

# MIE524 Data Mining **Network Analysis**

Slides Credits:

Slides from Leskovec, Rajaraman, Ullman (<http://www.mmds.org>), Leskovec & Ghashami

# Announcements

- Thursday: quiz + new lab/assignment
- Midterm viewing during lab on Thursday (after quiz/lab)
- Brief review of midterm – tomorrow
- Final exam date announced
- Additional practice questions will be posted before the final exam
- Dedicated office hours before the final exam will be announced

# MIE524: Course Topics (Tentative)

## Large-scale Machine Learning

Learning Embedding  
(NN / AE)

Decision Trees

Ensemble Models  
(GBTs)

## High-dimensional Data

Locality sensitive hashing

Clustering

Dimensionality reduction

## Graph Data

Processing Massive Graphs

PageRank, SimRank

Graph Representation Learning

## Applications

Recommender systems

Association Rules

Neural Language Models

Computational Models:

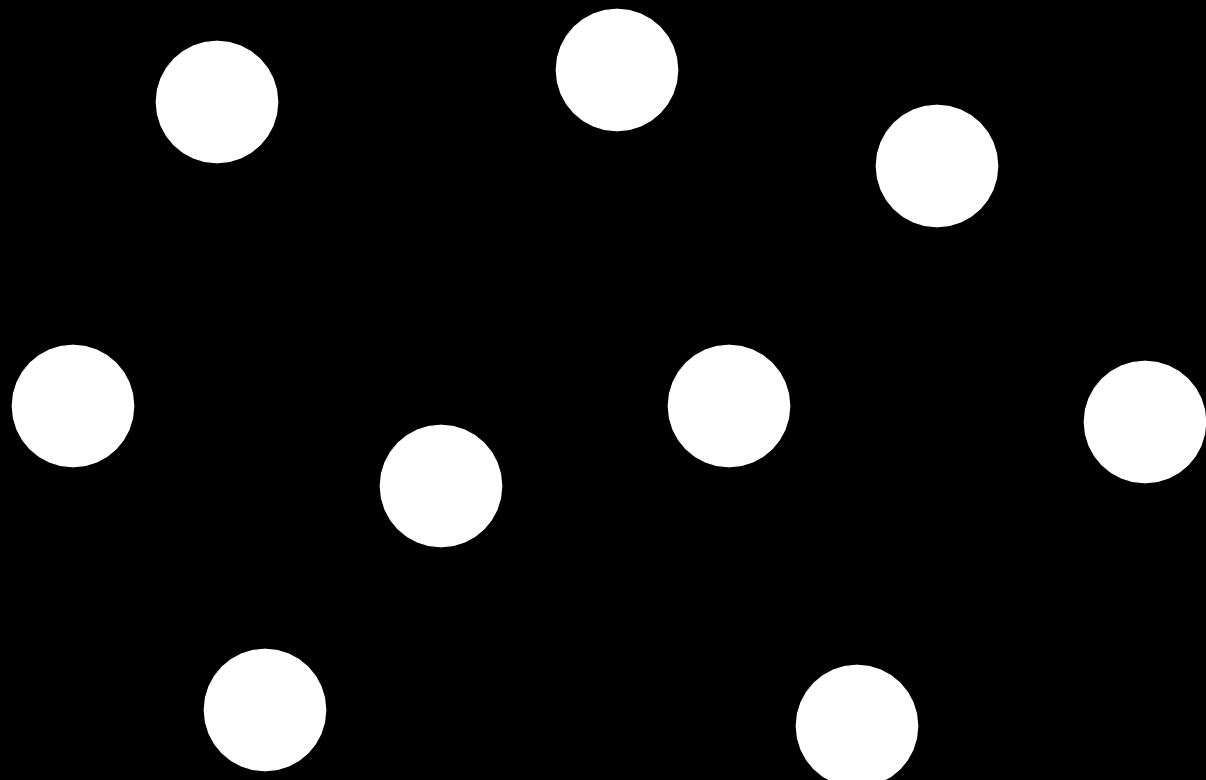
Single Machine

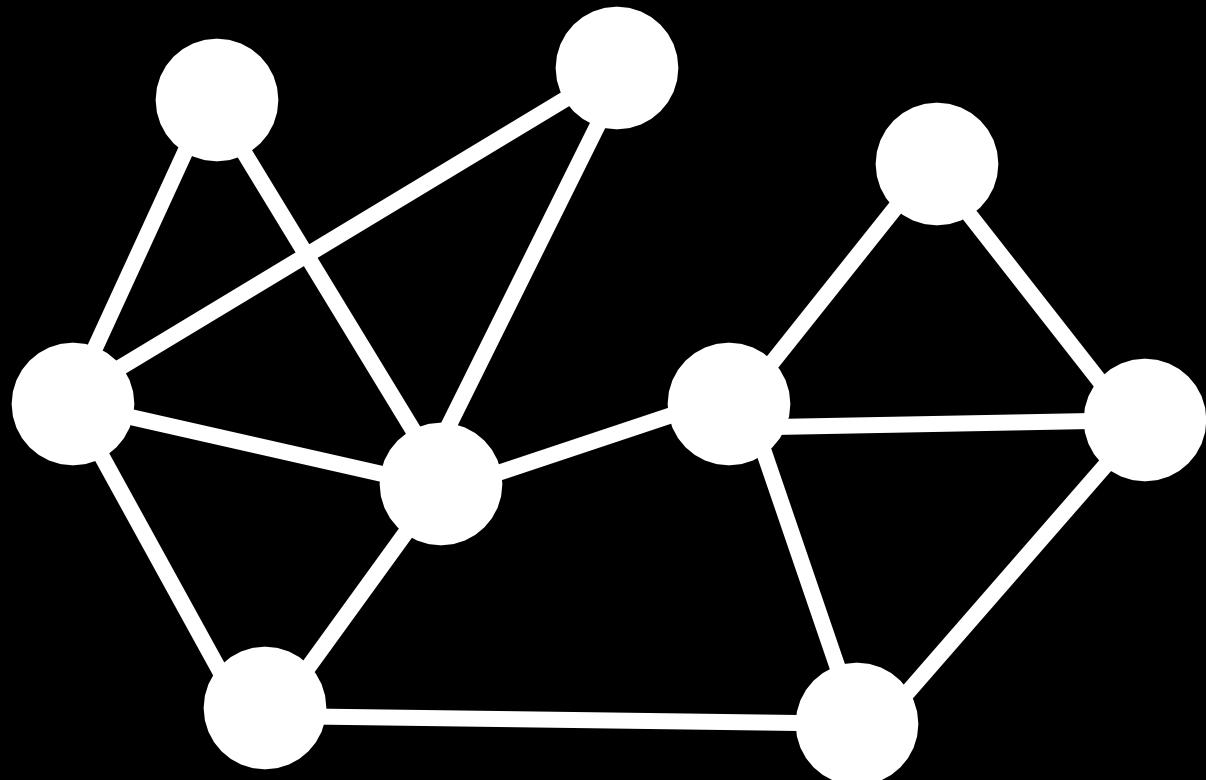
MapReduce/Spark

GPU

# 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?**

# Networks!!

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

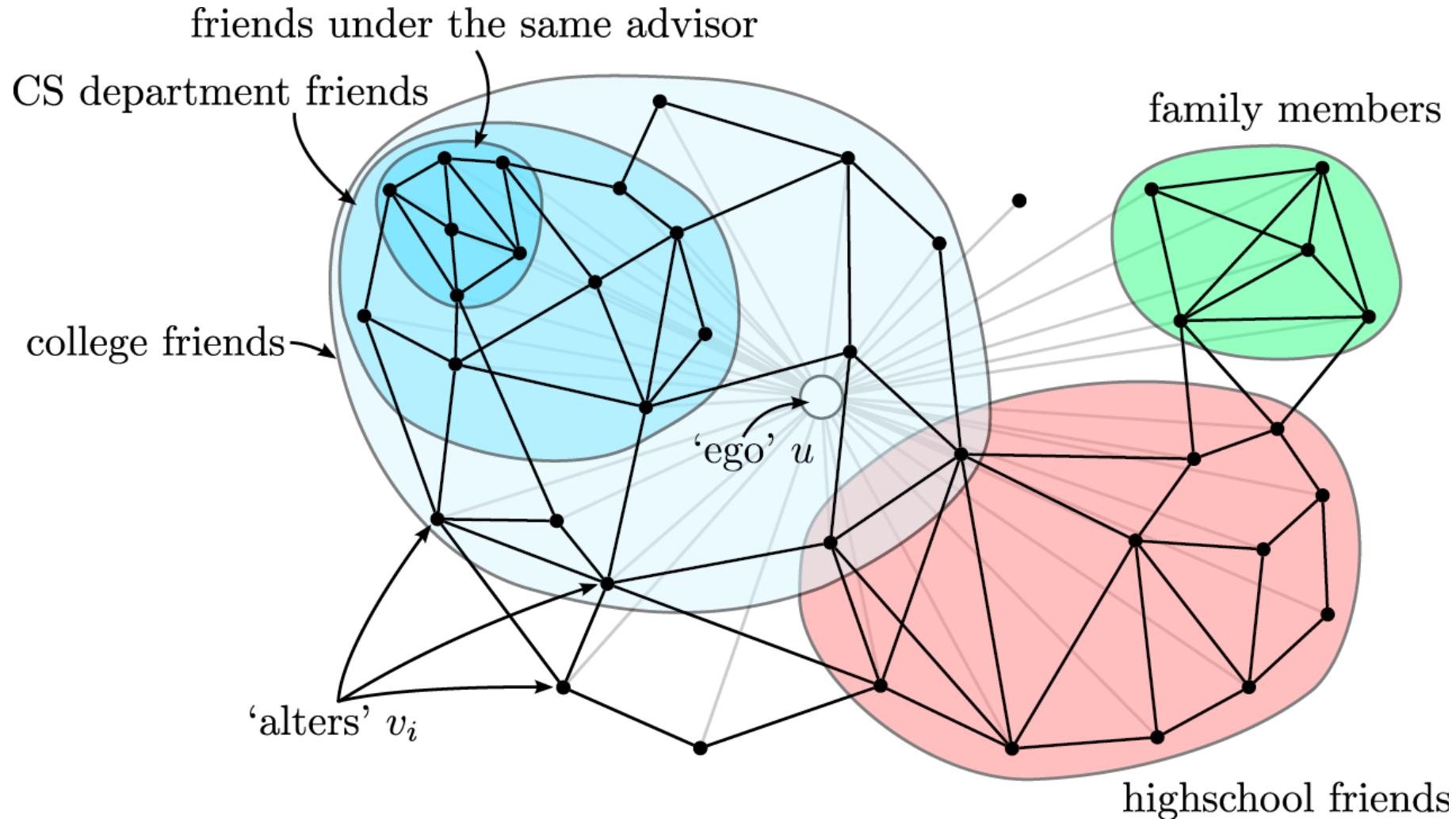
**We will never understand these systems unless we understand the networks behind them!**

