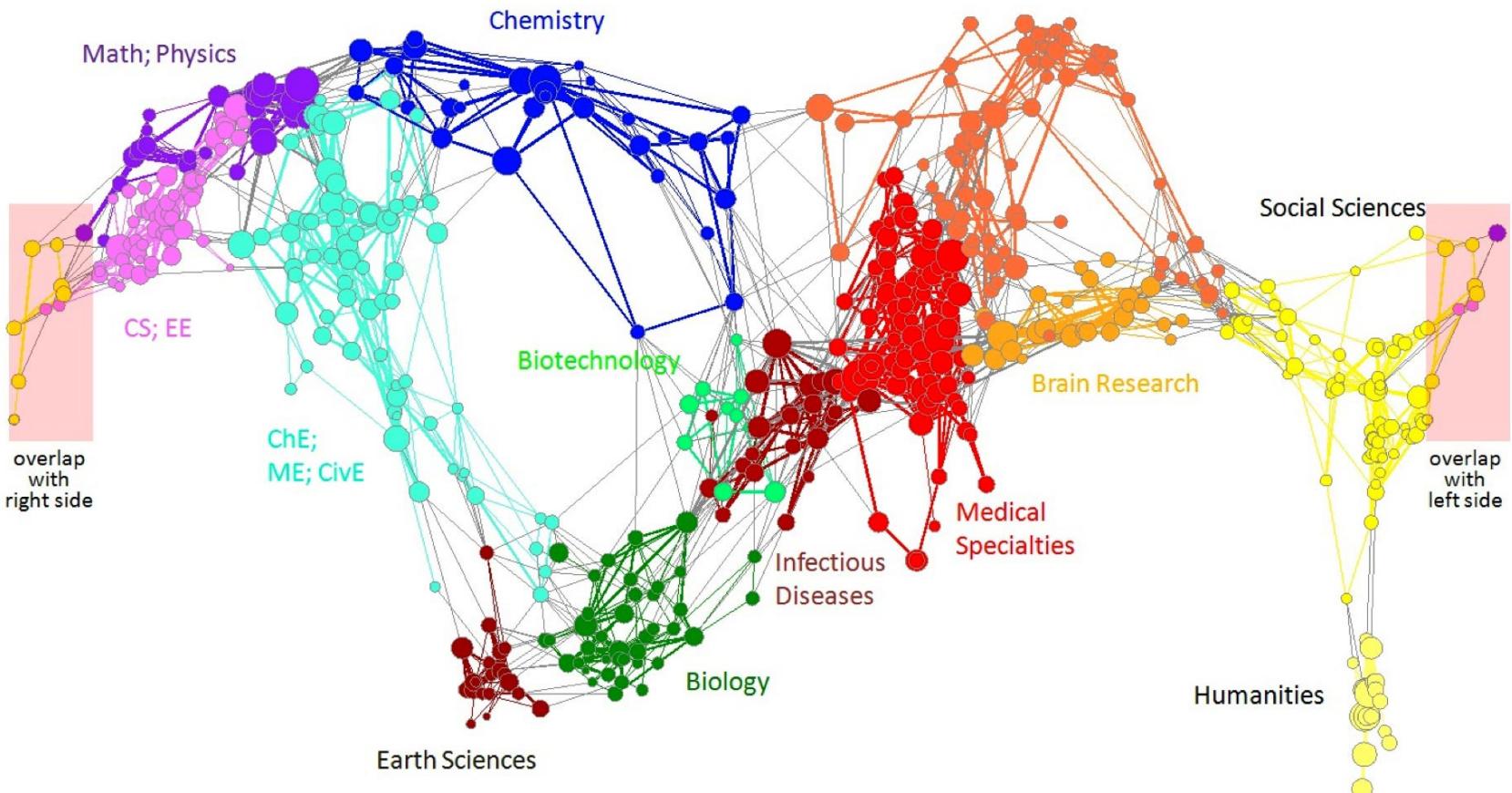
# Application: Product Adoption



**60-90% of LinkedIn users signed up due to an invitation from another user.**  
[Global Diffusion via Cascading Invitations: Structure, Growth, and Homophily](#).

Anderson et al., WWW '15.

# (4) Networks: Information, Knowledge

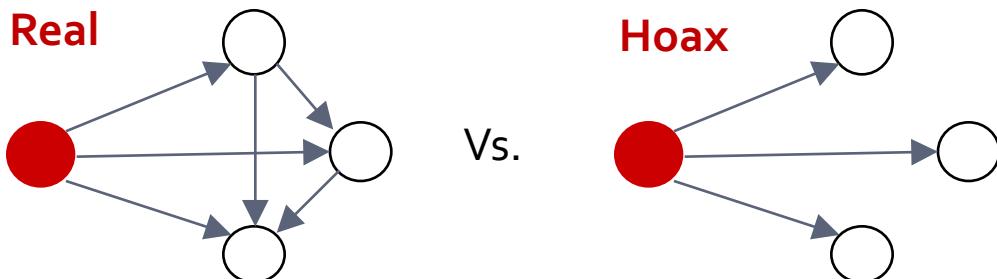


Citation networks and maps of science  
[Börner et al., 2012]

# Application: Misinformation

- Q: Is a given Wikipedia article a hoax?

- Real articles link more coherently:



Hoax article detection performance:

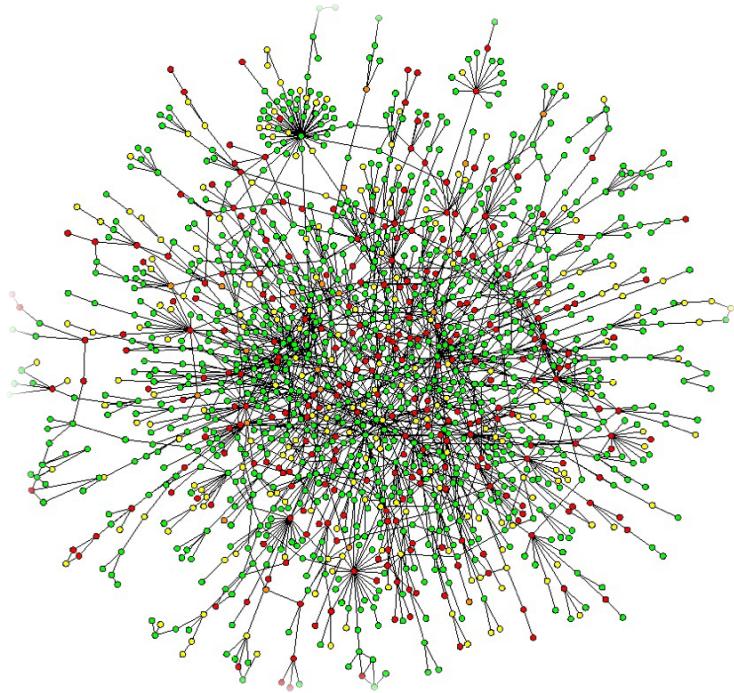
50%	66%	86%
Random	Human	WWW '16

Disinformation on the Web: Impact, Characteristics, and Detection of Wikipedia Hoaxes. Kumar et al. WWW '16.

This screenshot shows a Wikipedia page for 'Balboa French Creole'. The page title is 'Wikipedia:List of hoaxes on Wikipedia/Balboa French Creole'. A red box highlights a warning message: 'This article does not cite any references (sources). Please help improve this article by adding citations to reliable sources. Unsourced material may be challenged and removed. (January 2010)'.

The page content describes Balboa French Creole as a language used in Balboa Island, California, originating from a blend of French, English, Spanish, and German. It notes that the language is highly incomprehensible to most French speakers and is virtually extinct, with only 14 people remaining. The sidebar provides information about the language's native speakers, region, and ISO codes.

