# Social Media Analysis (1 of 4)

Lexing Xie
Computer Science, ANU

# Networks … of Documents
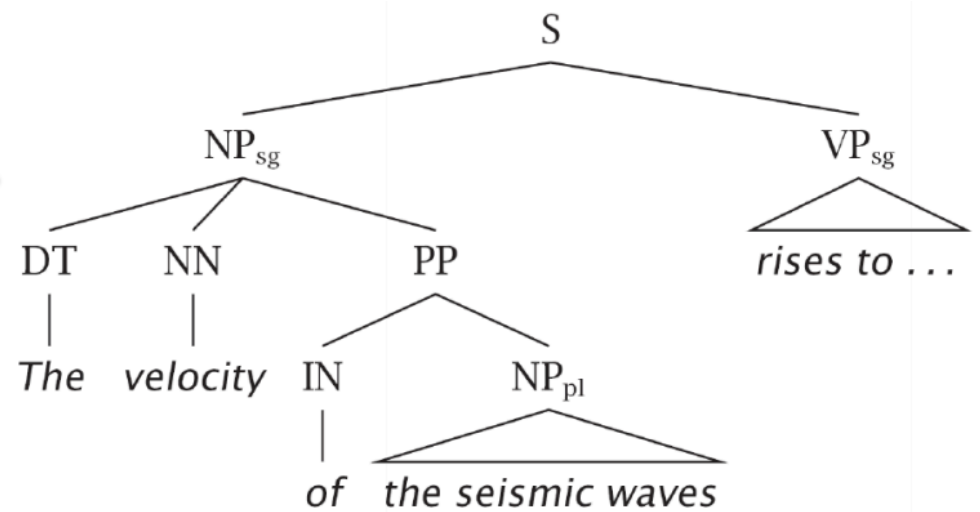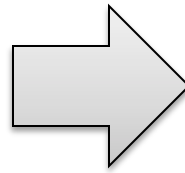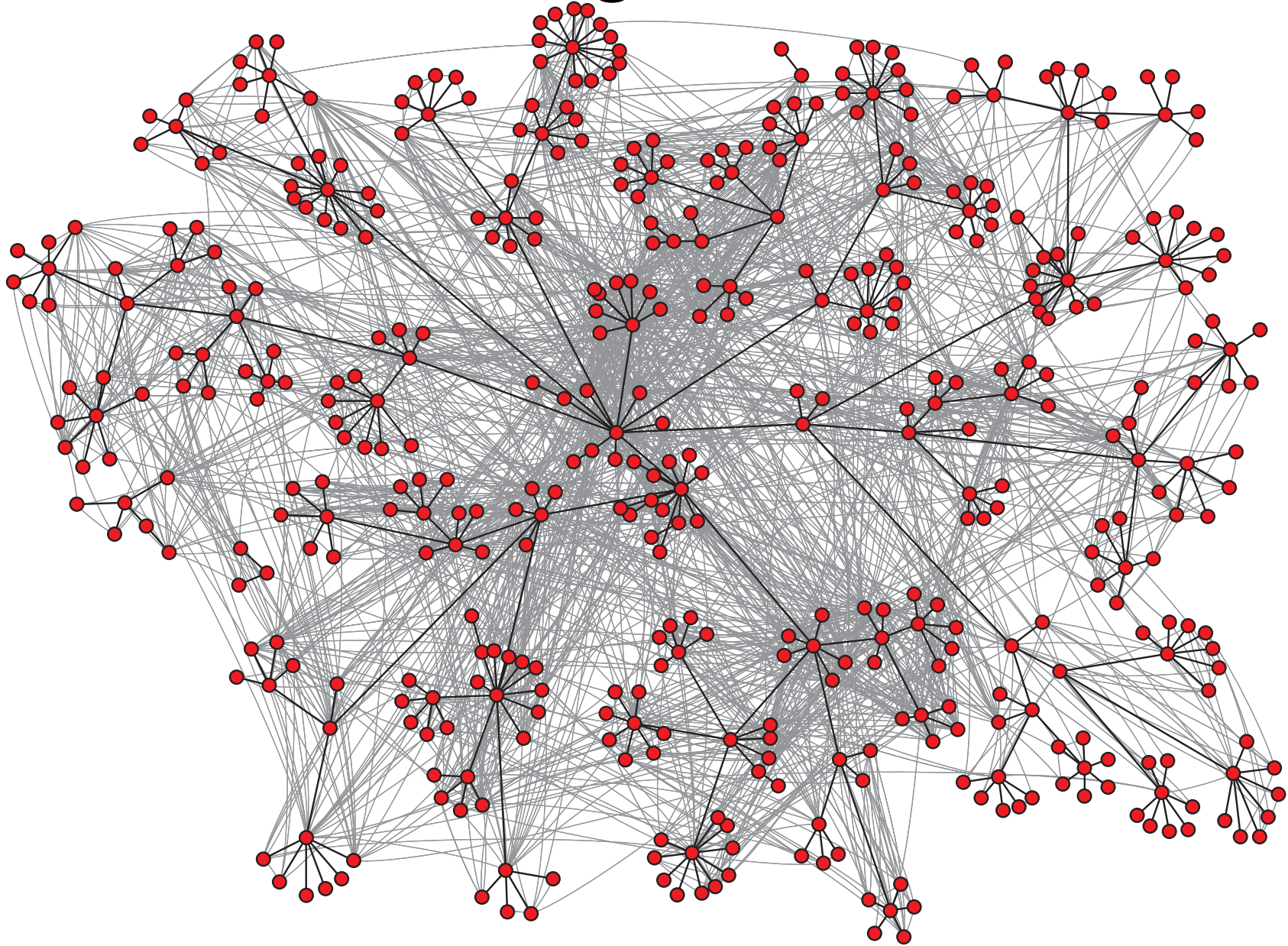
- The rest of this class covers
  - Extracting document elements: nouns/verbs, subjects/object, web page structure …
  - Labelling and matching documents: IR, document classification, clustering …

- What if we look at …
  - Documents and their relations to each other?
  - Document elements and their relations to each other?
  - People and their relations to each other