# 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 and Applications**

# (1) Networks: Social



## Discover circles and why they exist

# (2) Networks: Infrastructure



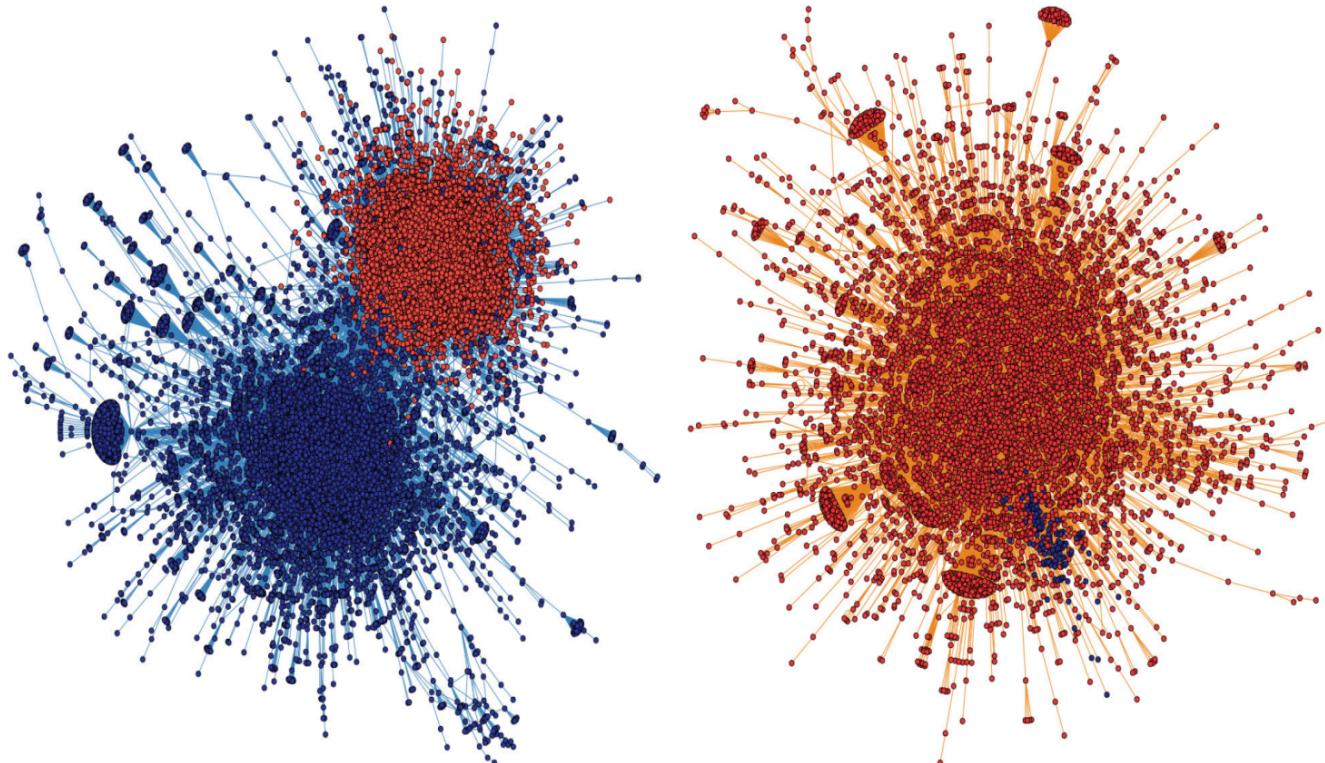
Water supply distribution networks



Airline networks

# (3) Networks: Online Media

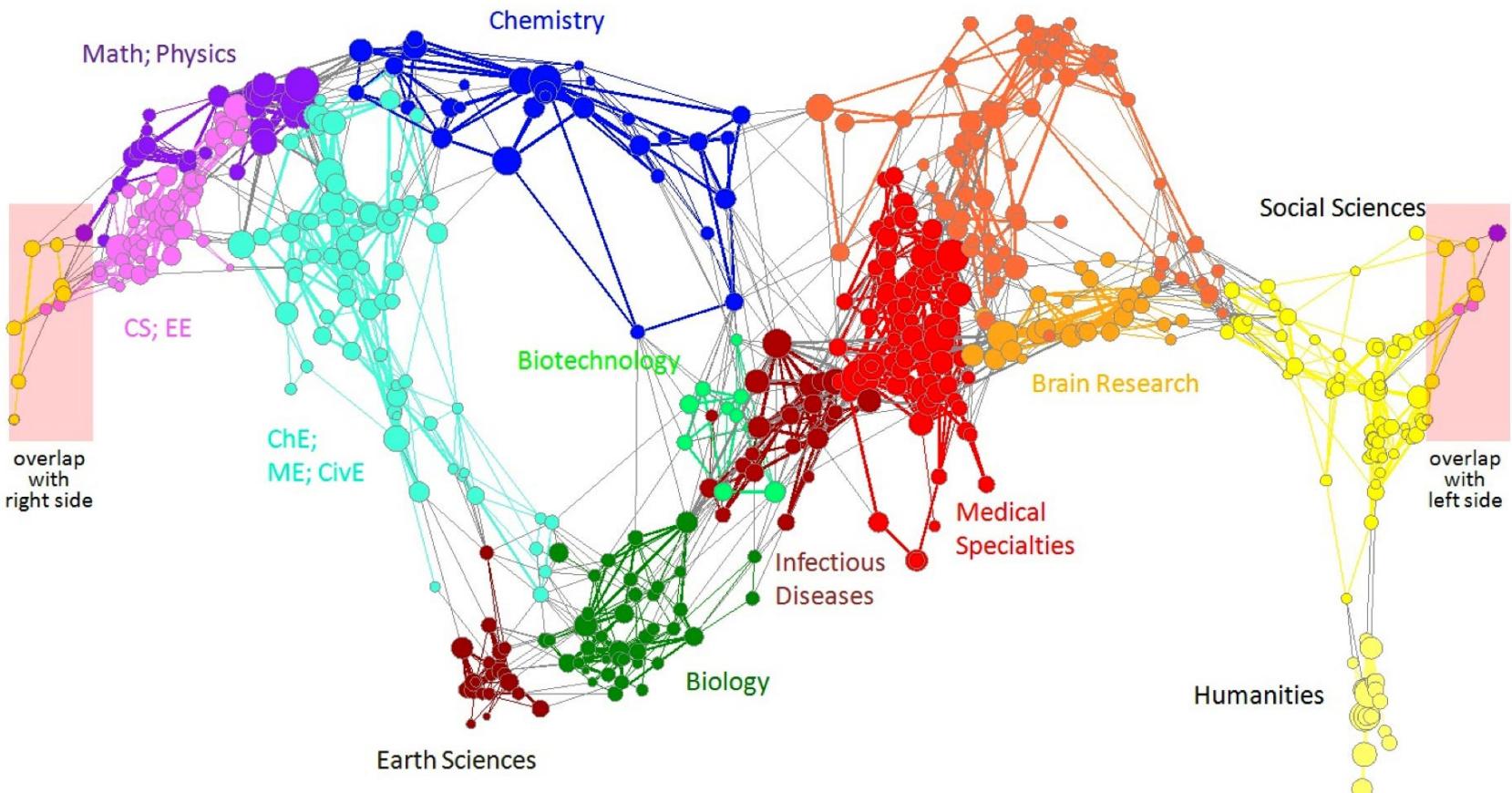
## Polarization on Twitter



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

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

# (4) Networks: Information, Knowledge

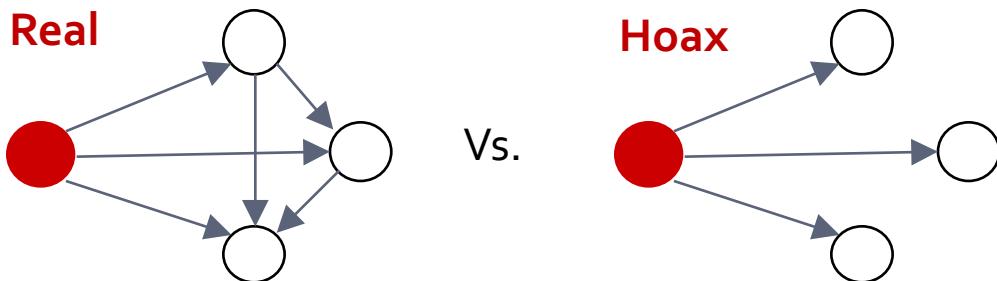


Citation networks and maps of science  
[Börner et al., 2012]

# Application: Misinformation

- Q: Is a given Wikipedia article a hoax?