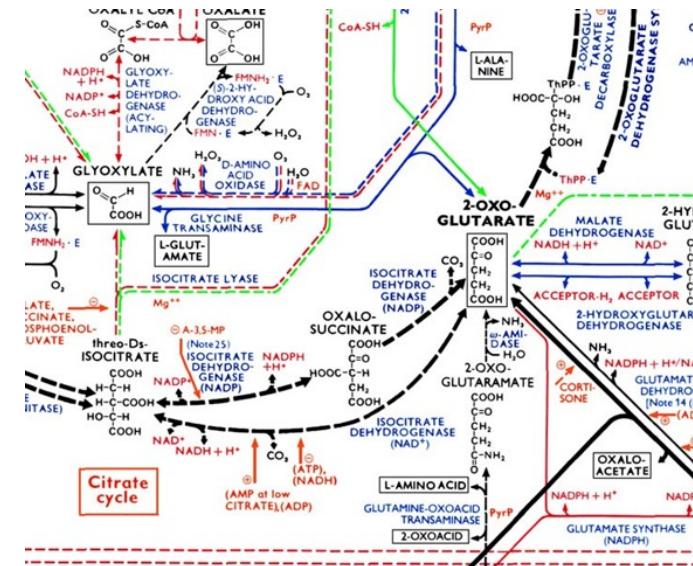
# (5) Networks: Biology



**Protein-protein interaction (PPI) networks:**

Nodes: Proteins

Edges: 'Physical' interactions



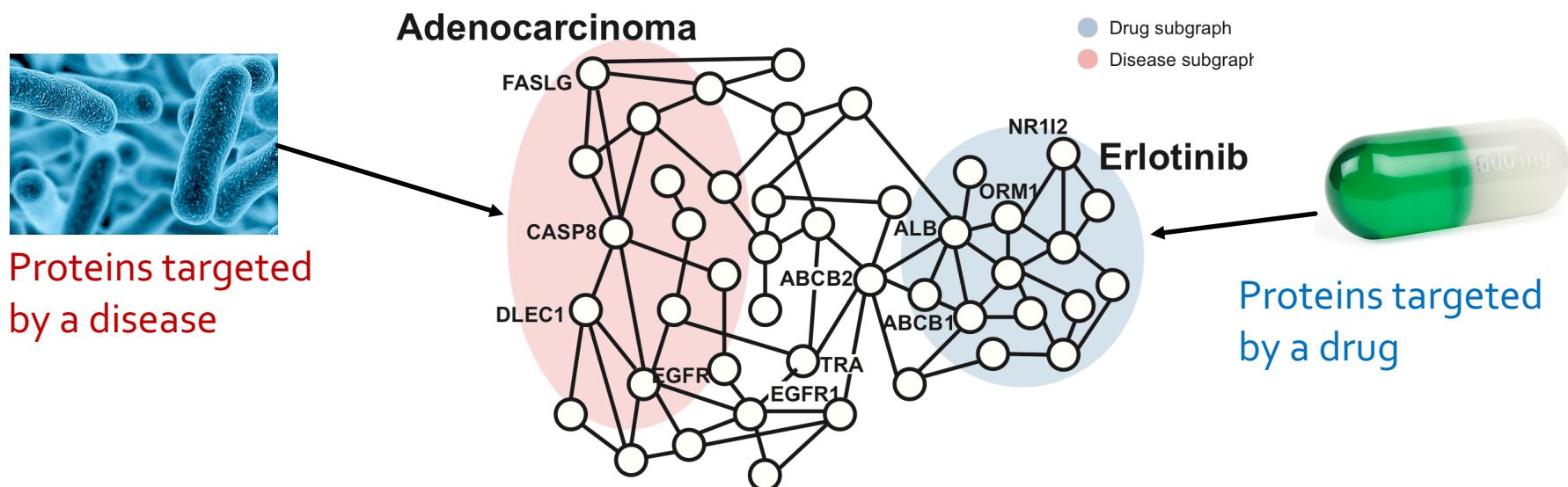
**Metabolic networks:**

Nodes: Metabolites and enzymes

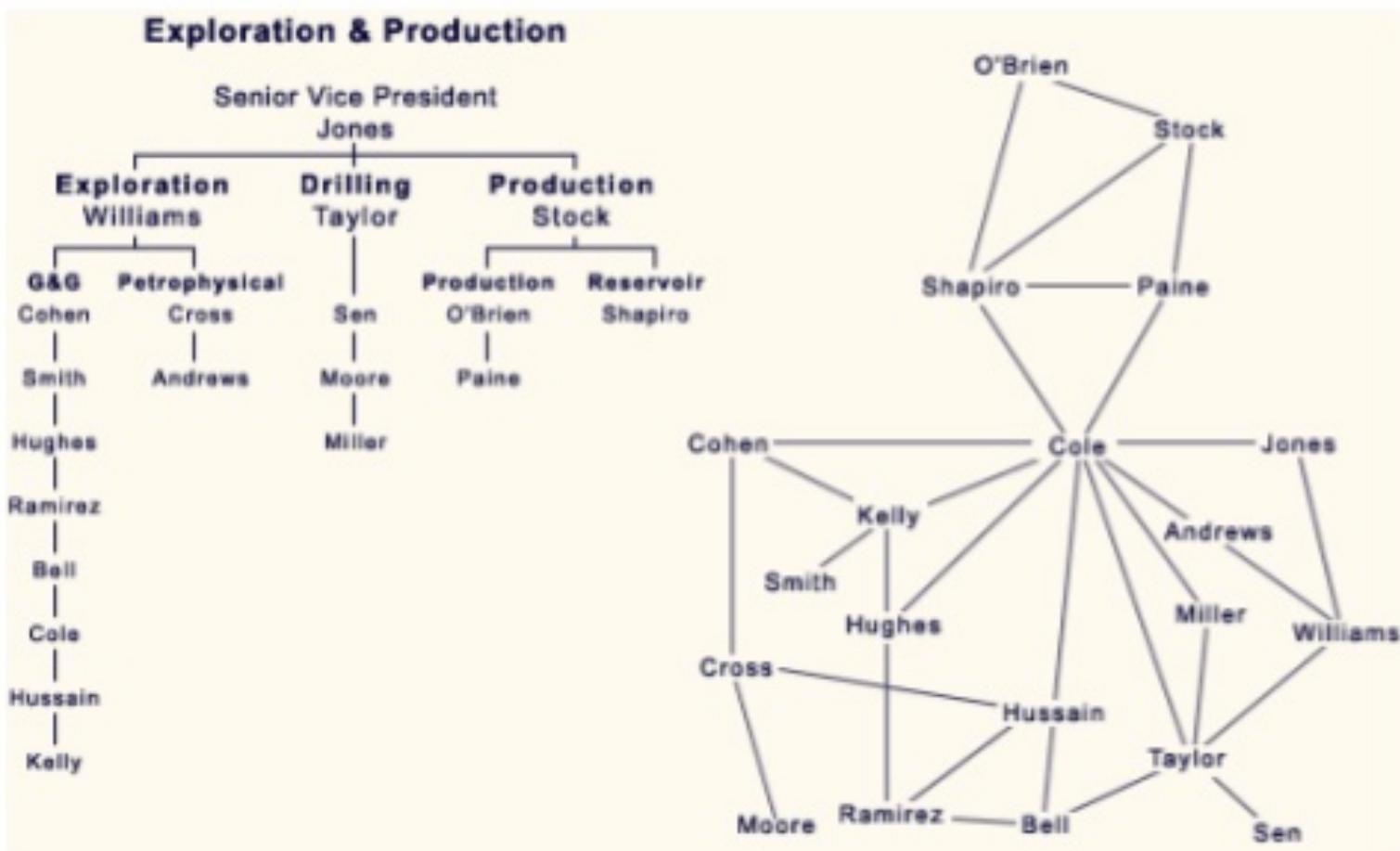
Edges: Chemical reactions

# Application: Drug Repurposing

- Q: Can we predict therapeutic uses of a drug?
- **Insight:** Proteins are worker molecules in a cell.  
Protein interaction networks capture how the cell works.