# Earlier in This Course …

Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man prüfen, wie gut die Sch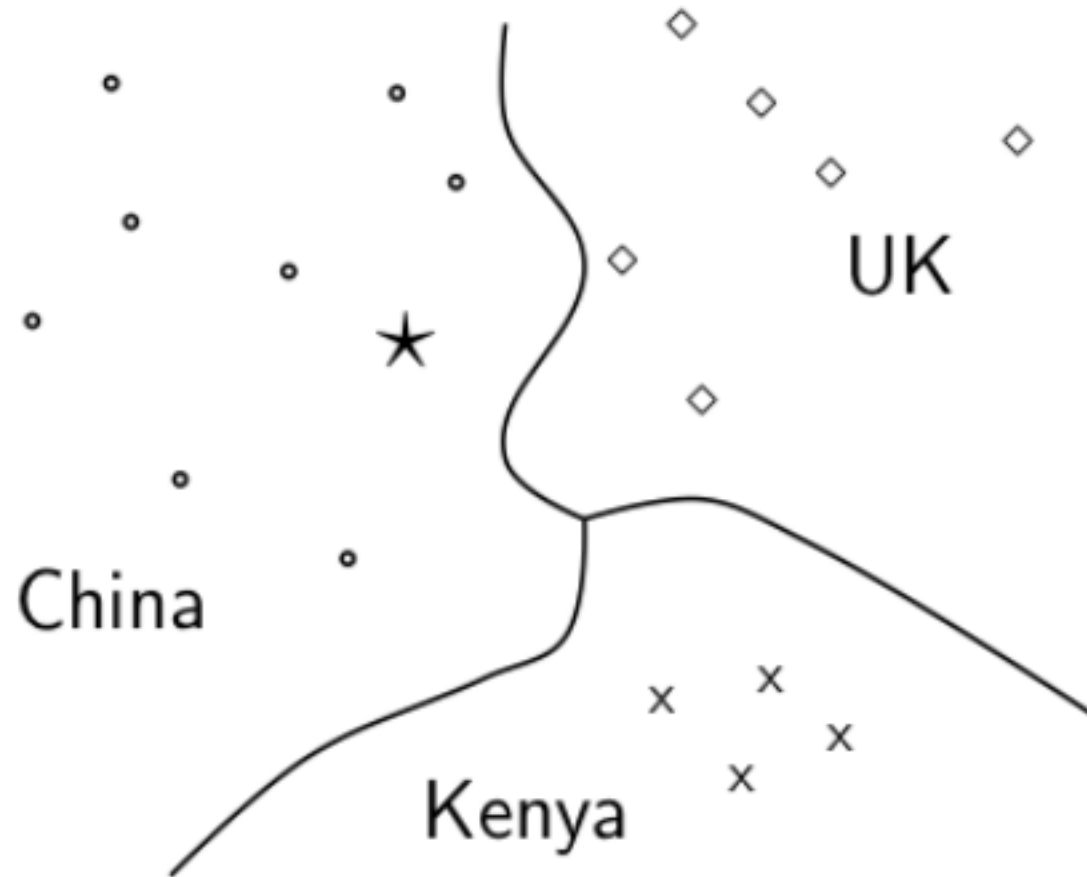rift zu lesen ist und wie sie auf den Leser wirkt. Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man prüfen, wie gut die Schrift zu lesen ist und wie sie auf den Leser wirkt.
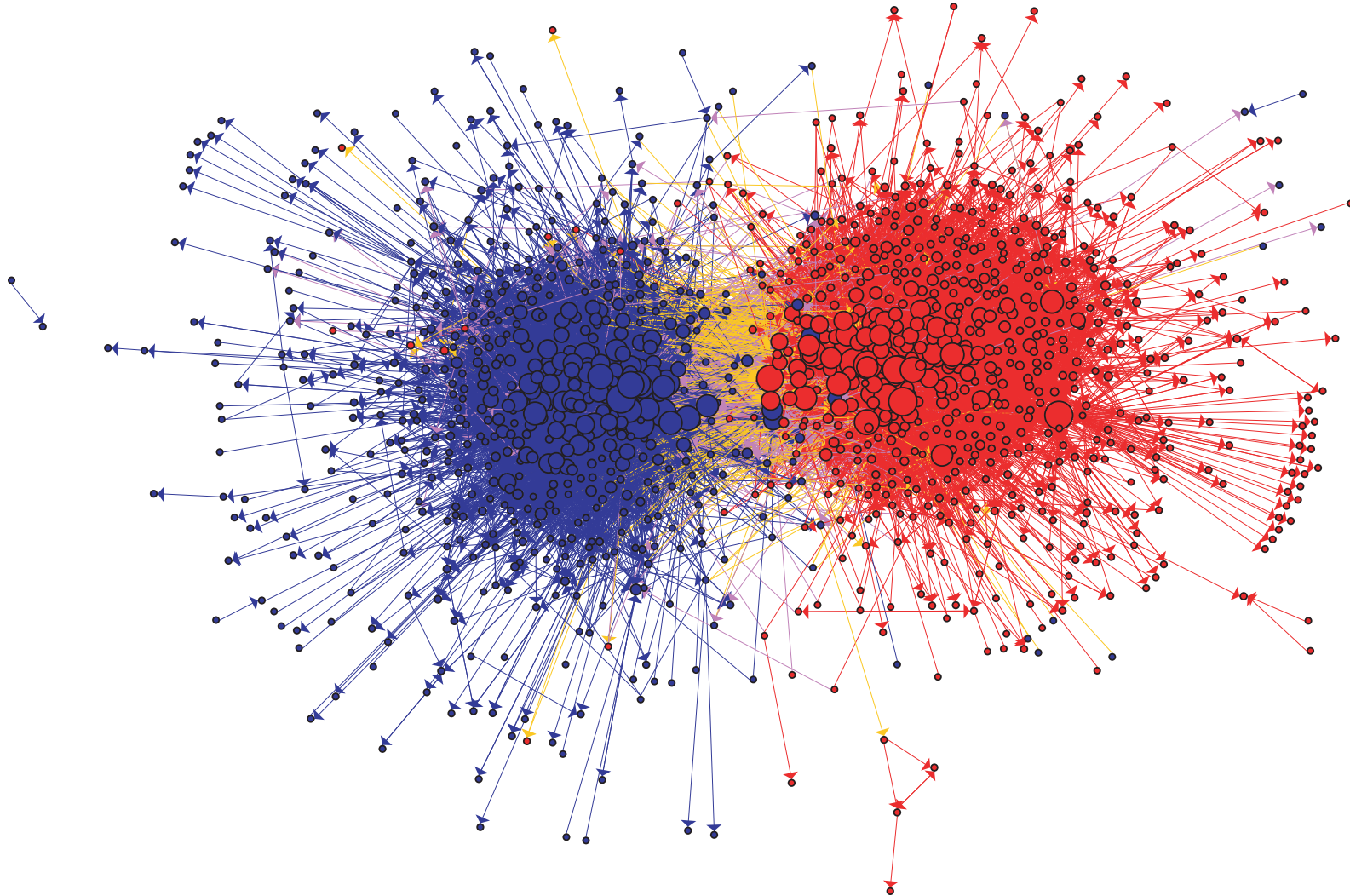
# Networks: An organization - HP Labs

# Earlier in this Course: classifying documents

# Political blogs



Lada A. Adamic and Natalie Glance. 2005. The political blogosphere and the 2004 U.S. election: divided they blog. In Proceedings of the 3rd international workshop on Link discovery (LinkKDD '05)    by Lada Adamic, U Michigan
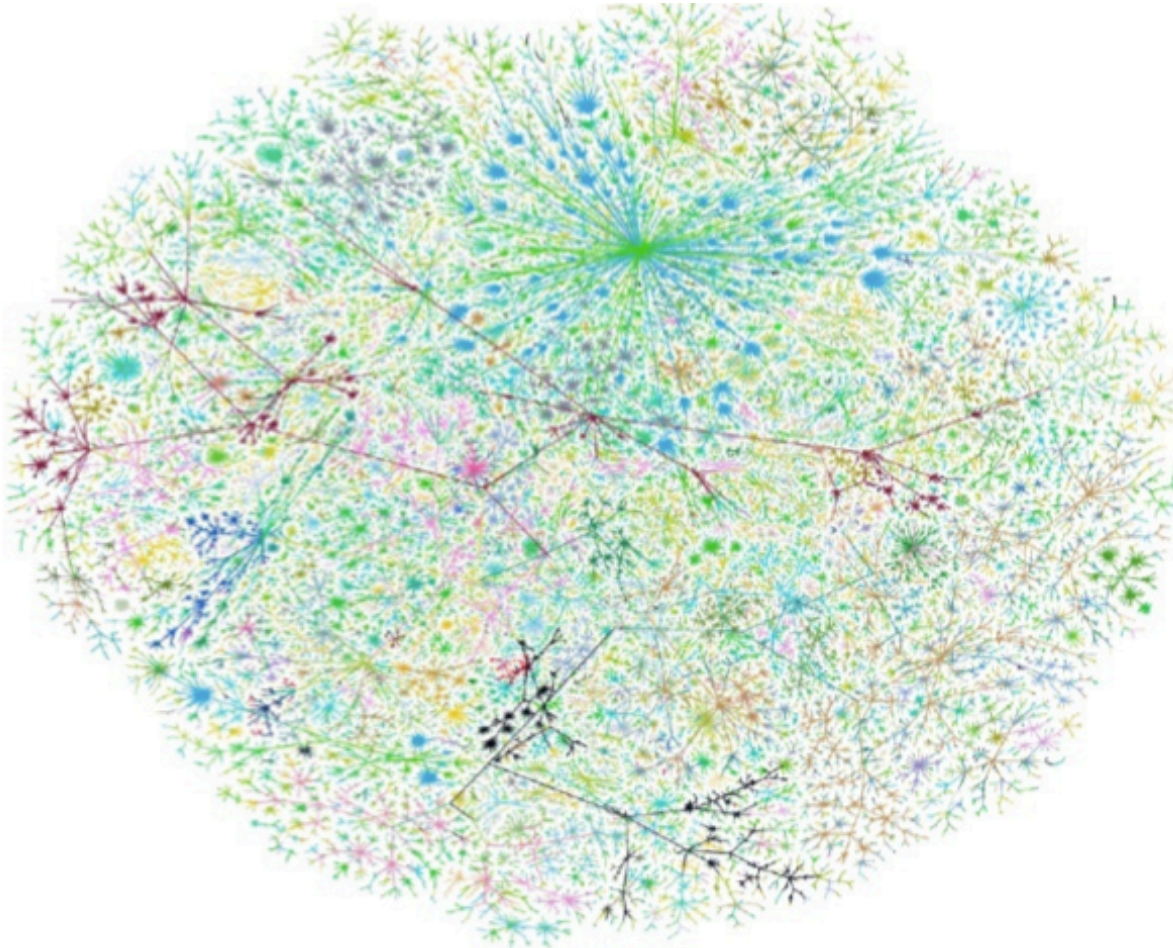
# A View of Facebook via 10 M links



By Paul Butler, https://www.facebook.com/note.php?note_id=469716398919

# Why Networks?

- Behind each of these complex systems there is an intricate wiring diagram,
  a network
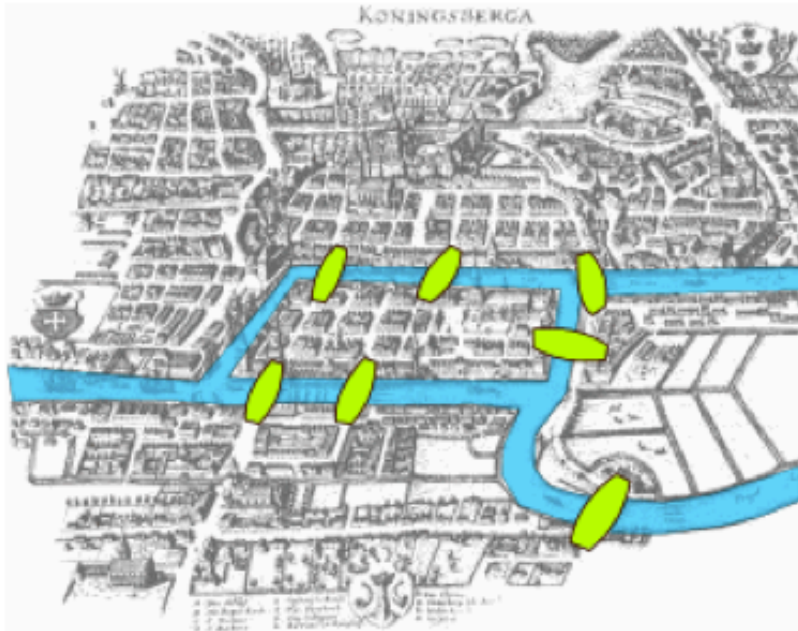  that defines the interactions between system components.