- Real articles link more coherently:



Hoax article detection performance:

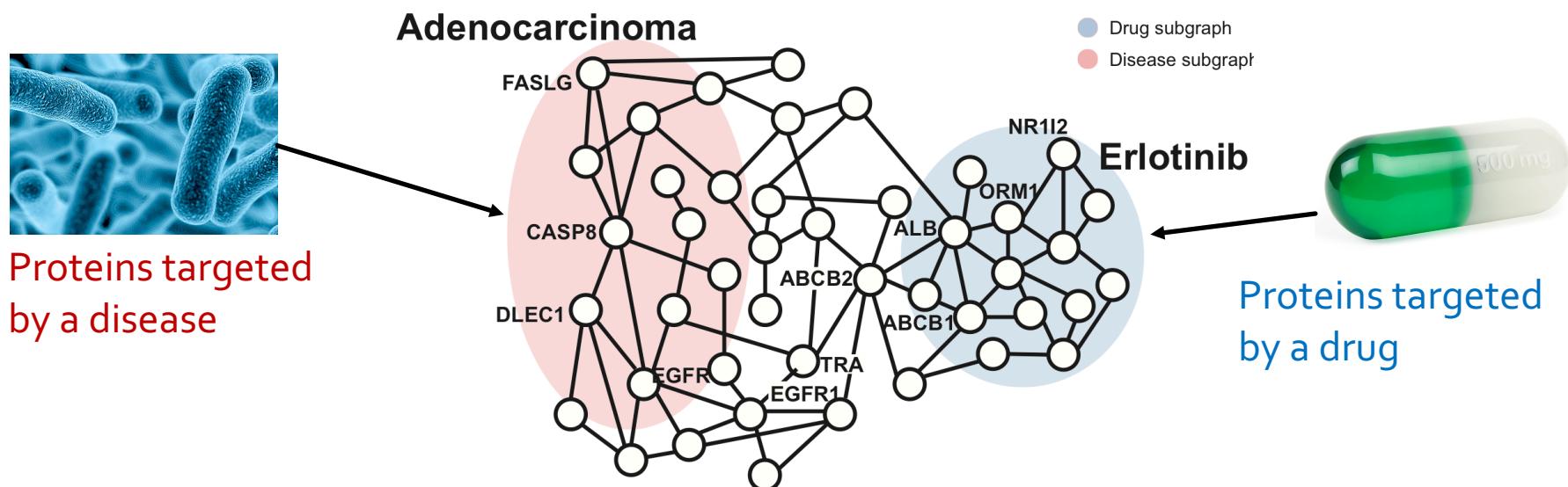
50%	66%	86%
Random	Human	WWW '16

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

A screenshot of a Wikipedia page titled "Wikipedia:List of hoaxes on Wikipedia/Balboa Creole French". The page shows 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)". To the right, there is a detailed sidebar for "Balboa Creole French" with sections for Native to, Region, Native speakers, Language family, and Language codes. The sidebar notes that it is a Creole language used in Balboa Island, California, and is virtually extinct, with only 14 people remaining who can speak it.