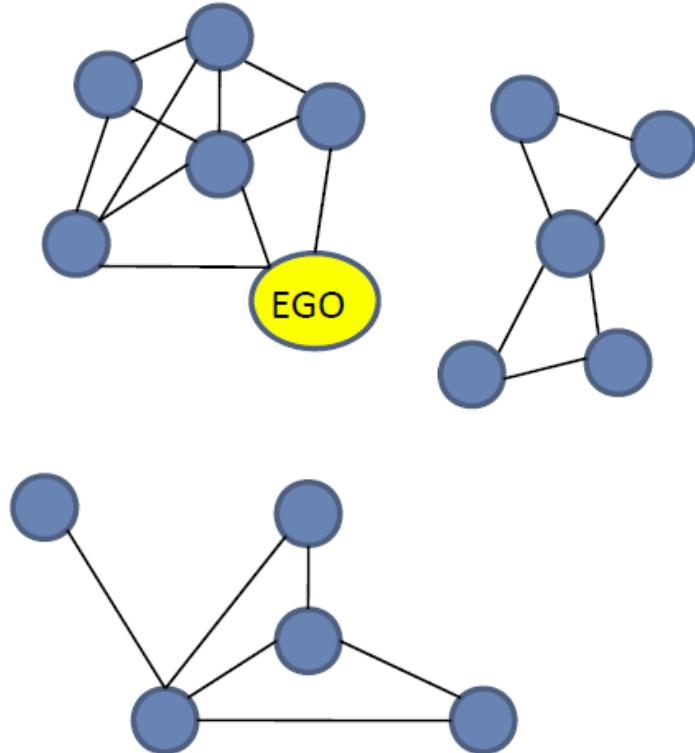


# (6) Networks: Organizations

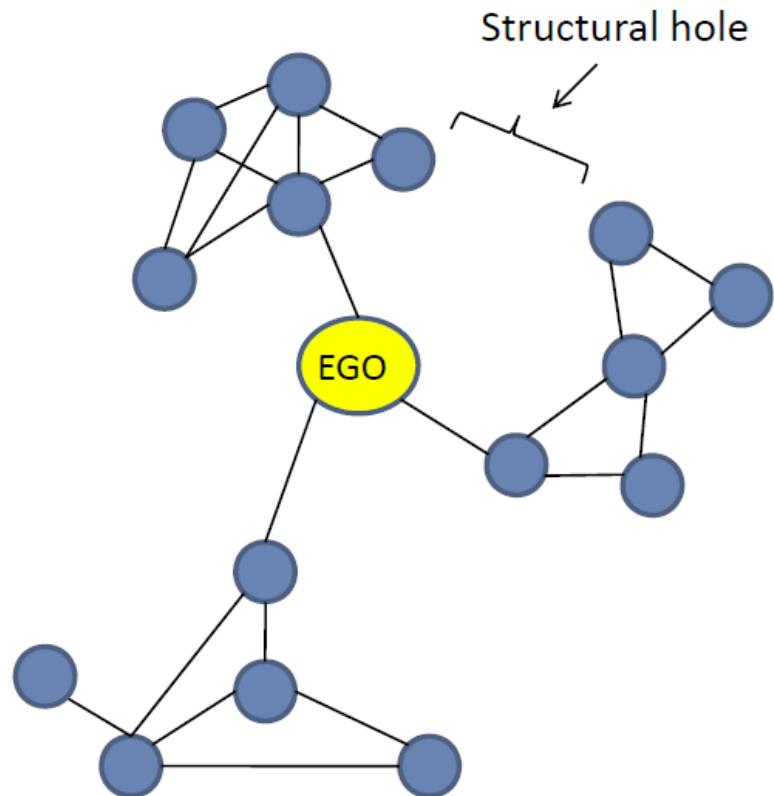


# Who are the central nodes in the organization? [Krebs, 2002]

# Application: Employee Success



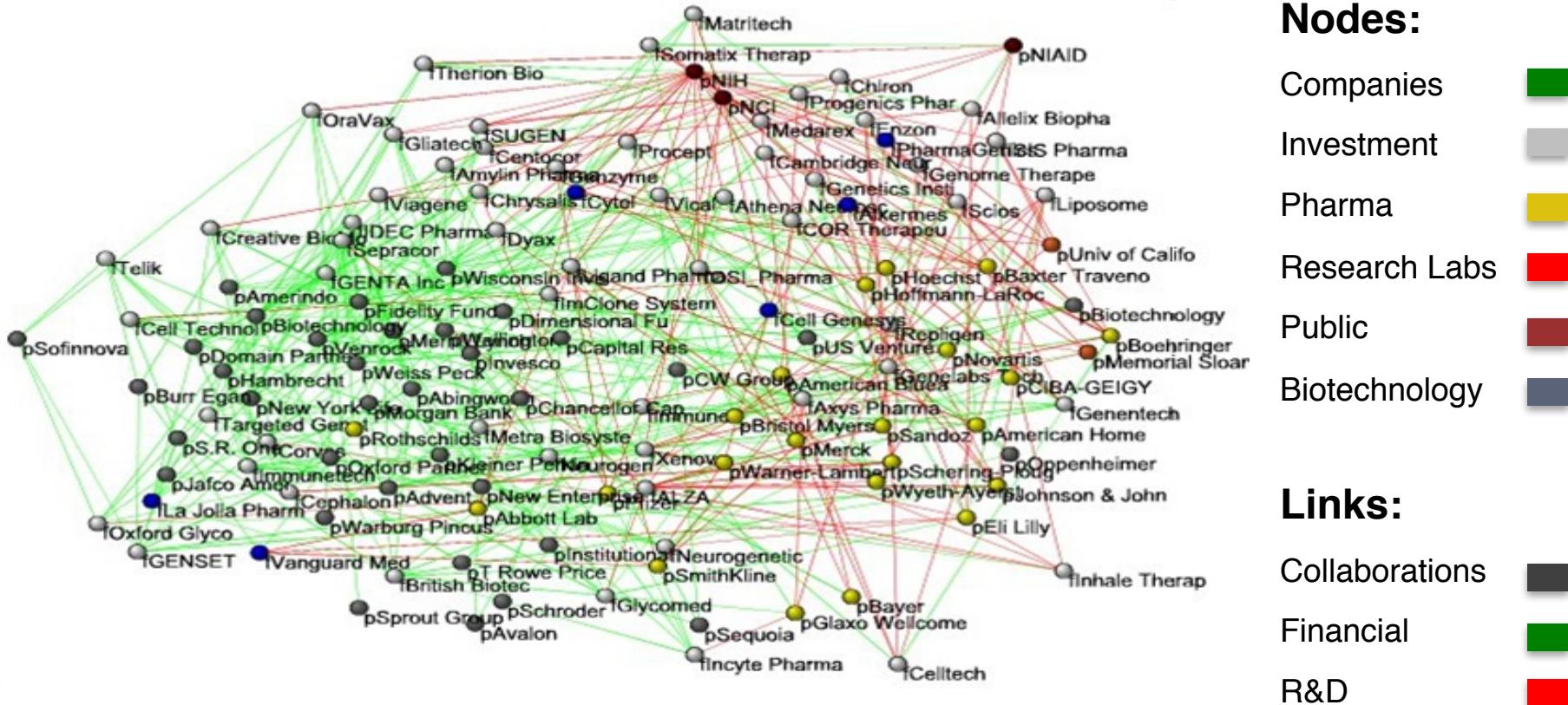
Few structural holes



Many structural holes

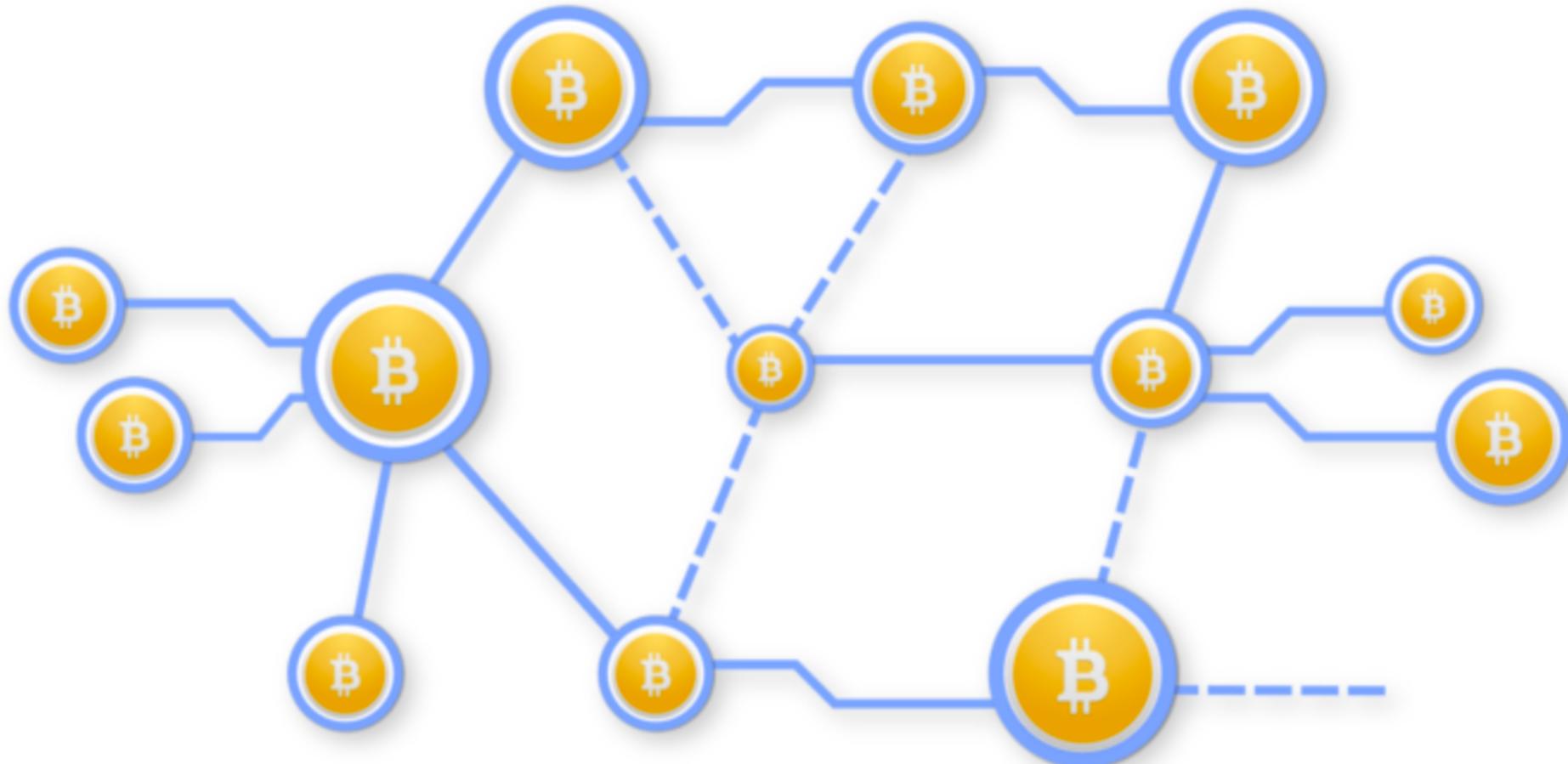
Structural Holes provide ego with access  
to novel information, power, freedom