Understanding the network is key to understand the behaviors of such complex system.

# Networks: Communications
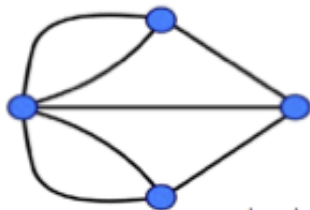


**Graph of the Internet
(Autonomous Systems)**

# Networks: Transportation



**Seven Bridges of Königsberg
(Euler 1735)**

Return to the starting point by traveling each
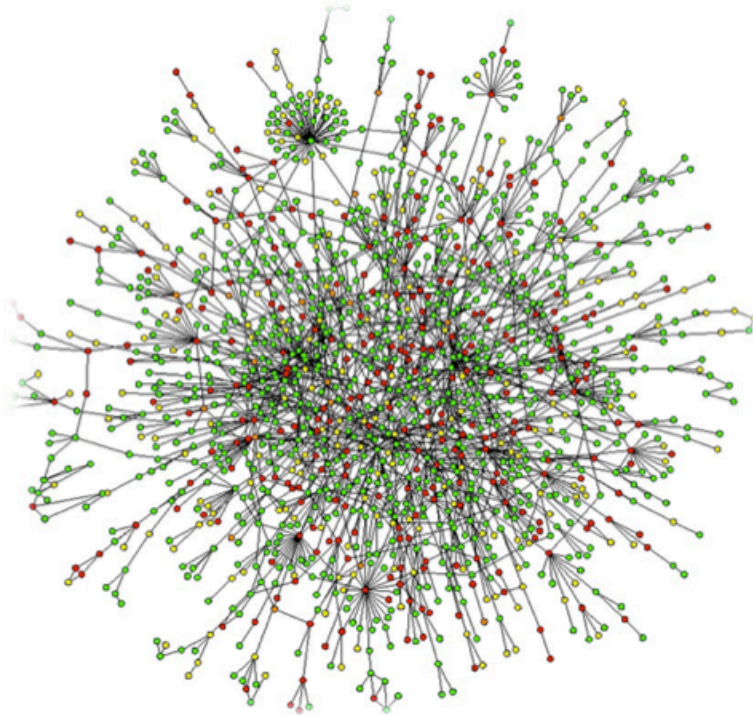link of the graph once and only once.



**London Underground**

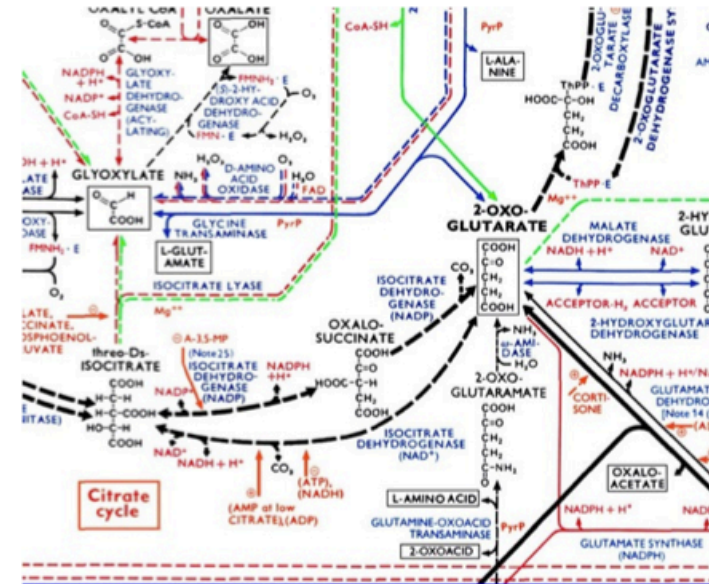# Networks: Brain



**Human brain has between
10-100 billion neurons**

# Networsk: Cells



**Protein-Protein Interaction Networks:**
Nodes: Proteins
Edges: 'physical' interactoins



**Metabolic networks:**
Nodes: Metabolites and enzymes
Edges: Chemical reactions

# Ingredient networks



Recipe recommendation using ingredient networks. Chun-Yuen Teng, Yu-Ru Lin, Lada A. Adamic, WebSci 2011.

# Why Networks?

- Behind each of these complex systems there is an intricate wiring diagram,
  a network
  that defines the interactions between system components.

Understanding the network is key to understand the behaviors of such complex system.

# Application domains in network analysis

This class

- Social (people-people) networks ✔
- Information networks ✔
- Organization and political networks ✔


- Computer networking
- Biology
- Transportation networks

# CSS: Computational challenges

This class

- Machine Learning and applied statistics:
  Predictive modeling, models for network structure,
  models for dynamic events, optimization ✔

- Natural Language Processing:
  linguistic styles and variations, document similarity, ✔
  information retrieval

- Visual analytics, visualization ✔

- Computer systems: scalability, reliability,
  programming languages and tools ✗

- …

## Visualizing Friendships

by Paul Butler for Facebook Engineering (Notes) on Tuesday, December 14, 2010 at 10:16am
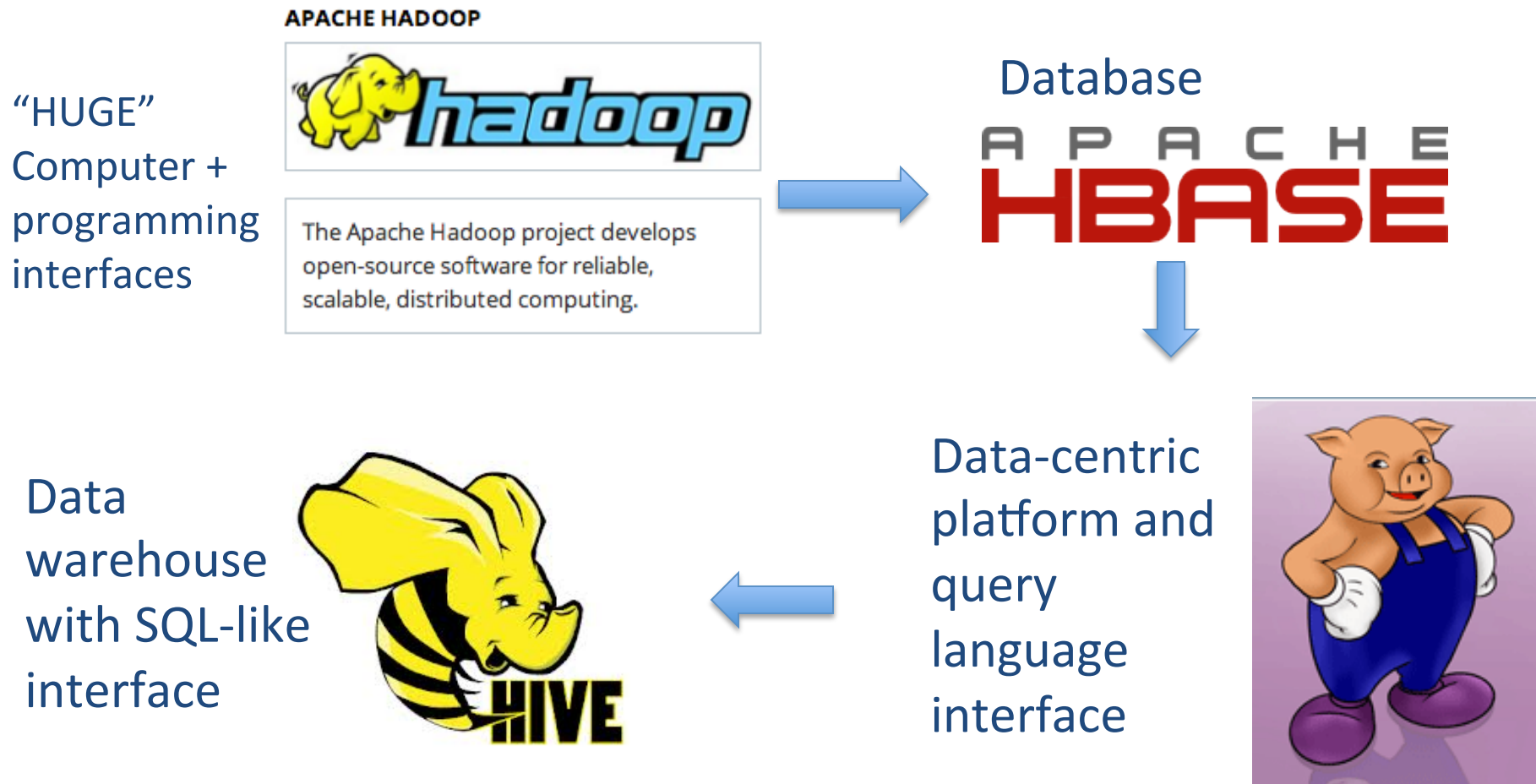
Visualizing data is like photography. Instead of starting with a blank canvas, you manipulate the lens used to present the data from a certain angle.