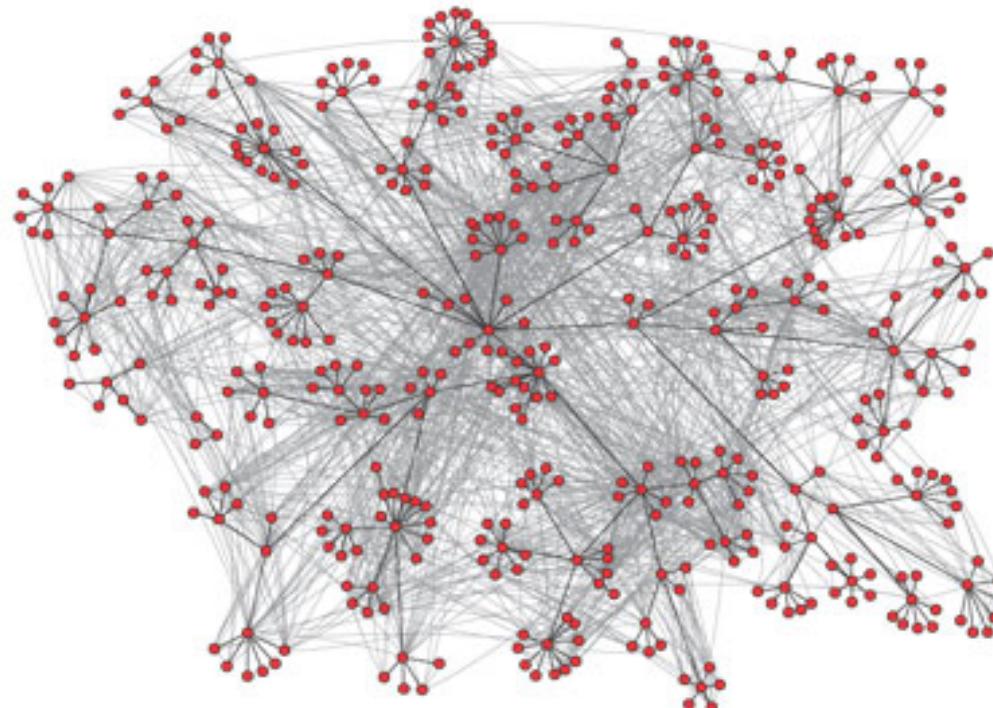
# Application: Drug Repurposing

- 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.



# **Structure of Networks**

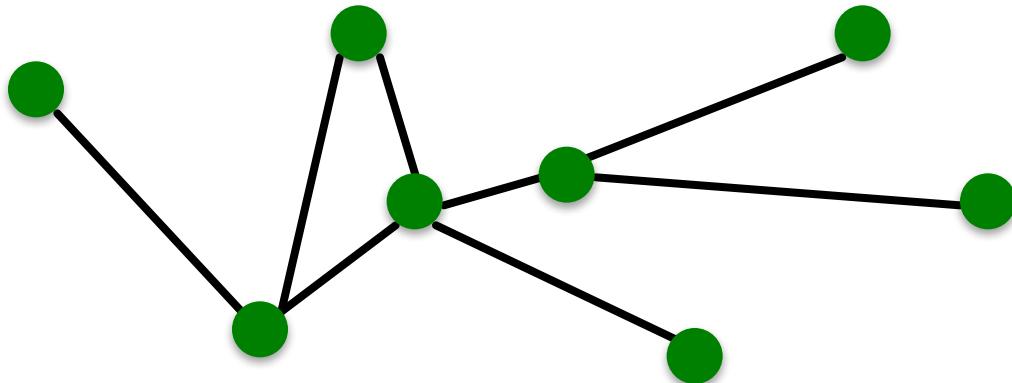
# Structure of Networks?



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?

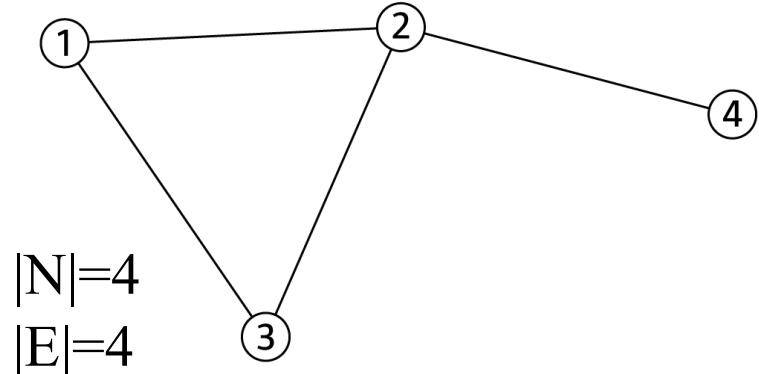
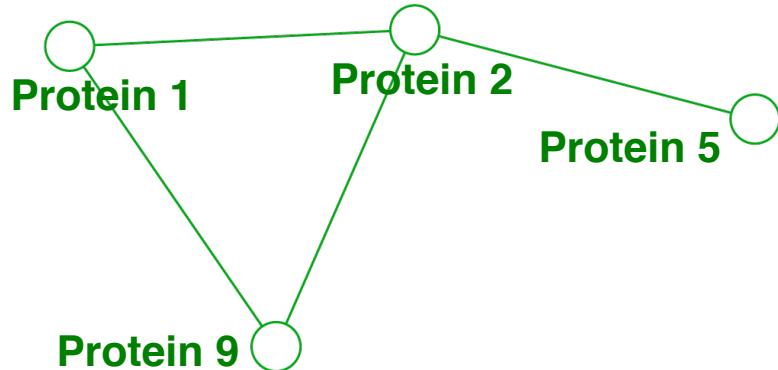
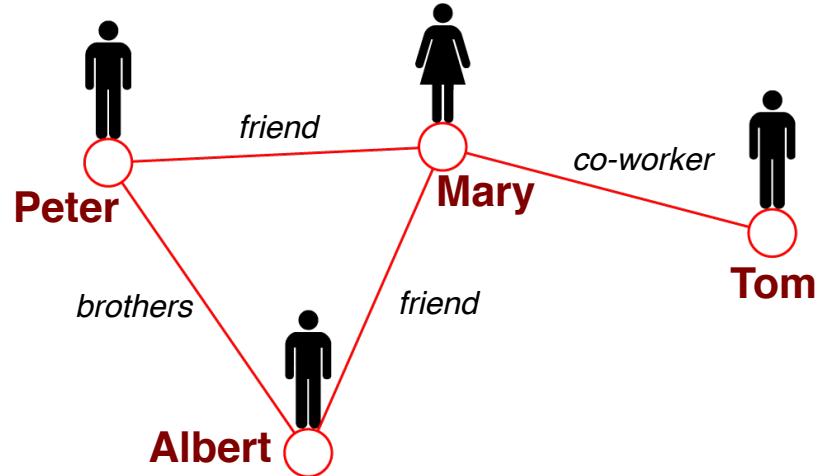
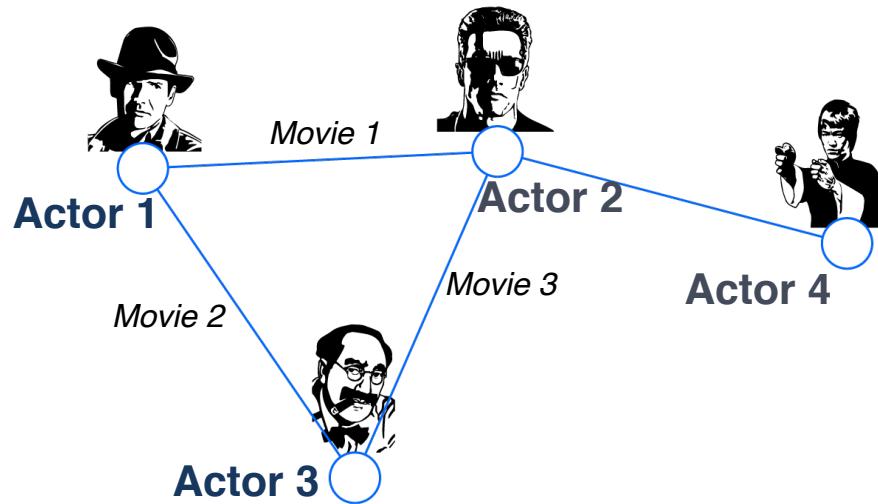
- **Network** often refers to real systems
  - Web, Social network, Metabolic network