# Networks: Economic Networks



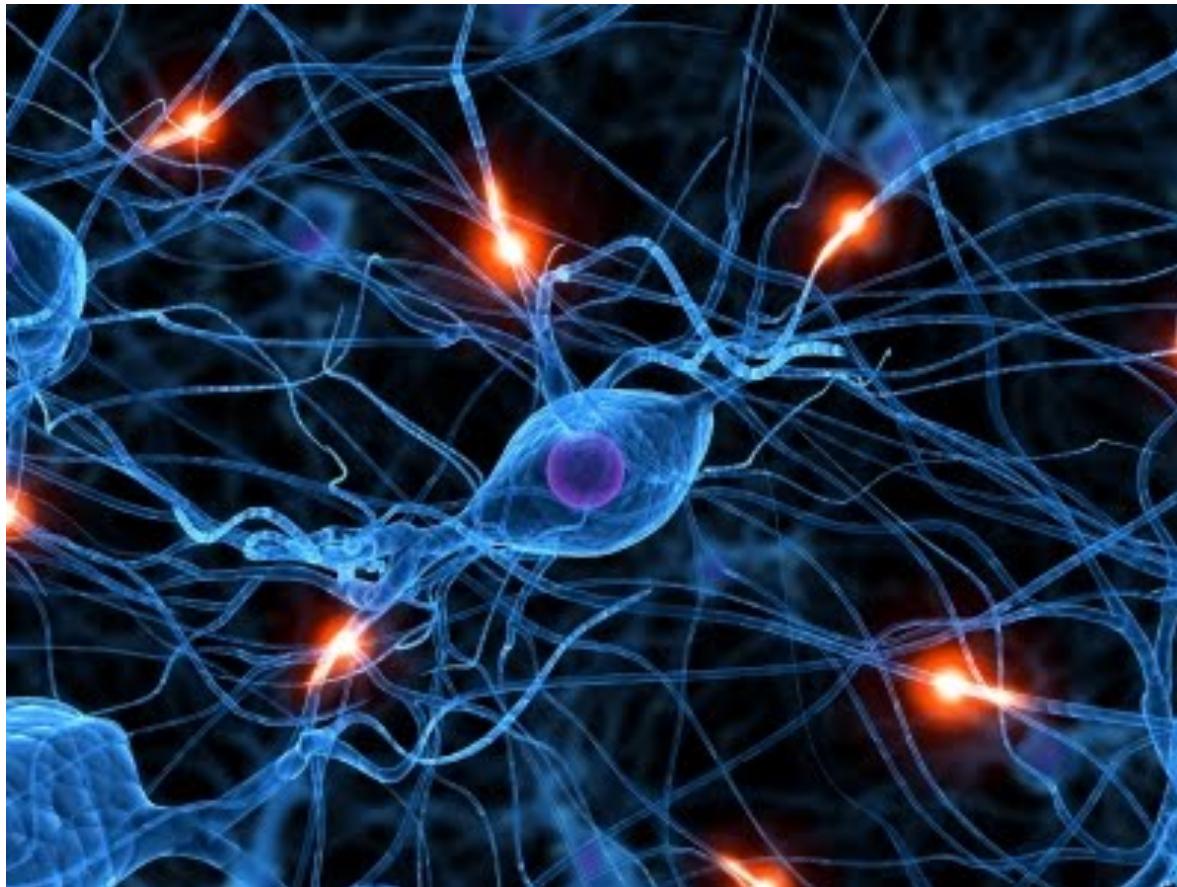
Bio-tech companies: Why companies succeed?  
[Powell-White-Koput, 2002]

# Networks: Transactions



- Detecting fraud and money laundering

# Networks: Brain



**Human brain has between  
10-100 billion neurons**  
**[Sporns, 2011]**

# Networks Really Matter

- If you want to understand the spread of diseases, **you need to figure out who will be in contact with whom**
- If you want to understand the structure of the Web, **you have to analyze the ‘links’.**
- If you want to understand dissemination of news or evolution of science, **you have to follow the flow.**

# About CS224W

# Reasoning about Networks

- What do we hope to achieve from studying networks?
  - Patterns and statistical **properties** of network data
  - **Design principles** and **models**
  - **Understand** why networks are organized the way they are
    - Predict behavior of networked systems

# Reasoning about Networks

- **How do we reason about networks?**
  - **Empirical:** Study network data to find organizational principles
    - How do we measure and quantify networks?
  - **Mathematical models:** Graph theory and statistical models
    - Models allow us to understand behaviors and distinguish surprising from expected phenomena
  - **Algorithms** for analyzing graphs
    - Hard computational challenges

# Networks: Structure & Process

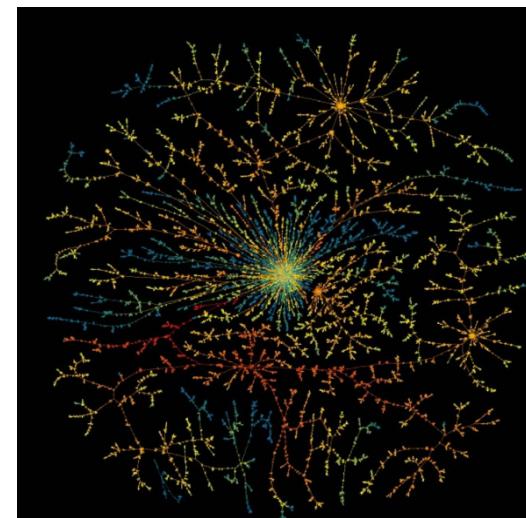
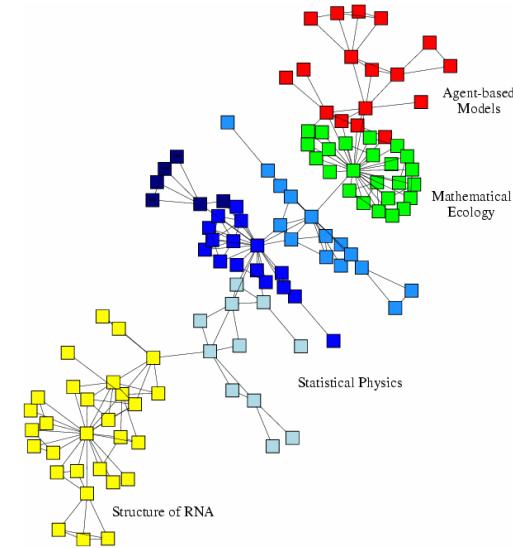
## What do we study in networks?

### ■ Structure and evolution:

- What is the structure of a network?
- Why and how did it come to have such structure?

### ■ Processes and dynamics:

- Networks provide “skeleton” for spreading of information, behavior, diseases
- How do information and diseases spread?



# How It All Fits Together

## Properties

Small diameter,  
Edge clustering

Scale-free

Strength of weak ties,  
Core-periphery

Densification power law,  
Shrinking diameters

Patterns of signed edge  
creation

Information virality,  
Memetracking

## Models

Small-world model,  
Erdős-Renyi model

Preferential attachment,  
Copying model

Kronecker Graphs

Microscopic model of  
evolving networks

Structural balance,  
Theory of status

Independent cascade model,  
Game theoretic model

## Algorithms

Decentralized search

PageRank, Hubs and  
authorities

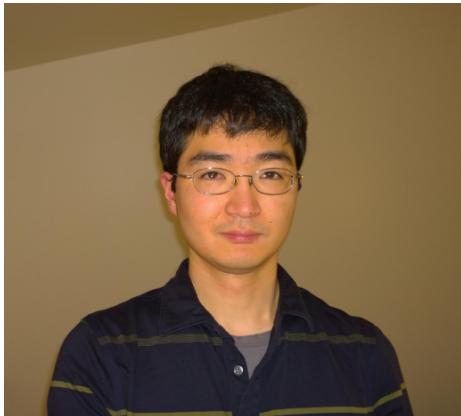
Community detection:  
Girvan-Newman, Modularity

Link prediction,  
Supervised random walks

Models for predicting  
edge signs

Influence maximization,  
Outbreak detection, LIM

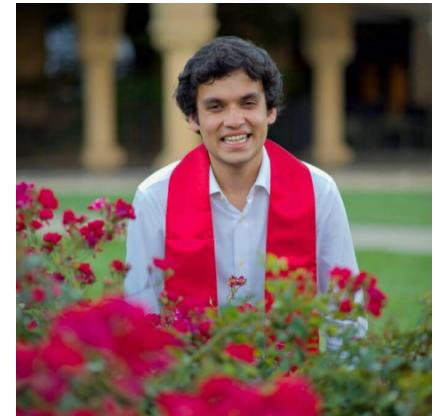
# Logistics: Course Assistants



Anthony Kim (Head TA)



Ziyi Yang



Anunay Kulshrestha



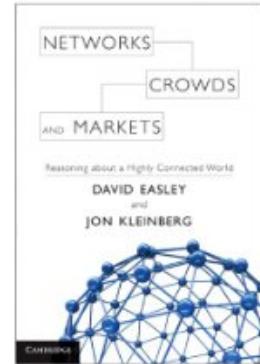
Silviana Ilcus



Praty Sharma

# Logistics: Website

- <http://cs224w.stanford.edu>
  - Slides posted the night before the class
- **Readings:**
  - Chapters from Easley&Kleinberg
  - Papers
- **Optional readings:**
  - Papers and pointers to additional literature
  - **This will be very useful for project proposals**



# Logistics: Communication

- **Piazza Q&A website:**
  - <http://piazza.com/stanford/fall2017/cs224w>
    - Use access code “snap”
  - **Please participate and help each other!  
(2% of the grade)**
- **For e-mailing course staff, always use:**
  - [cs224w-aut1718-staff@lists.stanford.edu](mailto:cs224w-aut1718-staff@lists.stanford.edu)
- We will post course announcements to Piazza  
(make sure you check it regularly)