1. Sample 10 million links from Hive
2. Get city/lat/lon for each node
3. Compute connection strength between cities
4. Render links among city pairs in R
5. Tweak rendering and path drawing until satisfactory

https://www.facebook.com/note.php?note_id=469716398919

# Storing and Manipulating 500M Links



"HUGE" Computer + programming interfaces

**APACHE HADOOP**

The Apache Hadoop project develops open-source software for reliable, scalable, distributed computing.

Database

**APACHE HBASE**

Data-centric platform and query language interface

Data warehouse with SQL-like interface

**HIVE**

http://developer.yahoo.com/blogs/hadoop/posts/2010/08/pig_and_hive_at_yahoo/

# Why Social Computing?

* The next generation could be the one with access to an unprecedented amount of behavioral data

* This can solve real problems

  ... not just finding a movie or a restaurant

  o ensuring energy efficiency
  o monitoring our environment
  o reduce inequality
  o informing social decision

7

# Why should you learn about it? #1

* These concepts matter to the companies that you want a job from (including your startup!)
  o Social information is becoming the web's hottest commodity (Google, Facebook, IBM, Telcos, Media)
  o Users's data are company's key differentiating factor
* You (not me) are the social media generation!
  o The game is just starting; it gets harder
  o CS deals with "complexity" deeply and elegantly
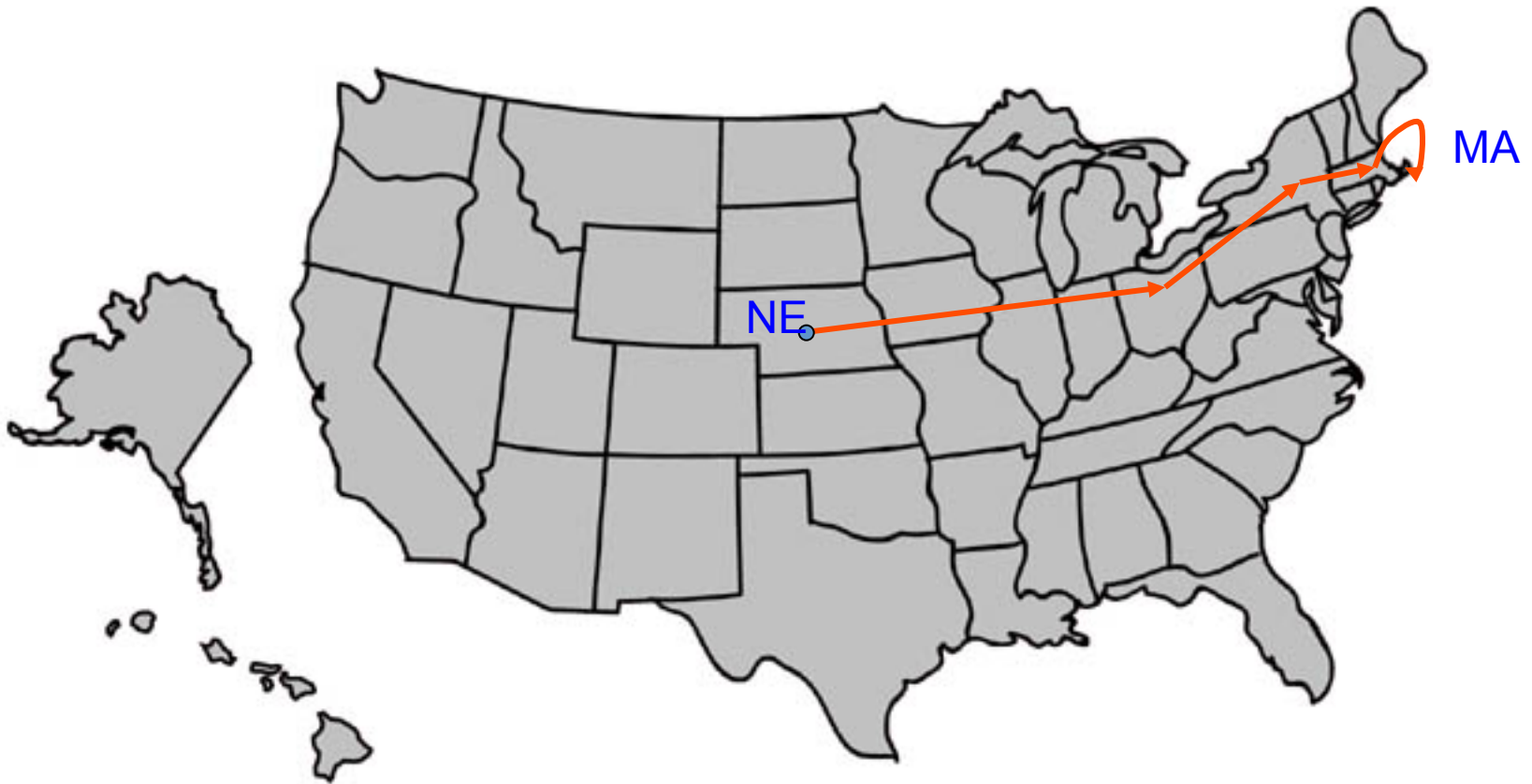  o learning foundational concepts adds to your assets

# Why should you learn about it? #2

\* This topic is fertile for research, here and at large
  - o Big data is everywhere, especially in public-funding
  - o Many of these data are networks connecting people

- Many opportunities in/around ANU
  - Big data research theme at CS
  - Machine learning group at NICTA
  - Master degree with Network Science concentration at CASS
  - Overall, great topic to look for an academic job!

What are the classic, solve, and open questions in CSS?

# PROBLEMS IN SOCIAL NETWORK

# Small-world phenomena: Milgram's Experiment



http://en.wikipedia.org/wiki/Small_world_experiment          1967 -- 1969

# Milgram's experiment

**Instructions:**
Given a target individual (stockbroker in Boston), pass the message to a person you correspond with who is "closest" to the target.

**Outcome:**

**20% of initiated chains reached target**
**average chain length = 6.5**

- "Six degrees of separation"

# Milgram's experiment repeated

email experiment
Dodds, Muhamad, Watts,
Science 301, (2003)
         (optional reading)

- 18 targets
- 13 different countries

- 60,000+ participants
- 24,163 message chains
- 384 reached their targets
- average path length 4.0

Slide by Lada Adamic, U Michigan

- Is 6 is a *surprising* number?

  In the 1960s? Today? Why?


  What is the mechanism behind "small-world" networks?

# Link prediction

"Given a snapshot of a social network, can we infer which new interactions among its members are likely to occur in the near future? "

What are the measurements?
-- "proximity" and "similarity" between two unconnected nodes
What are application domains?
-- social networks, friend recommendation; product / webpage recommendation;  predicting academic collaborations; predicting merger and acquisitions ...
How to measure performance?
-- use future held-out data, conversion rate, ...

Liben-Nowell, D. and Kleinberg, J. The link-prediction problem for social networks. Journal of the American Society for Information Science and Technology, 58(7) 1019{1031 (2007)

# Language Use in Social Media



Echoes of power:  Language effects and power differences in social interaction
Cristian Danescu-Niculescu-Mizil, Lillian Lee, Bo Pang and Jon Kleinberg.
Proceedings of WWW, 2012.