**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



# How do you define a network?

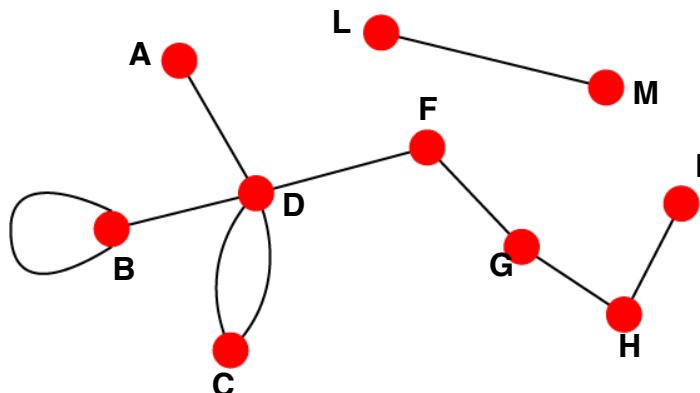
- **How to build a graph:**
  - What are nodes?
  - What are edges?
- **Choice of the proper network representation of a given domain/problem determines our ability to use networks successfully:**
  - In some cases there is a unique, unambiguous representation
  - In other cases, the representation is 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

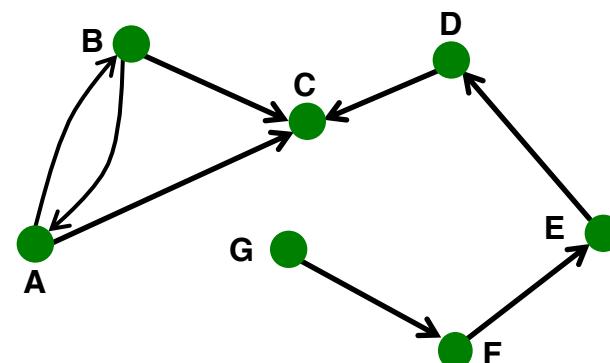
## Undirected

- Links: undirected  
(symmetrical, reciprocal)



## Directed

- Links: directed  
(arcs)



## Examples:

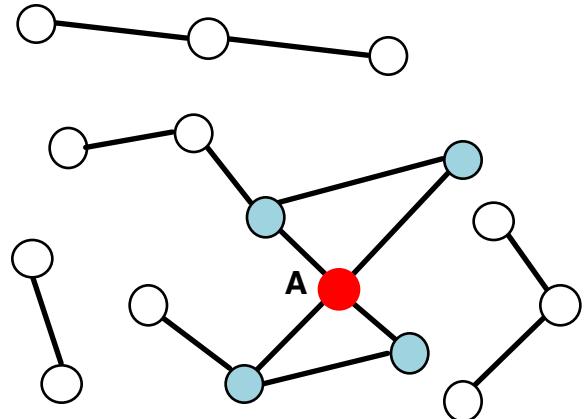
- Collaborations
- Friendship on Facebook

## Examples:

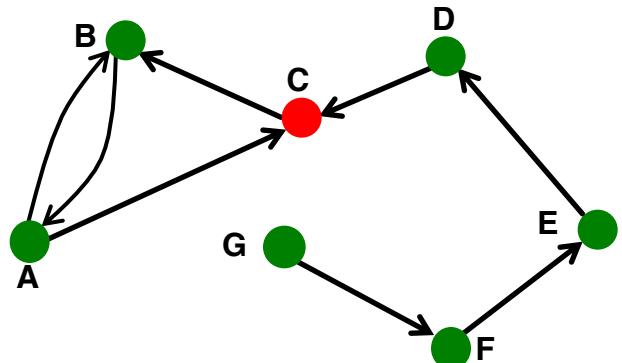
- Phone calls
- Following on Twitter

# 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.

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

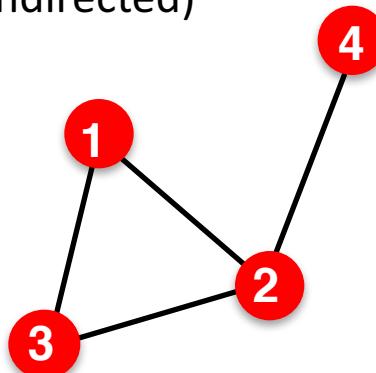
$$\bar{k} = \frac{E}{N}$$

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

# More Types of Graphs

## ■ Unweighted

(undirected)