# Homework, Write-ups

- **Assignments are long and take time (10-20h)**  
**Start early!**
  - A combination of data analysis, algorithm design, and math
- **How to submit?**
  - Upload via Gradescope (<http://gradescope.com>)
    - To register use the code MRW7PY
      - Use your Stanford email (if non-SCPD) and include your Stanford ID # (everyone)
    - **IMPORTANT:** One answer per page!
  - **Code and project write-ups** (proposal, milestone, final report) have to also be uploaded at <http://snap.stanford.edu/submit/>
- **Total of 2 late periods for the quarter:**
  - Late period expires on Monday at 23:59 Pacific Time
  - You can use at most 1 late period per assignment
  - No late periods for submissions related to final project

# Course Projects

- **Substantial course project:**
  - **Experimental evaluation** of algorithms and models on an interesting network dataset
  - A **theoretical project** that considers a model, an algorithm and derives a rigorous result about it
  - Develop **scalable algorithms** for massive graphs
- **Performed in groups of up to 3 students**
  - Fine to have groups of 1 or 2. The team size will be taken under consideration when evaluating the scope of the project in breadth and depth. But 3 person teams can be more efficient.
  - Project is the **main work** for the class
  - We will help with ideas, data and mentoring
  - Start thinking about this now!
  - Ok to combine projects. Clearly indicate, which part of the project is done for CS224W and which part is done for the other class.
- Poster session with many external visitors
- **Read:** <http://cs224w.stanford.edu/info.html#proj>

# Course Schedule

Week	Assignment	Due on (23:59 PST)
2	<b>Homework 0</b>	October 5
3	<b>Homework 1</b>	October 12
4	<b>Project proposal</b>	October 19 (no late periods!)
5	<b>Homework 2</b>	October 26
6	<b>Work on the project</b>	
7	<b>Homework 3</b>	November 9
8	<b>Project milestone</b>	November 16 (no late periods!)
9	<b>Thanksgiving break</b>	
10	<b>Homework 4</b>	November 30
11	<b>Project report</b>	<b>Sun December 10</b> (no late periods!)
	<b>Poster session</b>	<b>Mon, December 11</b> <b>12:15-3:15pm</b>

# Work for the Course & Grading

- **Final grade will be composed of:**
  - **Homework: 48%**
    - Homework 1,2,3,4: 11.75% each, HW0: 1%
  - **Substantial class project: 50%**
    - Proposal: 20%
    - Project milestone: 20%
    - Final report: 50%
    - Poster presentation: 10%
  - **Piazza participation, snap code contribution: 2%**
    - Students between grades get extra credit for Piazza participation

# Prerequisites

- **No single topic in the course is too hard by itself**
- **But we will cover and touch upon many topics and this is what makes the course hard**
  - **Good background in:**
    - Algorithms and graph theory
    - Probability and statistics
    - Linear algebra
  - **Programming:**
    - You should be able to write non-trivial programs (in Python)
  - **2 recitation sessions:**
    - Review of Probability, Linear Algebra, and Proof Techniques:  
**Thursday, 9/28 (4:30-5:20pm, Gates B03)**
    - SNAP.PY review and installation party:  
**Friday, 9/29 (4:30-5:20pm, Gates B03)**

# Network Analysis Tools

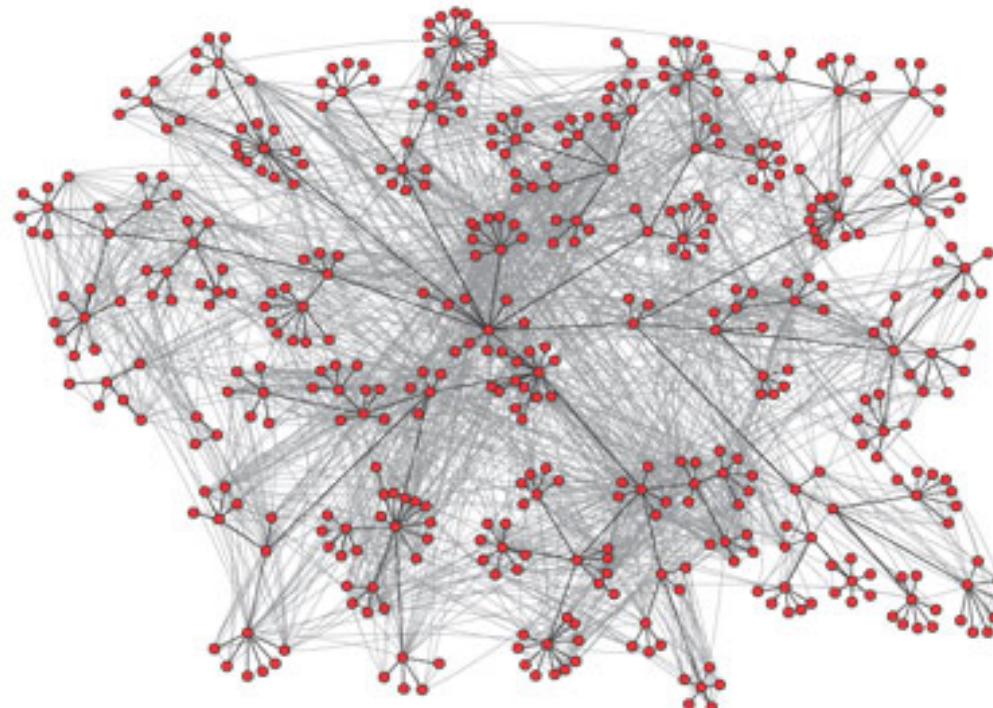
- We highly recommend **SNAP**:
  - **SNAP.PY**: Python ease of use, most of C++ scalability
    - HW0 asks you to do some very basic network analysis with `snap.py`
      - If you find HW0 difficult, this class is probably not for you
  - **SNAP C++**: more challenging but more scalable
  - Other tools include NetworkX, iGraph

SNAP.PY review and installation party:  
**Friday, 9/29 (4:30-5:20pm, Gates B03)**

# **Starter Topic:**

# **Structure of Graphs**

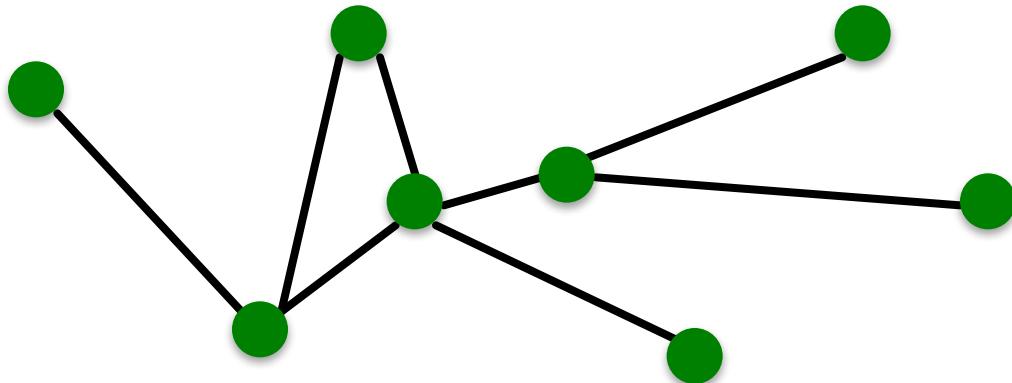
# Structure of Networks?



A network is a collection of objects where some pairs of objects are connected by links

**What is the structure of the network?**

# Components of a Network



- **Objects:** nodes, vertices  $N$
- **Interactions:** links, edges  $E$
- **System:** network, graph  $G(N,E)$