# Tracking Memes

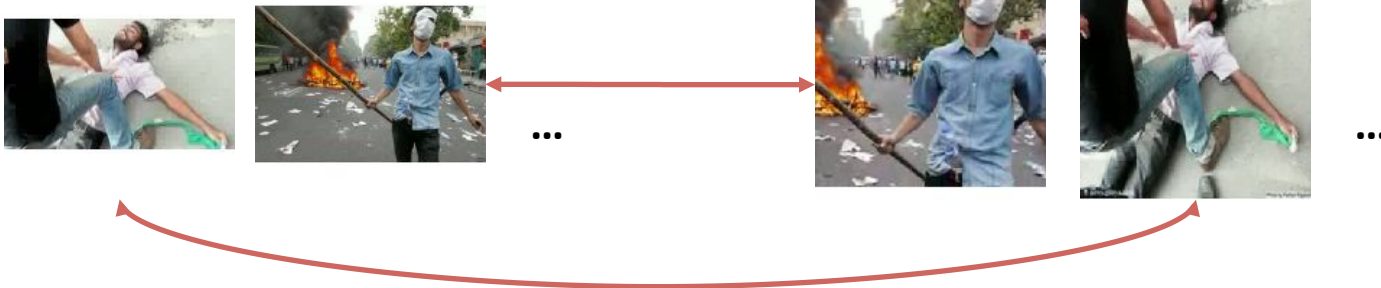# Remixing on YouTube – "Iran" topic
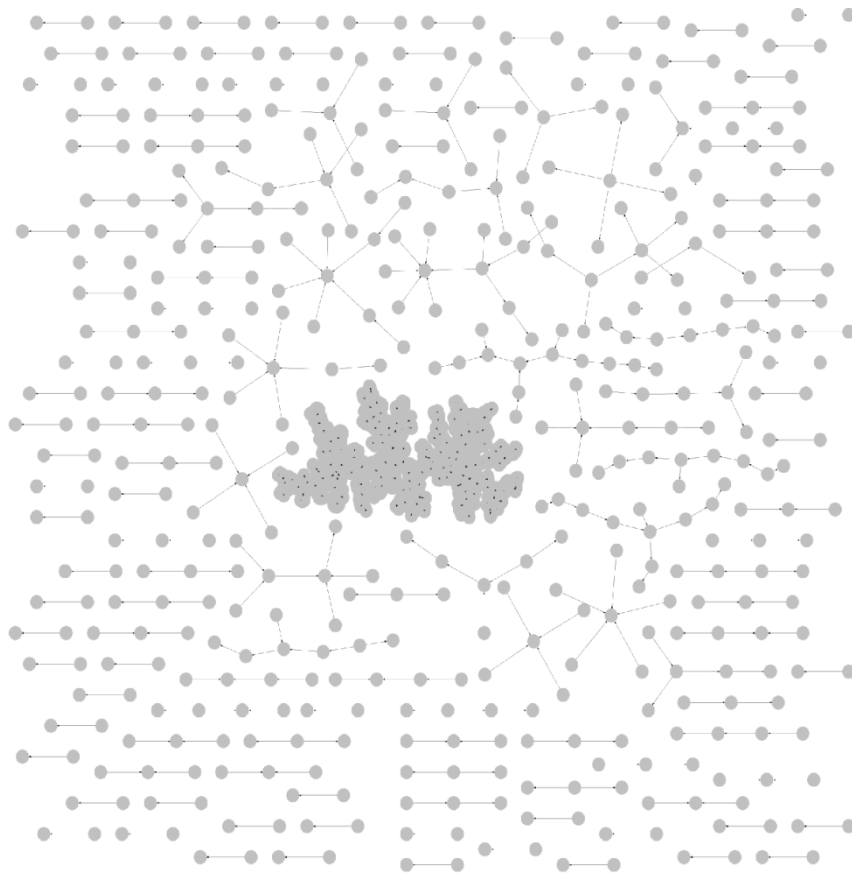


Video A                    Video B
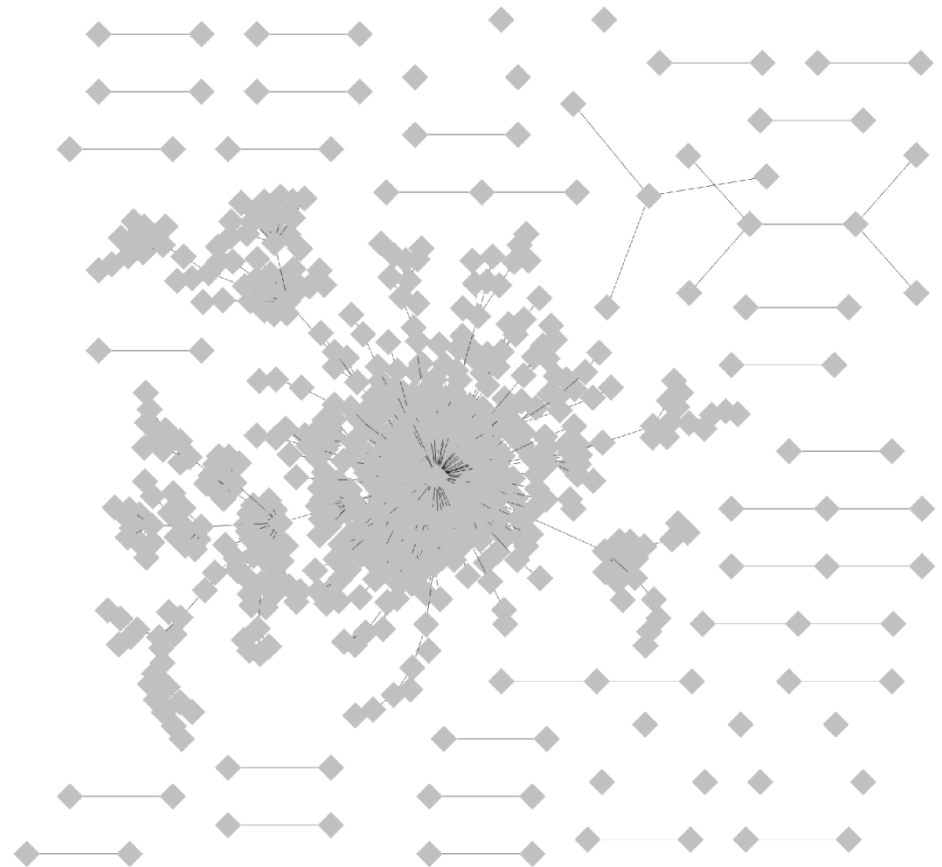
Youtube
Video page

Meme
shot
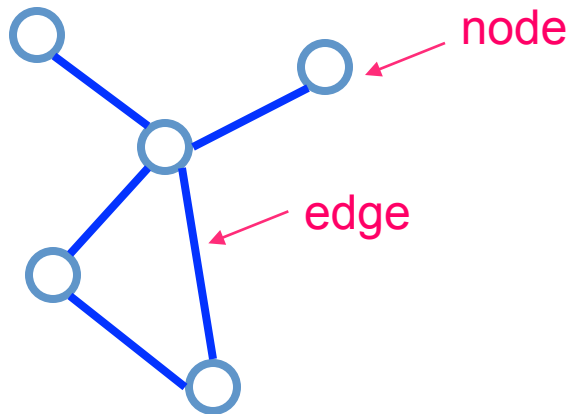examples

# YouTube Remix Network

Video graph

Author graph



Visual Memes in Social Media: Tracking Real-world News in YouTube Videos (2011), Lexing Xie, Apostol Natsev, Matthew Hill, John Kender, John R Smith, ACM Multimedia 2011, Scottsdale, AZ, USA, Nov 2011

# What are networks?

- Networks are sets of nodes connected by edges.

"Network" ≡ "Graph"



| points | lines | |
|--------|-------|--|
| vertices | edges, arcs | math |
| nodes | links | computer science |
| sites | bonds | physics |
| actors | ties, relations | sociology |

Slide by Lada Adamic, U Michigan

The attached image shows 5 streets (A and B streets, and 1st, 2nd, and 3rd Avenue). How can a network be constructed from these streets?
(Check all that apply)

☐ Roads (A St., B St., 1st Ave, ...) are nodes and an edge is drawn between every pair of roads that intersect.

☐ Intersections are nodes (e.g. A St. and 1st Ave, B St. and 2nd Ave), and an edge is drawn between any two intersections that are directly connected by a segment of street with no intervening intersections.

☐ Street blocks are nodes (e.g. the block between A and B, and 2nd and 3rd), and blocks that are adjacent (i.e. across the street from each other) have edges.

1st Ave.

2nd Ave.

3rd Ave.

A St.

B St.

