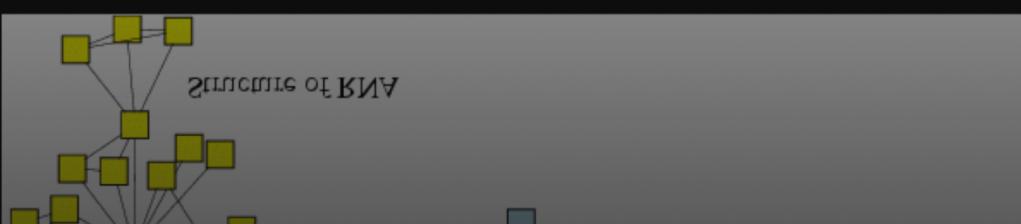
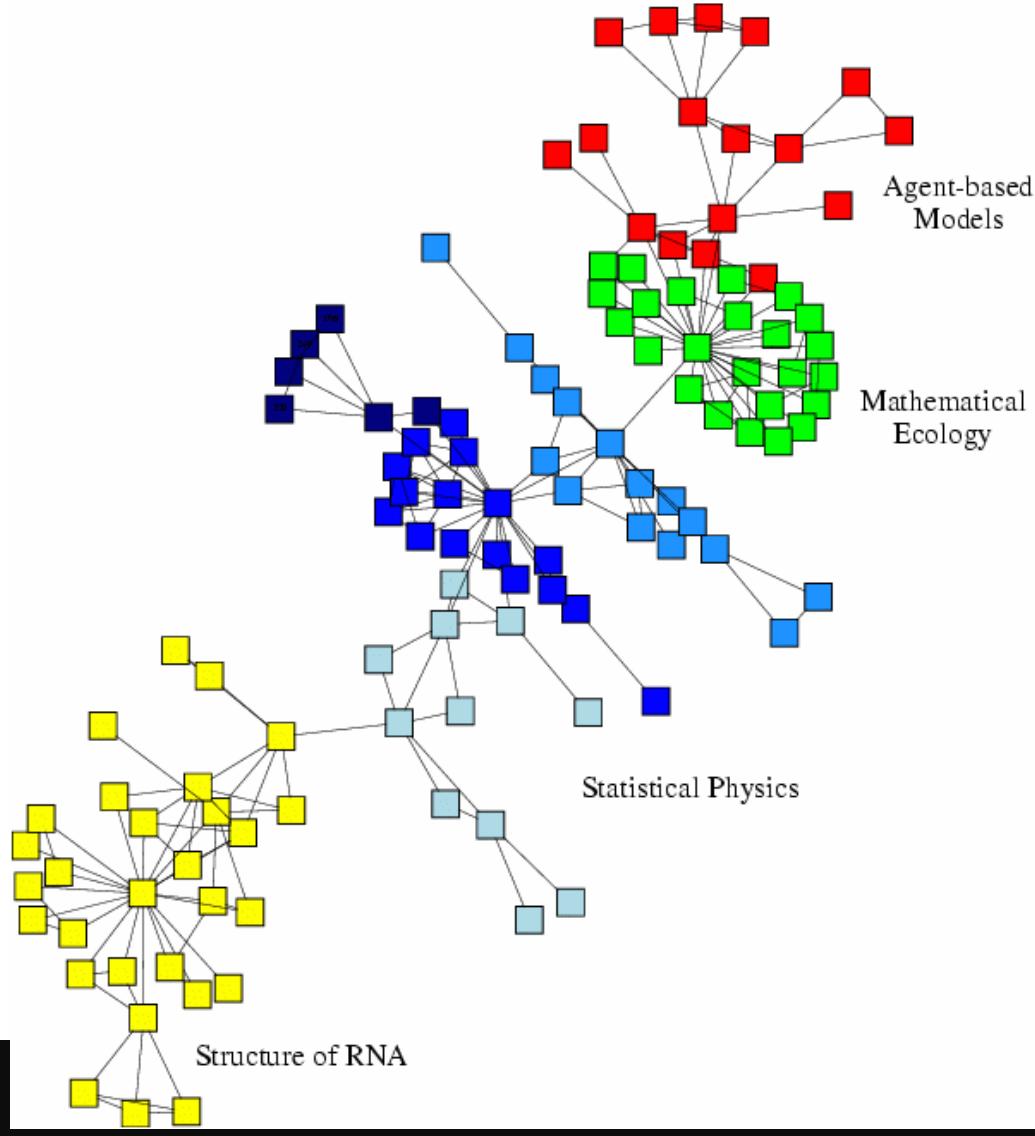