# Networks or Graphs?

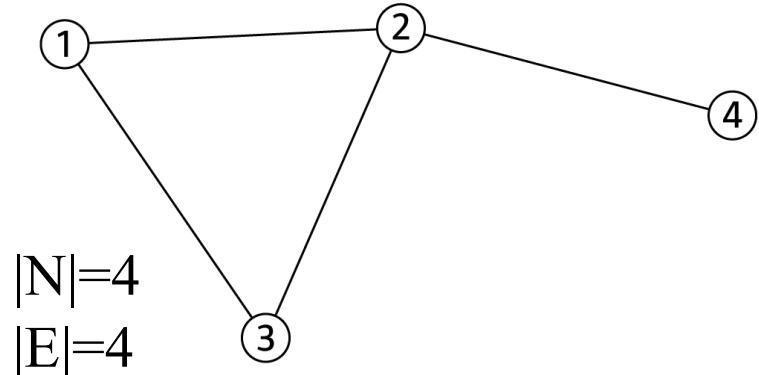
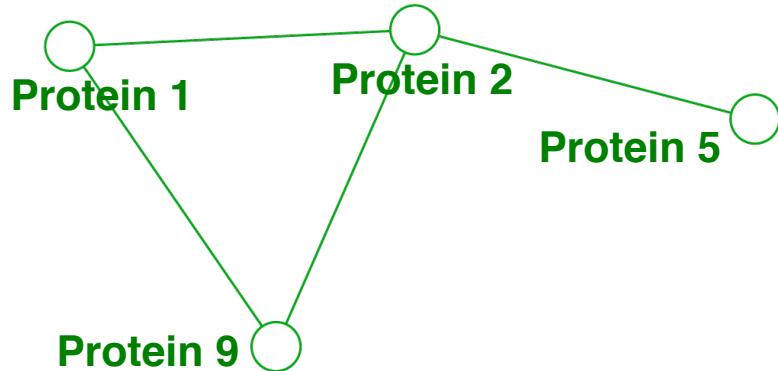
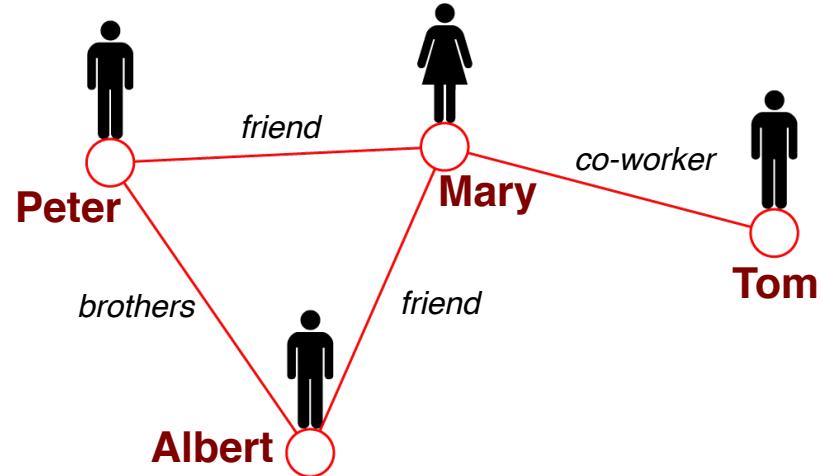
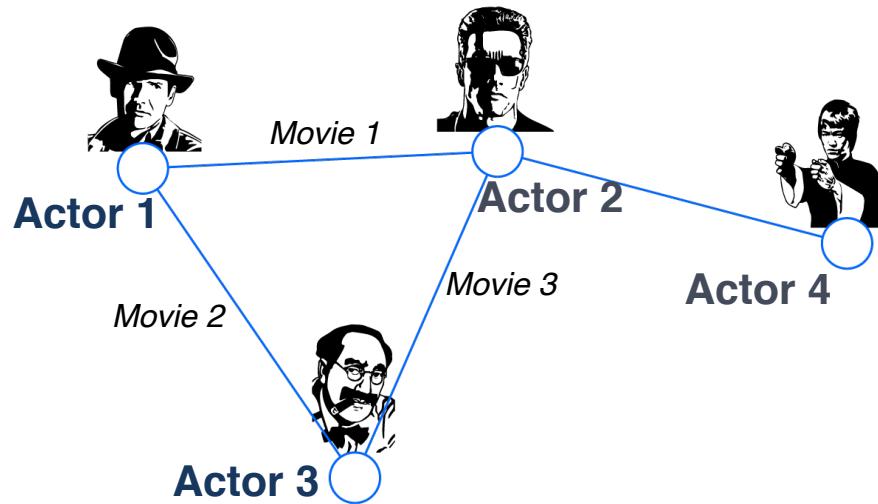
- **Network** often refers to real systems
  - Web, Social network, Metabolic network

**Language:** Network, node, link
- **Graph** is a mathematical representation of a network
  - Web graph, Social graph (a Facebook term)

**Language:** Graph, vertex, edge

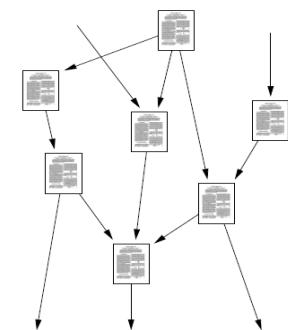
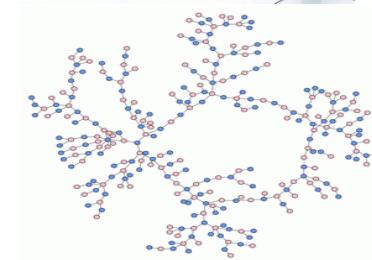
We will try to make this distinction whenever it is appropriate, but in most cases we will use the two terms interchangeably

# Networks: Common Language



# Choosing Proper Representations

- If you connect individuals that work with each other, you will explore a **professional network**
- If you connect those that have a sexual relationship, you will be exploring **sexual networks**
- If you connect scientific papers that cite each other, you will be studying the **citation network**
- **If you connect all papers with the same word in the title, you will be exploring what?** It is a network, nevertheless



# How do you define a network?

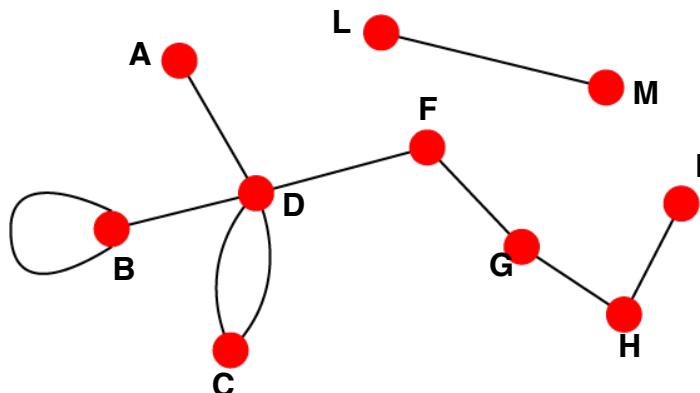
- **How to build a graph:**
  - What are nodes?
  - What are edges?
- **Choice of the proper network representation of a given domain/problem determines our ability to use networks successfully:**
  - In some cases there is a unique, unambiguous representation
  - In other cases, the representation is by no means unique
  - The way you assign links will determine the nature of the question you can study

# **Choice of Network Representation**

# Directed vs. Undirected Graphs

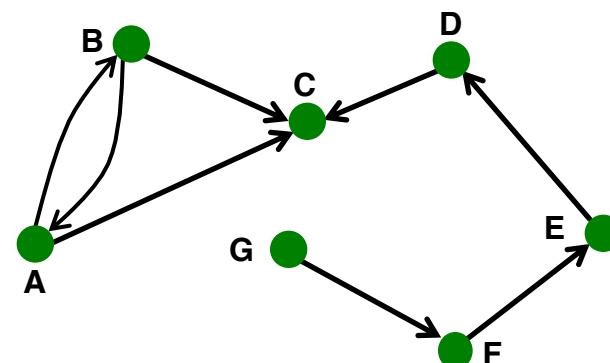
## Undirected

- Links: undirected  
(symmetrical, reciprocal)



## Directed

- Links: directed  
(arcs)



## Examples:

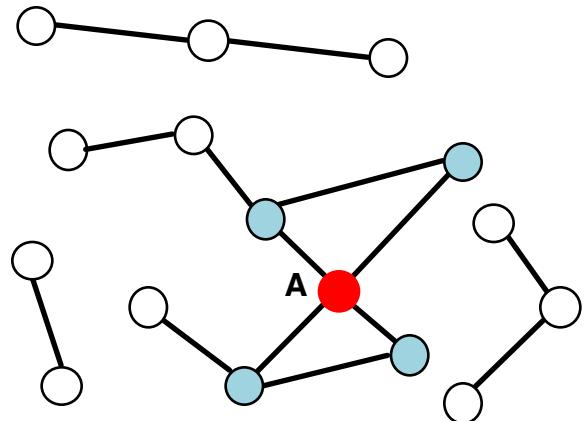
- Collaborations
- Friendship on Facebook

## Examples:

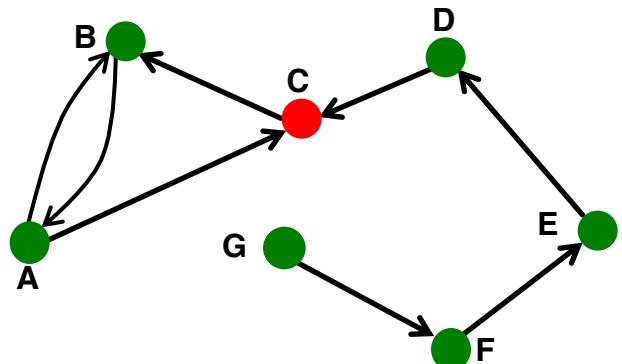
- Phone calls
- Following on Twitter

# Node Degrees

Undirected



Directed



**Source:** Node with  $k^{in} = 0$

**Sink:** Node with  $k^{out} = 0$

**Node degree,  $k_i$ :** the number of edges adjacent to node  $i$

$$k_A = 4$$

**Avg. degree:**  $\bar{k} = \langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i = \frac{2E}{N}$

In directed networks we define an **in-degree** and **out-degree**. The (total) degree of a node is the sum of in- and out-degrees.