Submit

Skip

# Network elements: edges

- Directed (also called arcs, links)
  - A -> B
    - A likes B, A gave a gift to B, A is B's child
- Undirected
  - A <-> B or A − B
    - A and B like each other
    - A and B are siblings
    - A and B are co-authors

# Edge attributes

- Examples
  - weight (e.g. frequency of communication)
  - ranking (best friend, second best friend…)
  - type (friend, relative, co-worker)
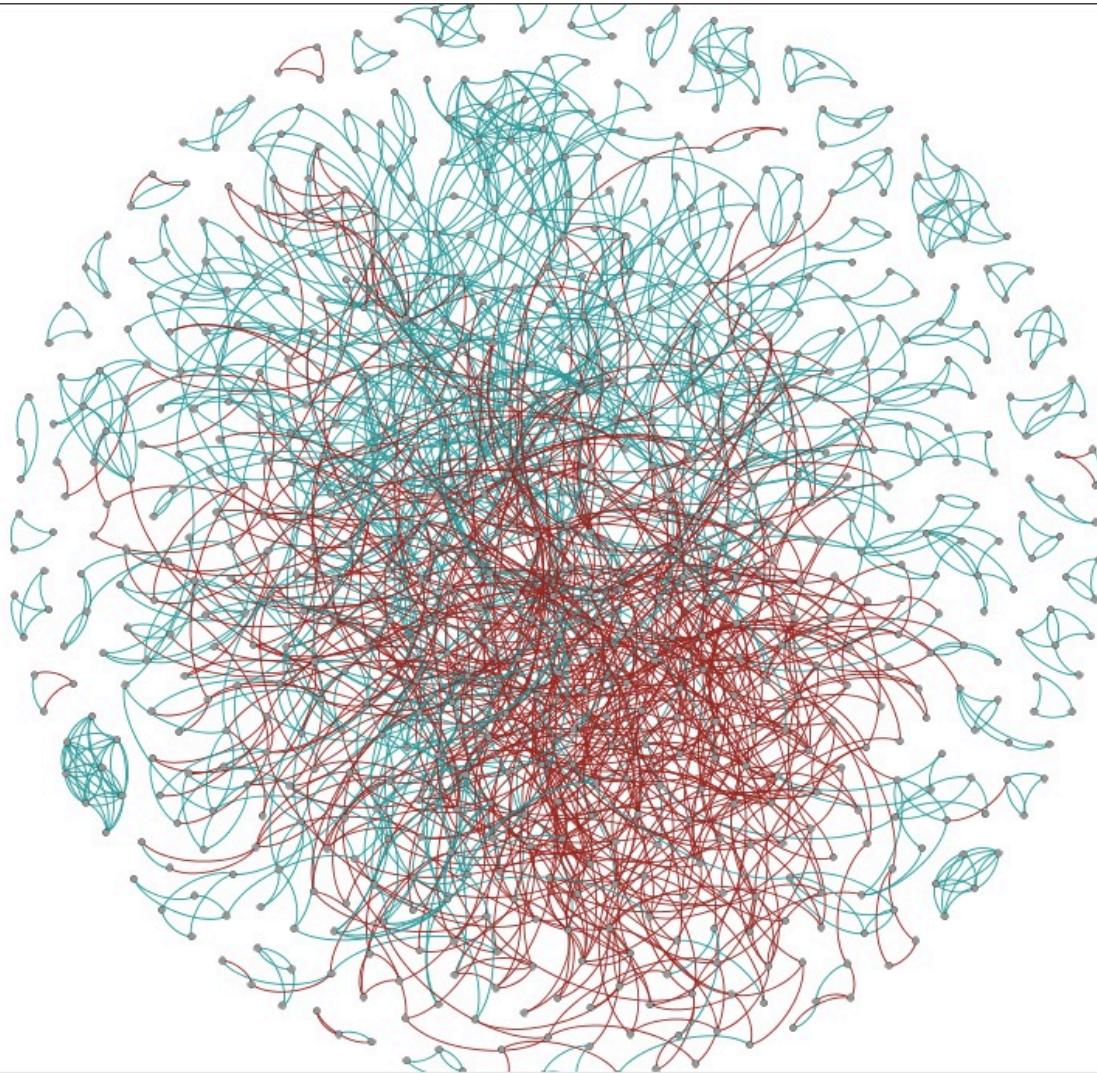  - properties depending on the structure of the rest of the graph: e.g. betweenness

by Lada Adamic, U Michigan

# Directed networks

- girls' school dormitory dining-table partners, 1st and 2nd choices
  (Moreno, *The sociometry reader*, 1960)

# Positive and negative weights



- e.g. one person trusting/distrusting another
  - Research challenge: How does one 'propagate' negative feelings in a social network? Is my enemy's enemy my friend?
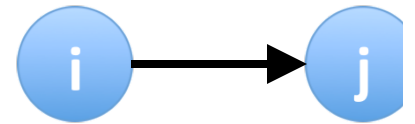
*sample of positive & negative ratings from Epinions network*
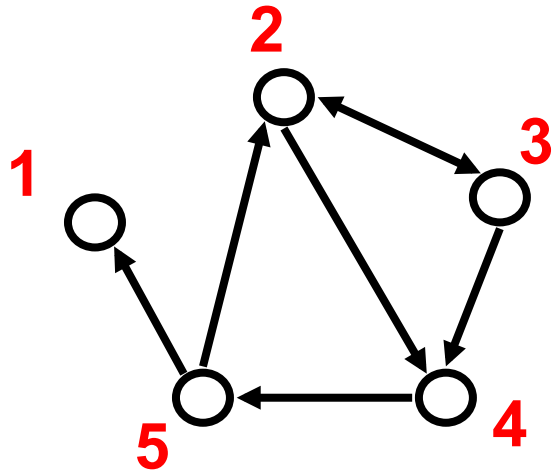by Lada Adamic, U Michigan

# Data representation

- adjacency matrix
- edgelist
- adjacency list

by Lada Adamic, U Michigan

# Adjacency matrices

- Representing edges (who is adjacent to whom) as a matrix
  - $A_{ij} = 1$ if node $i$ has an edge to node $j$
    $= 0$ if node $i$ does not have an edge to $j$

  

  - $A_{ii} = 0$ unless the network has self-loops

  

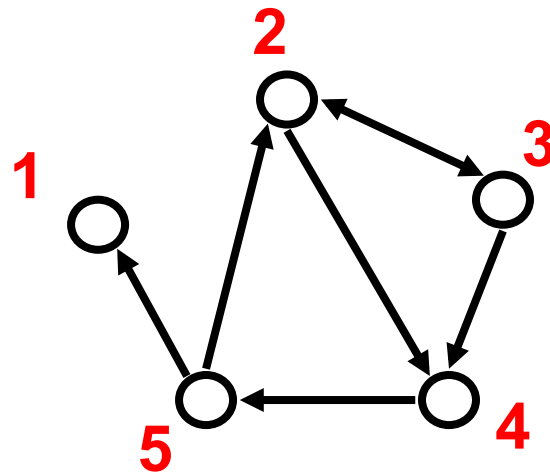  - $A_{ij} = A_{ji}$ if the network is undirected, or if $i$ and $j$ share a reciprocated edge

# Example adjacency matrix



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

# Edge list

- Edge list
  - 2, 3
  - 2, 4
  - 3, 2
  - 3, 4
  - 4, 5
  - 5, 2
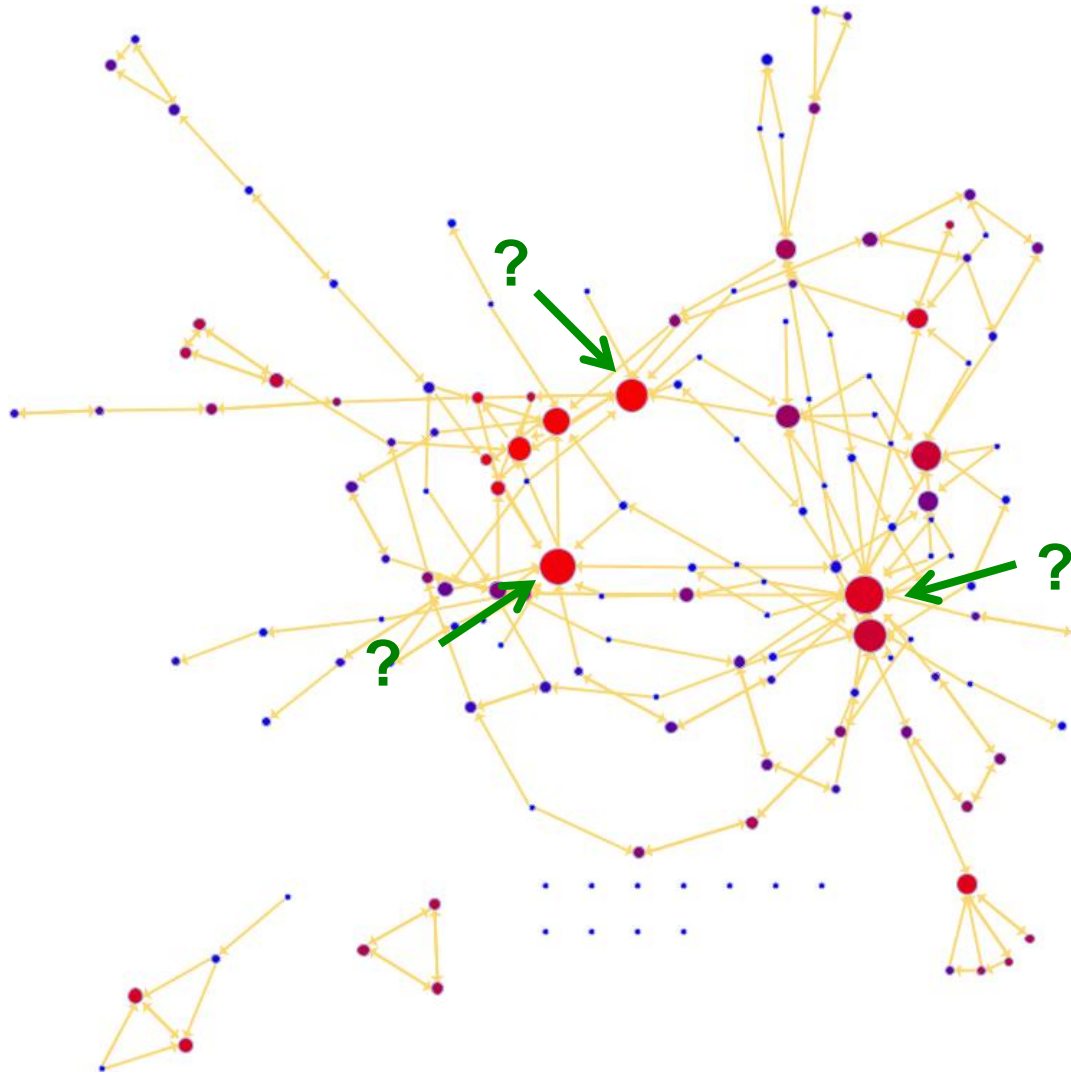  - 5, 1