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

$$A_{ii} = 0$$

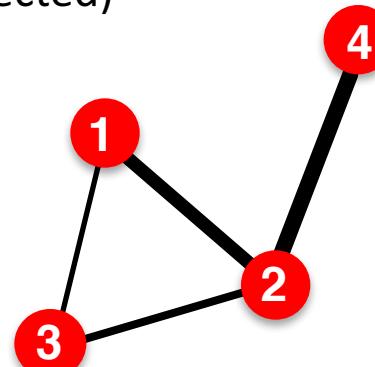
$$A_{ij} = A_{ji}$$

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

Examples: Friendship, Hyperlink

## ■ Weighted

(undirected)



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

$$A_{ii} = 0$$

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

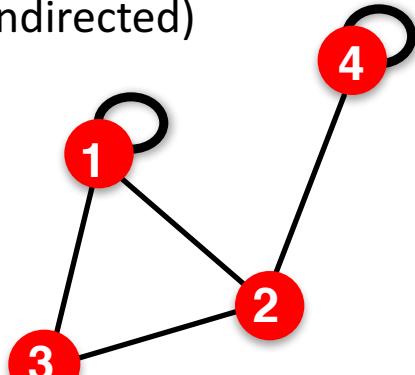
$$E = \frac{1}{2} \sum_{i,j=1}^N \text{nonzero}(A_{ij}) \quad \bar{k} = \frac{2E}{N}$$

Examples: Collaboration, Internet, Roads

# More Types of Graphs

## ■ Self-edges (self-loops)

(undirected)



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

$$A_{ii} \neq 0$$

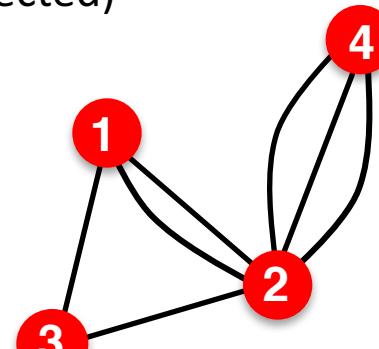
$$A_{ij} = A_{ji}$$

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

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

# Bipartite Graph

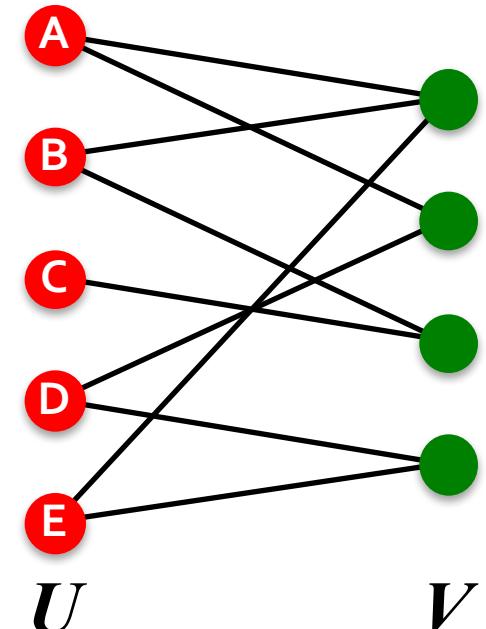
- **Bipartite graph** is a graph whose nodes can be divided into two disjoint sets  $U$  and  $V$  such that every link connects a node in  $U$  to one in  $V$ ; that is,  $U$  and  $V$  are **independent sets**

- **Examples:**

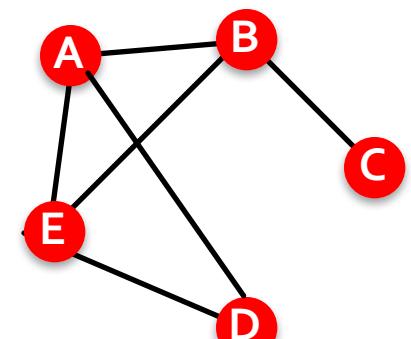
- Authors-to-papers (they authored)
- Actors-to-Movies (they appeared in)
- Users-to-Movies (they rated)

- **“Folded” networks:**

- Author collaboration networks
- Movie co-rating networks



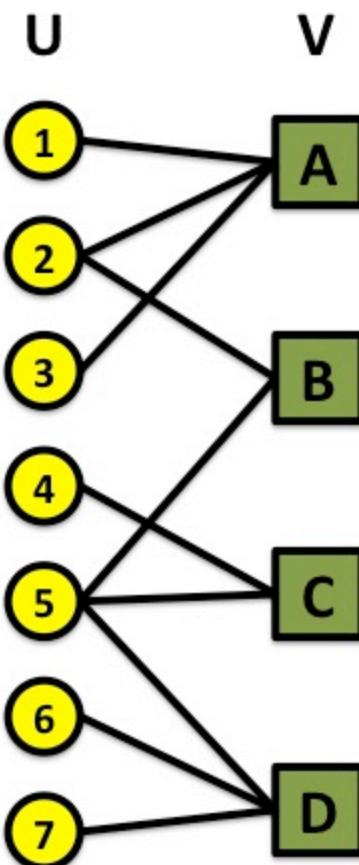
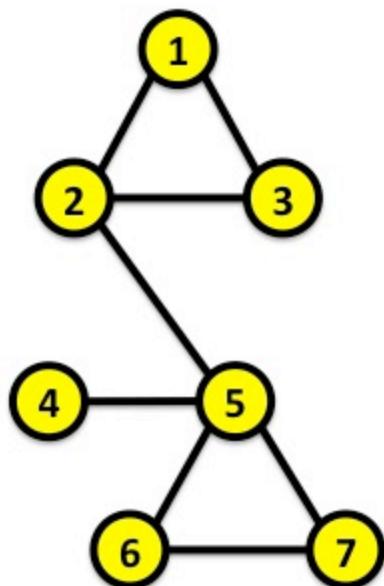
$U$        $V$