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

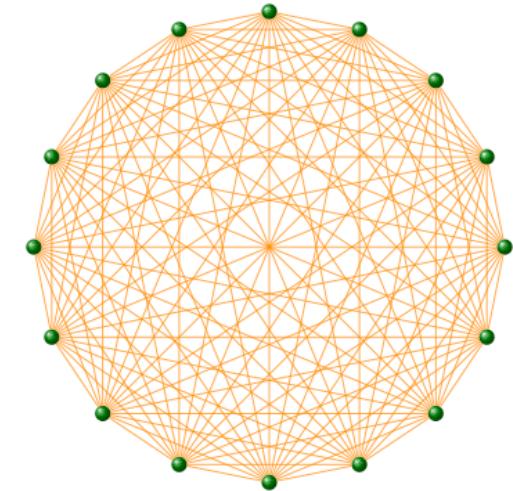
$$\bar{k} = \frac{E}{N}$$

$$\overline{k^{in}} = \overline{k^{out}}$$

# Complete Graph

The **maximum number of edges** in an undirected graph on  $N$  nodes is

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



An undirected graph with the number of edges  $E = E_{\max}$  is called a **complete graph**, and its average degree is  $N-1$

# Bipartite Graph

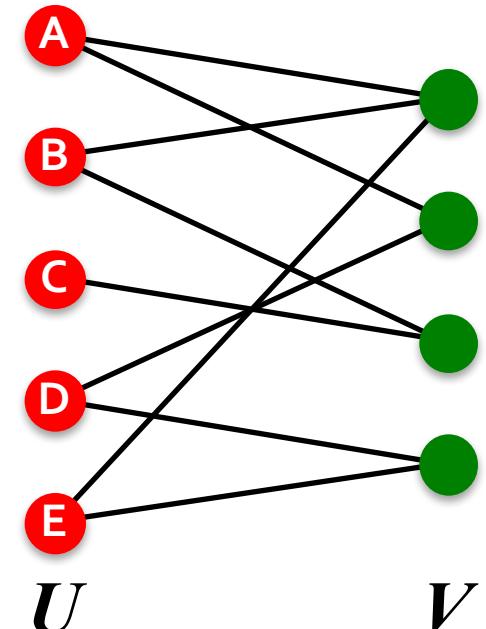
- **Bipartite graph** is a graph whose nodes can be divided into two disjoint sets  $U$  and  $V$  such that every link connects a node in  $U$  to one in  $V$ ; that is,  $U$  and  $V$  are **independent sets**

- **Examples:**

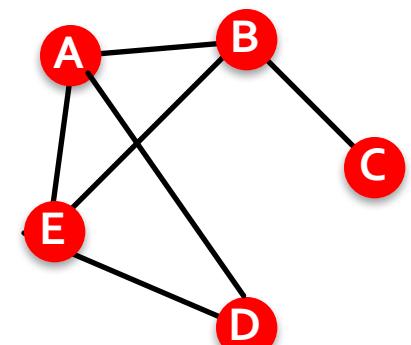
- Authors-to-papers (they authored)
- Actors-to-Movies (they appeared in)
- Users-to-Movies (they rated)

- **“Folded” networks:**

- Author collaboration networks
- Movie co-rating networks

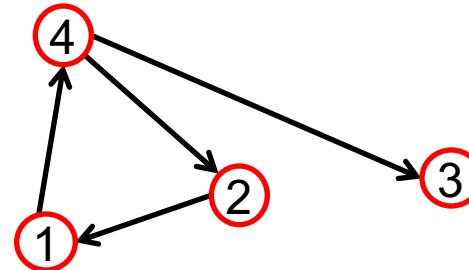
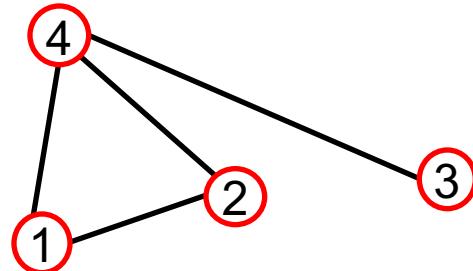


$U$        $V$



Folded version of the graph above

# Representing Graphs: Adjacency Matrix



$A_{ij} = 1$  if there is a link from node  $i$  to node  $j$

$A_{ij} = 0$  otherwise

$$A = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$

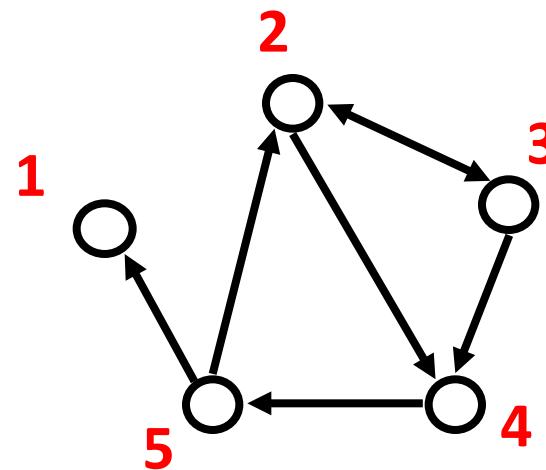
$$A = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{pmatrix}$$

Note that for a directed graph (right) the matrix is not symmetric.

# Representing Graphs: Edge list

- Represent graph as a set of edges:

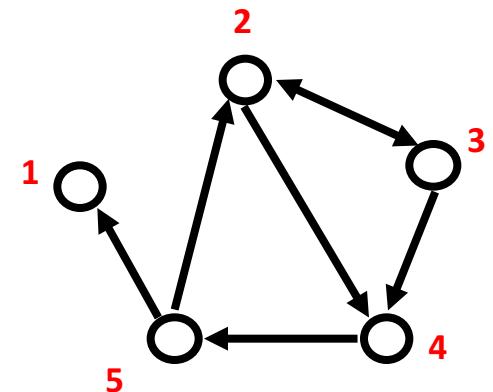
- (2, 3)
- (2, 4)
- (3, 2)
- (3, 4)
- (4, 5)
- (5, 2)
- (5, 1)



# Representing Graphs: Adjacency list

## ■ Adjacency list:

- Easier to work with if network is
  - Large
  - Sparse
- Allows us to quickly retrieve all neighbors of a given node
  - 1:
  - 2: 3, 4
  - 3: 2, 4
  - 4: 5
  - 5: 1, 2



# Networks are Sparse Graphs

Most real-world networks are **sparse**

$$E \ll E_{\max} \text{ (or } \bar{k} \ll N-1)$$

WWW (Stanford-Berkeley):	$N=319,717$	$\langle k \rangle = 9.65$
Social networks (LinkedIn):	$N=6,946,668$	$\langle k \rangle = 8.87$
Communication (MSN IM):	$N=242,720,596$	$\langle k \rangle = 11.1$
Coauthorships (DBLP):	$N=317,080$	$\langle k \rangle = 6.62$
Internet (AS-Skitter):	$N=1,719,037$	$\langle k \rangle = 14.91$
Roads (California):	$N=1,957,027$	$\langle k \rangle = 2.82$
Proteins (S. Cerevisiae):	$N=1,870$	$\langle k \rangle = 2.39$

(Source: Leskovec et al., Internet Mathematics, 2009)

**Consequence: Adjacency matrix is filled with zeros!**

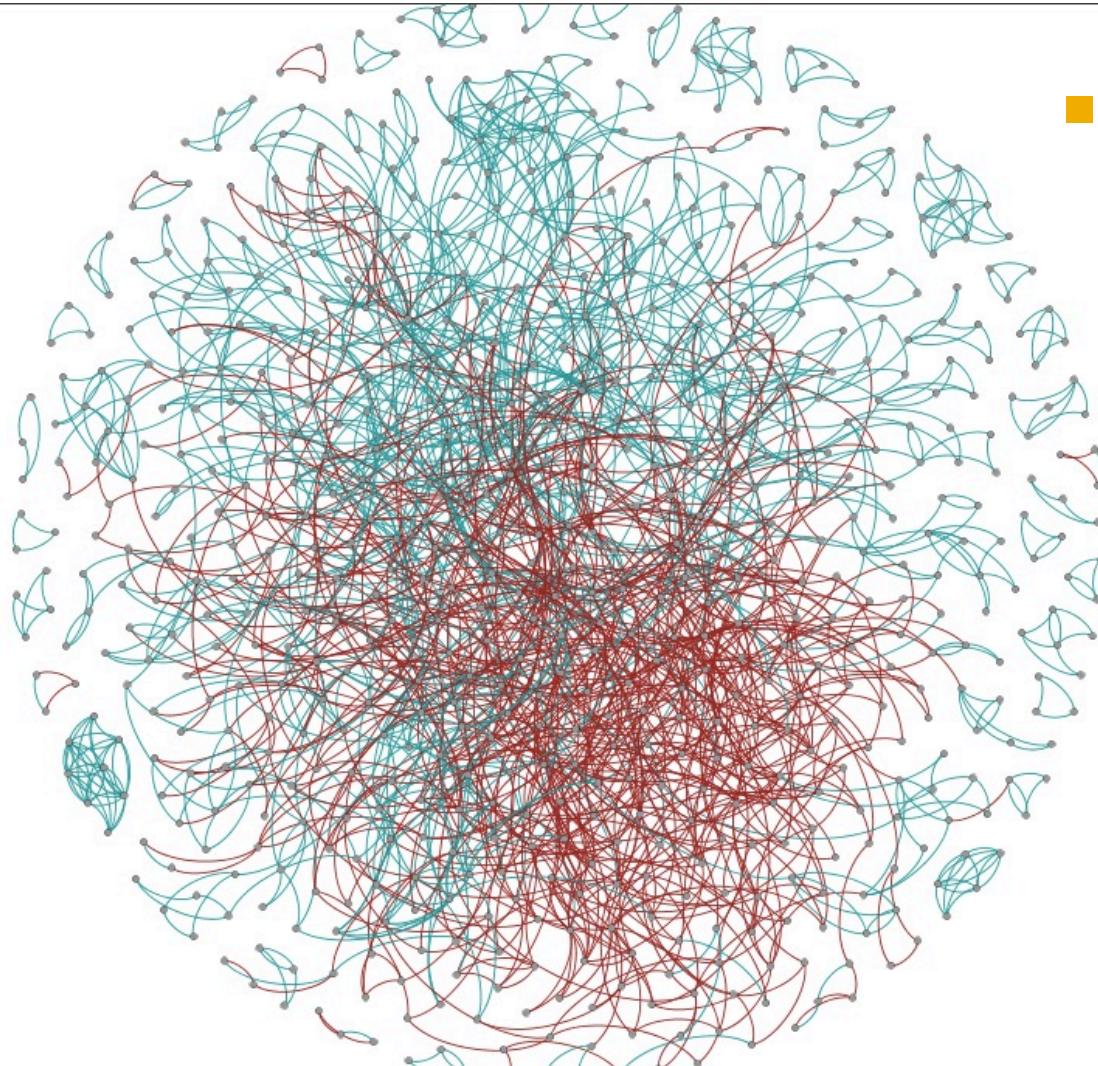
(Density of the matrix ( $E/N^2$ ): WWW= $1.51 \times 10^{-5}$ , MSN IM =  $2.27 \times 10^{-8}$ )