# Adjacency lists

- Adjacency list
  - is easier to work with if network is
    - large
    - sparse
  - quickly retrieve all neighbors for a node
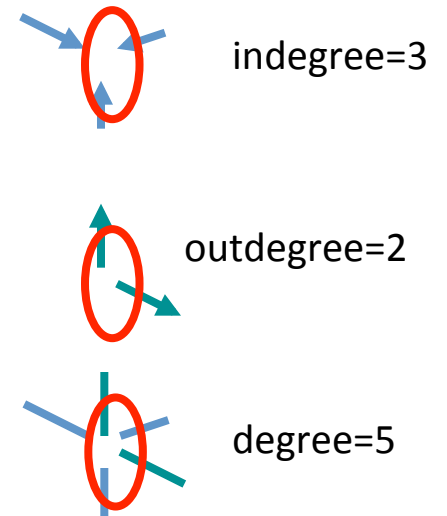    - 1:
    - 2: 3 4
    - 3: 2 4
    - 4: 5
    - 5: 1 2
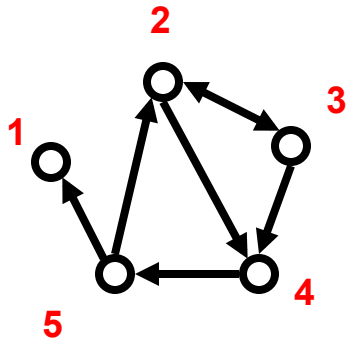
# Degree: which node has the most edges?

# Node degrees

- ## Node network properties
  - ### from immediate connections
    - #### indegree
      how many directed edges (arcs) are incident on a node

      indegree=3

    - #### outdegree
      how many directed edges (arcs) originate at a node

      outdegree=2

    - #### degree (in or out)
      number of edges incident on a node

      degree=5

  - ### from the entire graph
    - #### centrality (betweenness, closeness)

# Node degree from matrix values



- Outdegree = $\displaystyle\sum_{j=1}^{n} A_{ij}$

example: outdegree for node 3 is 2, which we obtain by summing the number of non-zero entries in the 3rd *row*

$$\sum_{j=1}^{n} A_{3j}$$

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

- Indegree = $\displaystyle\sum_{i=1}^{n} A_{ij}$

example: the indegree for node 3 is 1, which we obtain by summing the number of non-zero entries in the 3rd *column*

$$\sum_{i=1}^{n} A_{i3}$$

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

by Lada Adamic, U Michigan

# Network metrics: degree sequence and degree distribution

- Degree sequence: An ordered list of the (in,out) degree of each node
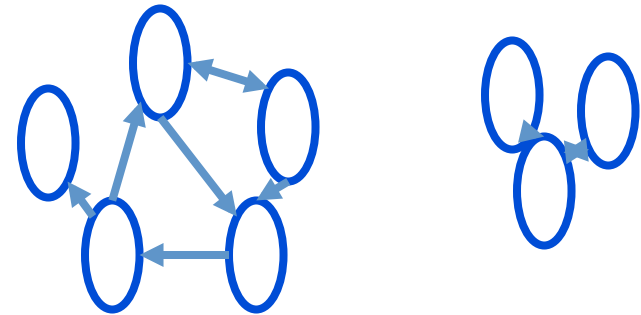
  - In-degree sequence:
    - [2, 2, 2, 1, 1, 1, 1, 0]
  - Out-degree sequence:
    - [2, 2, 2, 2, 1, 1, 1, 0]
  - (undirected) degree sequence:
    - [3, 3, 3, 2, 2, 1, 1, 1]