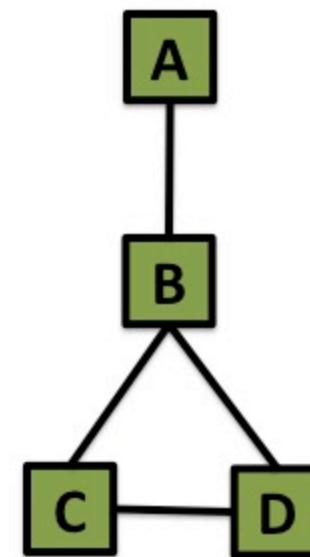
Folded version of the graph above

# Folded/Projected Bipartite Graphs

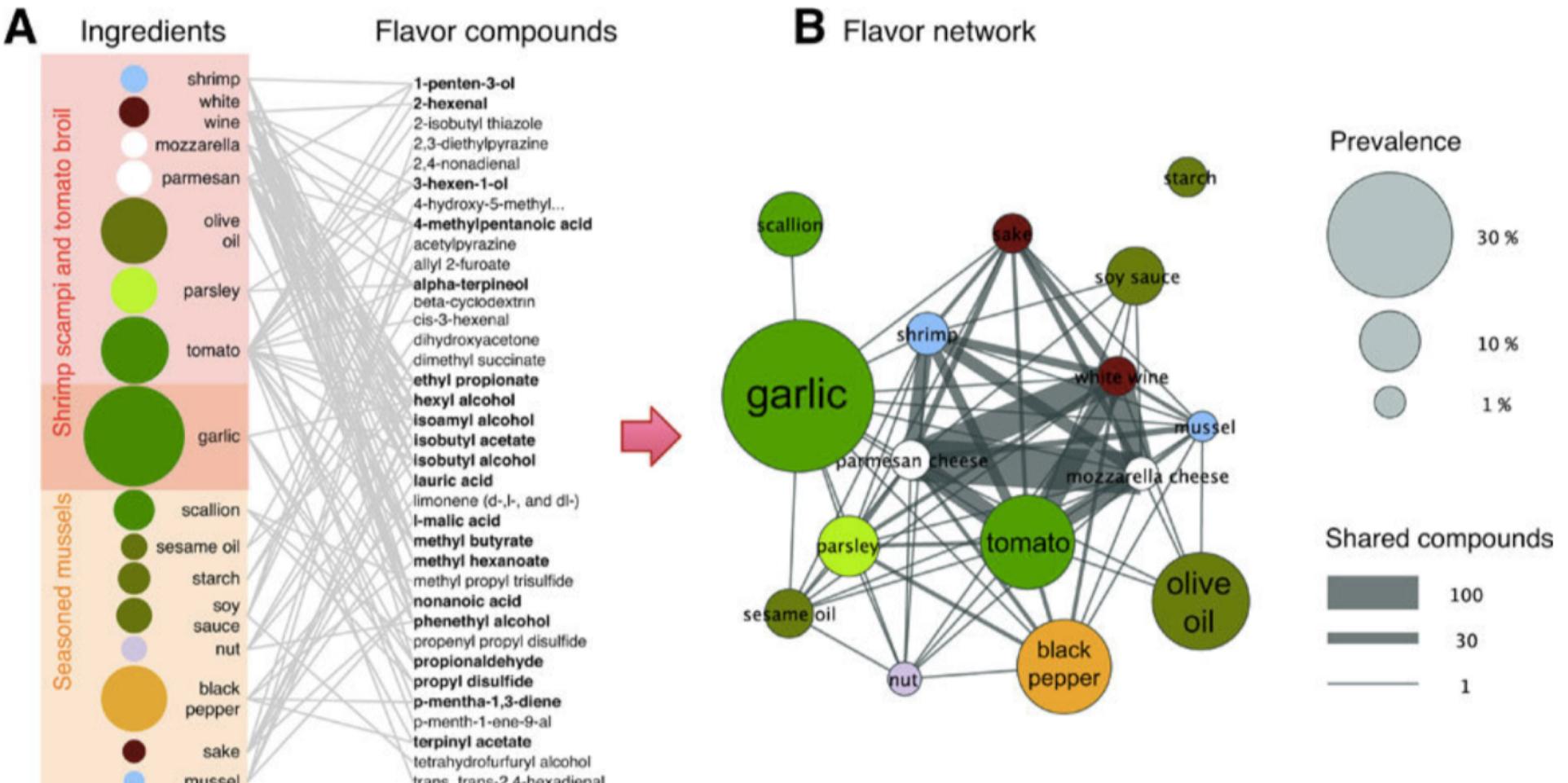
Projection U



Projection V



# Example: Ingredients and Flavors

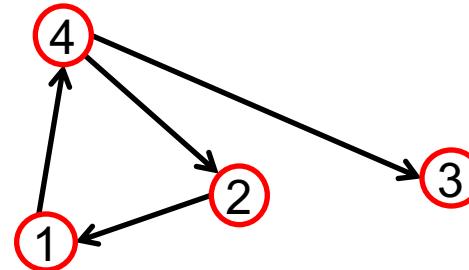
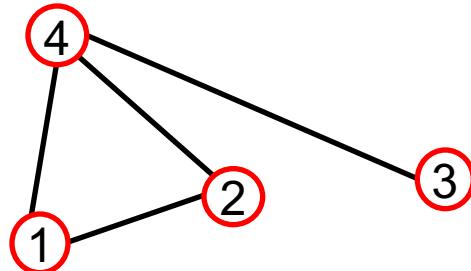


Y.-Y. Ahn, S. E. Ahnert, J. P. Bagrow, A.-L. Barabási

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

Network Science: Graph Theory

# 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.

# 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