# Edge Attributes

## Possible options:

- Weight (e.g. frequency of communication)
- Ranking (best friend, second best friend...)
- Type (friend, relative, co-worker)
- Sign: Friend vs. Foe, Trust vs. Distrust
- Properties depending on the structure of the rest of the graph: number of common friends

# Positive and Negative Weights



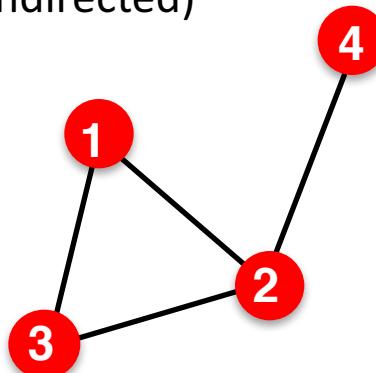
*sample of positive & negative ratings from Epinions network*

- One person trusting/distrusting another
  - Research challenge: How does one ‘propagate’ negative feelings in a social network? Is my enemy’s enemy my friend?

# More Types of Graphs

## ■ Unweighted

(undirected)



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

$$A_{ii} = 0$$

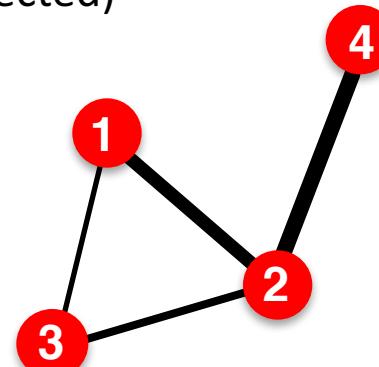
$$A_{ij} = A_{ji}$$

$$E = \frac{1}{2} \sum_{i,j=1}^N A_{ij} \quad \bar{k} = \frac{2E}{N}$$

Examples: Friendship, Hyperlink

## ■ Weighted

(undirected)



$$A_{ij} = \begin{pmatrix} 0 & 2 & 0.5 & 0 \\ 2 & 0 & 1 & 4 \\ 0.5 & 1 & 0 & 0 \\ 0 & 4 & 0 & 0 \end{pmatrix}$$

$$A_{ii} = 0$$

$$A_{ij} = A_{ji}$$

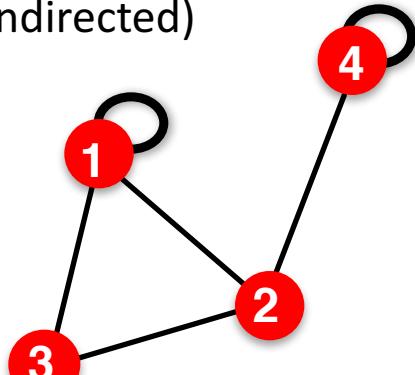
$$E = \frac{1}{2} \sum_{i,j=1}^N \text{nonzero}(A_{ij}) \quad \bar{k} = \frac{2E}{N}$$

Examples: Collaboration, Internet, Roads

# More Types of Graphs

## ■ Self-edges (self-loops)

(undirected)



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

$$A_{ii} \neq 0$$

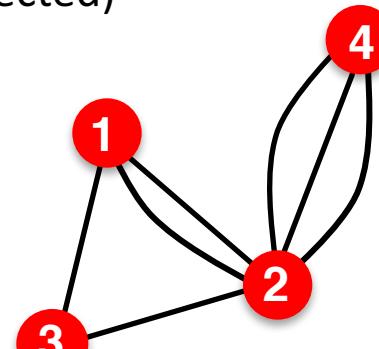
$$A_{ij} = A_{ji}$$

$$E = \frac{1}{2} \sum_{i,j=1, i \neq j}^N A_{ij} + \sum_{i=1}^N A_{ii}$$

Examples: Proteins, Hyperlinks

## ■ Multigraph

(undirected)



$$A_{ij} = \begin{pmatrix} 0 & 2 & 1 & 0 \\ 2 & 0 & 1 & 3 \\ 1 & 1 & 0 & 0 \\ 0 & 3 & 0 & 0 \end{pmatrix}$$

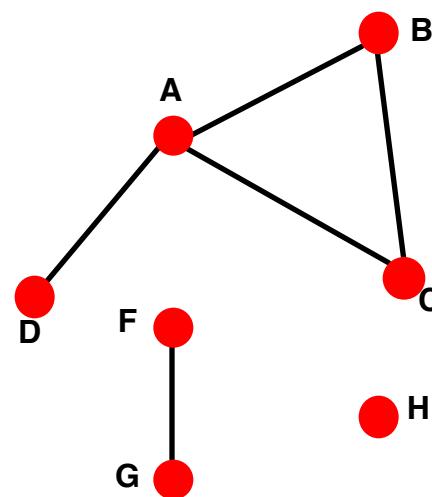
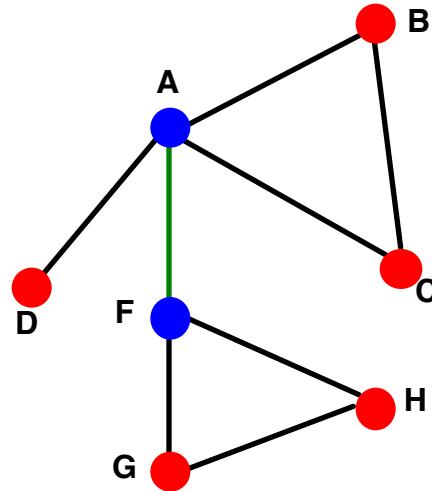
$$A_{ii} = 0$$

$$E = \frac{1}{2} \sum_{i,j=1}^N \text{nonzero}(A_{ij}) \quad \bar{k} = \frac{2E}{N}$$

Examples: Communication, Collaboration

# Connectivity of Undirected Graphs

- **Connected (undirected) graph:**
  - Any two vertices can be joined by a path
- A disconnected graph is made up by two or more connected components



Largest Component:  
**Giant Component**

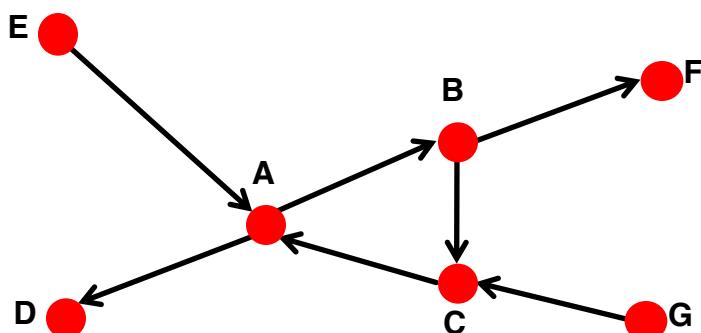
Isolated node (node H)

**Bridge edge:** If we erase the **edge**, the graph becomes disconnected.

**Articulation node:** If we erase the **node**, the graph becomes disconnected

# Connectivity of Directed Graphs

- **Strongly connected directed graph**
  - has a path from each node to every other node and vice versa (e.g., A-B path and B-A path)
- **Weakly connected directed graph**
  - is connected if we disregard the edge directions



Graph on the left is connected but not strongly connected (e.g., there is no way to get from F to G by following the edge directions).

# Network Representations

Email network >> directed multigraph with self-edges

Facebook friendships >> undirected, unweighted

Citation networks >> unweighted, directed, acyclic

Collaboration networks >> undirected multigraph or weighted graph

Mobile phone calls >> directed, (weighted?) multigraph

Protein Interactions >> undirected, unweighted with self-interactions