- Degree distribution: A frequency count of the occurrence of each degree
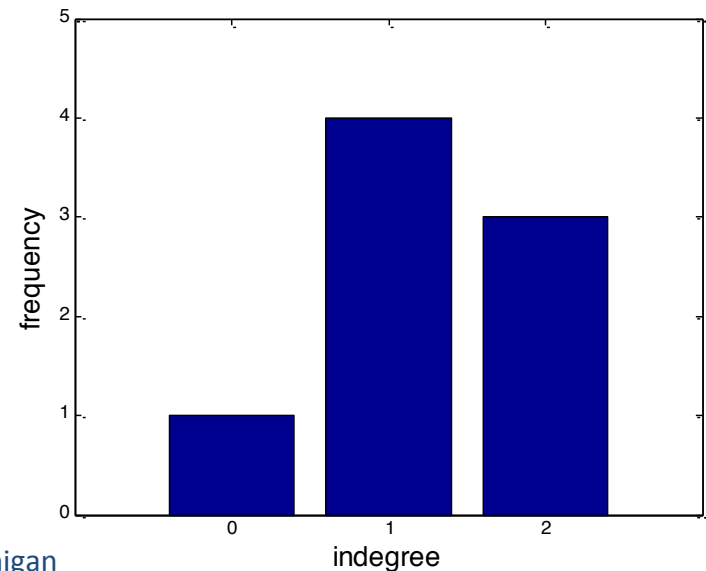  - In-degree distribution:
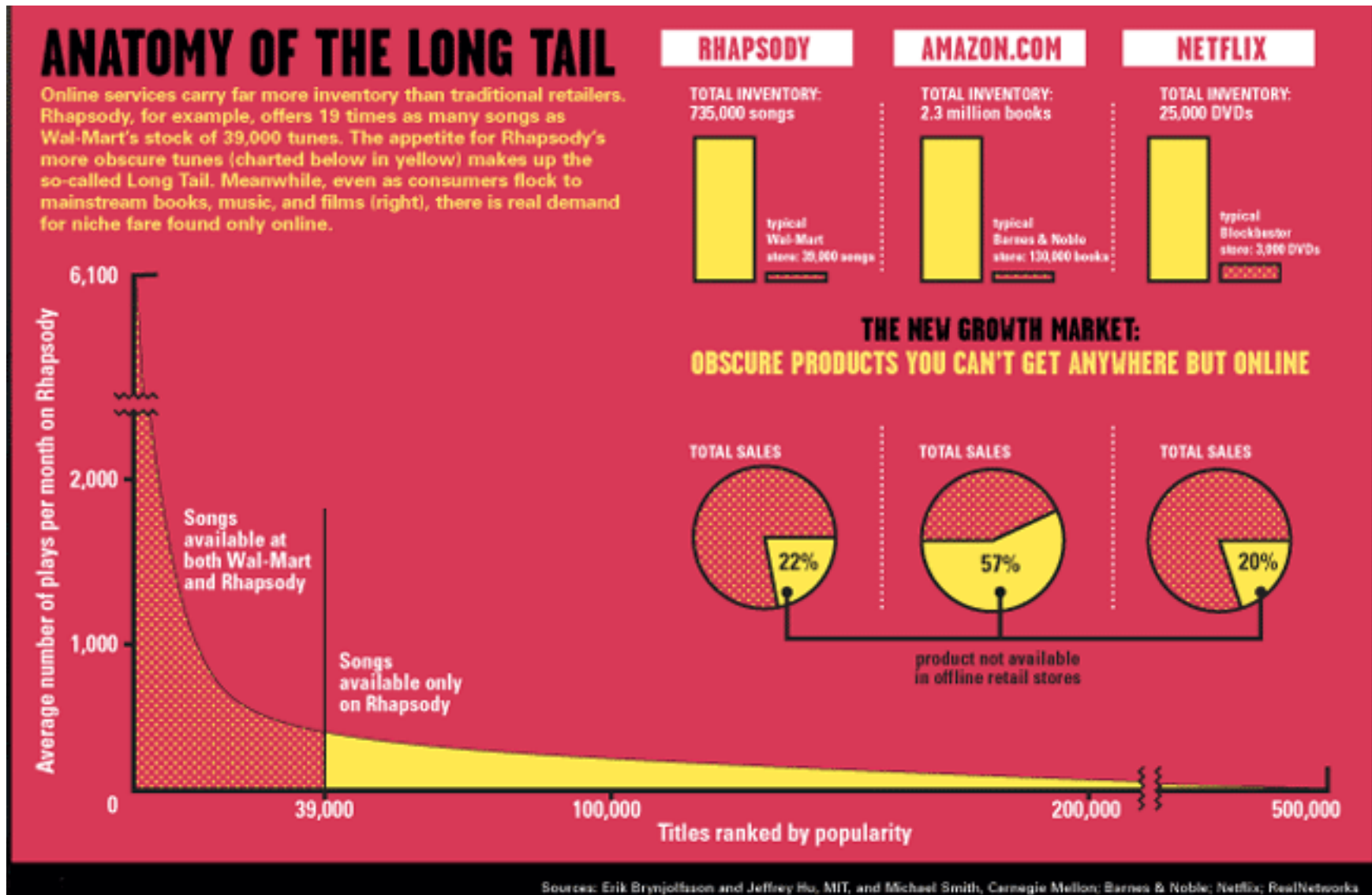    - [(2,3)  (1,4)  (0,1)]
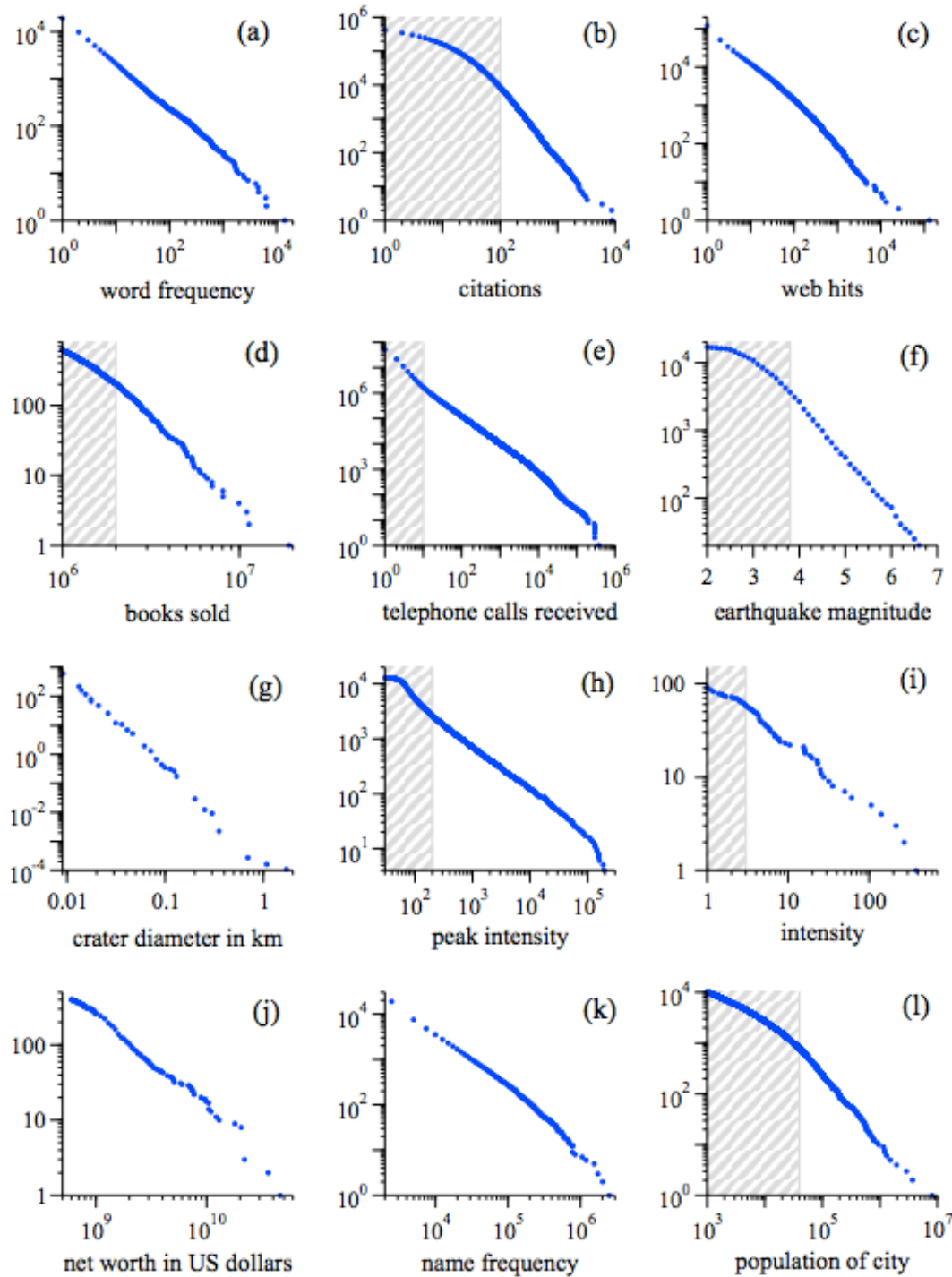  - Out-degree distribution:
    - [(2,4)  (1,3)  (0,1)]
  - (undirected) distribution:
    - [(3,3) (2,2) (1,3)]

# Long-Tail and power law

(a) word frequency
(b) citations
(c) web hits
(d) books sold
(e) telephone calls received
(f) earthquake magnitude
(g) crater diameter in km
(h) peak intensity
(i) intensity
(j) net worth in US dollars
(k) name frequency
(l) population of city

Power laws, Pareto distributions and Zipf's law. M. Newman, Cont. Physics (2005)

# Power Law Seems Ubiquitous

* Popularity of items: Amazon, YouTube, Flickr
* Degree of nodes in large graphs
    - links to website, connection between Internet Ases
    - Collaboration between actors, scientists
* Across a broad spectrum of natural system
    - # Species in genus, proteine seq. genome, words
* Time and space: characterizes human activities

# How does power laws arise?

- A process for generating web links:

1) Nodes (i.e. webpages) join the graph in sequence, creating edges linking to previous nodes.
2) Each node u creates one outgoing edge as follows
   a) Choose a node v uniformly
   b) With prob. p, edge u->v is created,
   c) With prob. (1-p), edge u->w for w a (uniformly chose) children of v

   c) is equivalent to: "with probability 1−p, page u chooses a page w with probability proportional to w's current number of in-links, and creates a link to w"

# Rich-Get-Richer Process

Initial condition

$$x_j(j) = 0$$

Rate of growth

$$\frac{dx_j}{dt} = \frac{p}{t} + \frac{(1-p)x_j}{t}.$$

Fraction of nodes with k incoming links

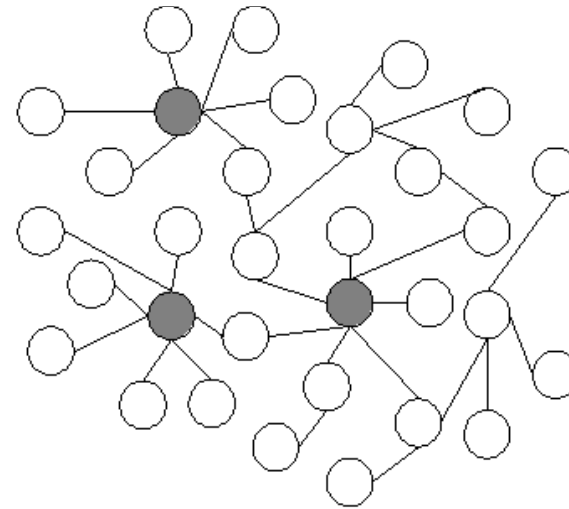~ k$^{-\alpha}$,  $\alpha$= 1 + 1/(1-p)

[Easly and Kleinberg, Chapter 18] Fig 18.2 and 18.3

# Random vs Power-law networks



(a) Random network

(b) Scale-free network

Terminologies that describe (different aspects of) the same process:
Power-law, Long-tail
Scale-free
Rich-get-richer, preferential attachment

# Networks from Document Collections

| Dataset | Nodes | Edges | Edge weight? | Directed? | Potential use of the network |
|---|---|---|---|---|---|
| "Divided they Blog" | blog | | | | |
| Email | people | | | | |
| Research Papers | authors | | | | |
| Research Papers | papers | | | | |
| Research Papers | Journals/ conferences | | | | |
| Recipe | | | | | |
| … | | | | | |

# Summary for Today

- What are networks of documents

- A collection of network analysis problems

- Network representation

- Power-law degree distribution and rich-get-richer processes


- HW: representing Twitter as network

- Coming up in Wed: network description

# Document Elements In Twitter

**User**

**Mention**

**Hashtag**

**Hyperlink**

# Networks from Twitter

| Nodes | Edges | Edge weight? | Directed? | Potential use of the network |
|---|---|---|---|---|
| User | | | | |
| User | | | | |
| Hashtag | | | | |
| Website | | | | |
| Website, User | | | | |
| User | | | | |