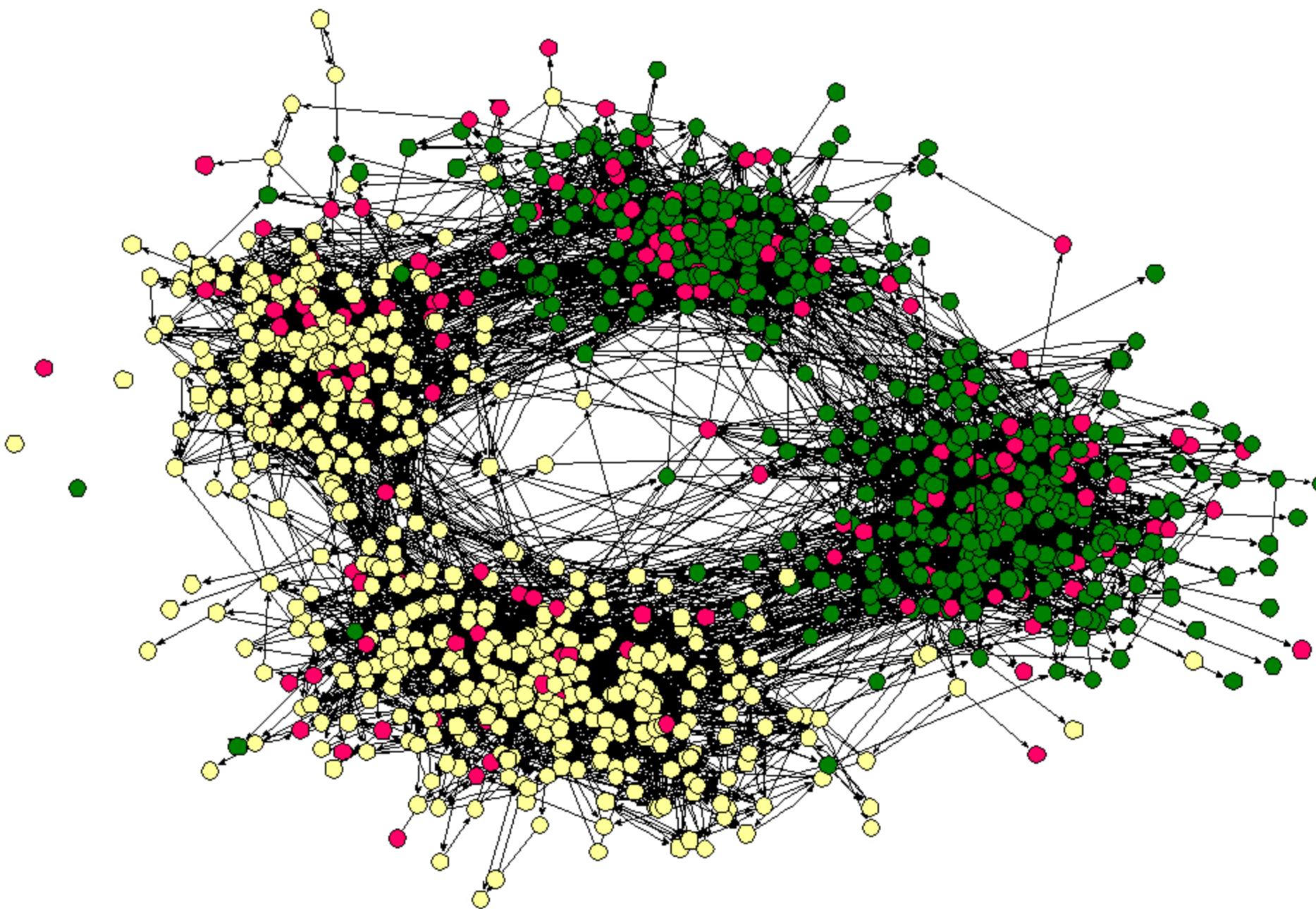


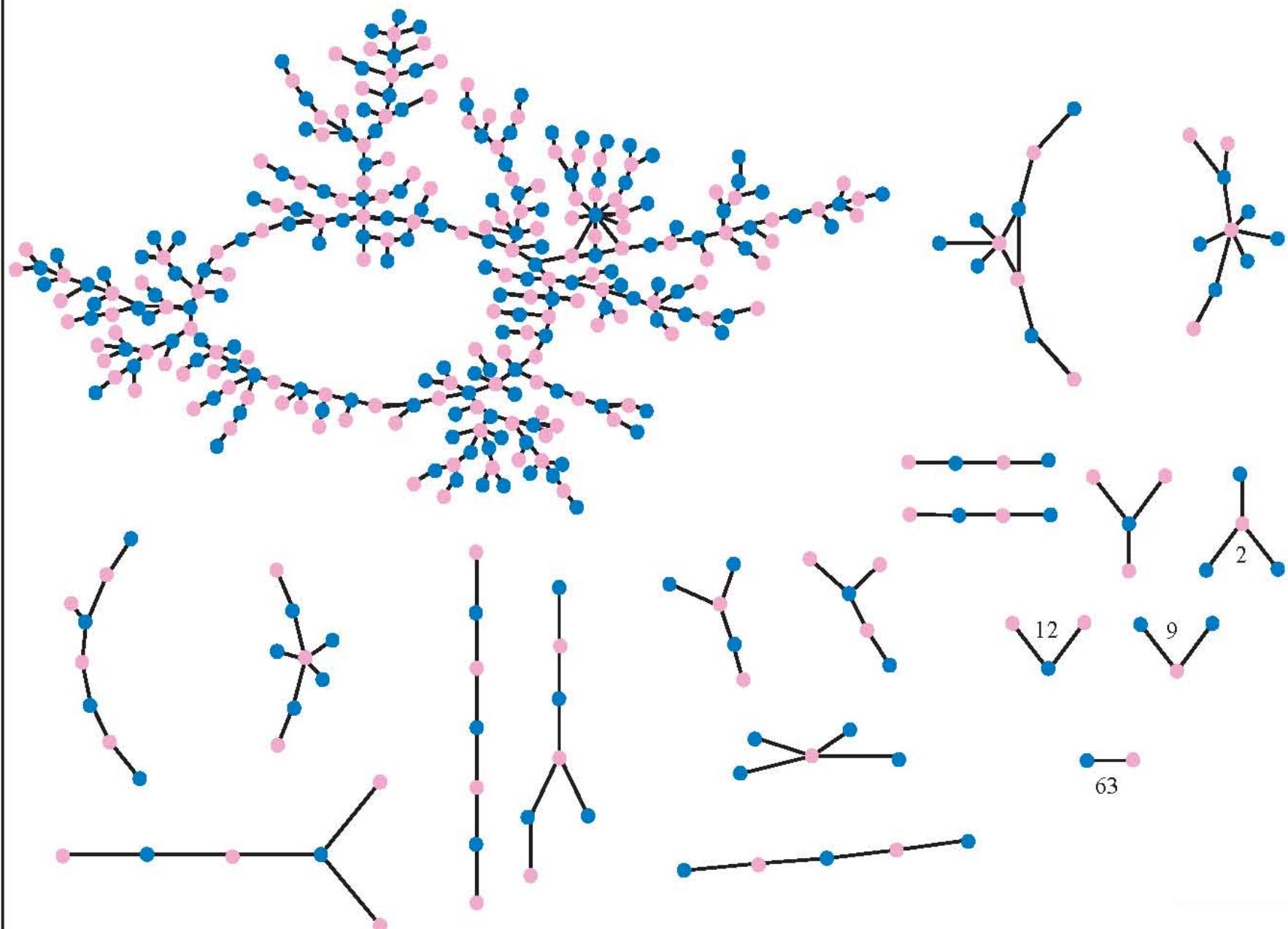
Foundations of Data Science

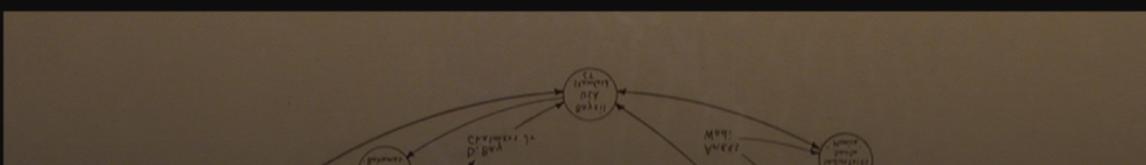
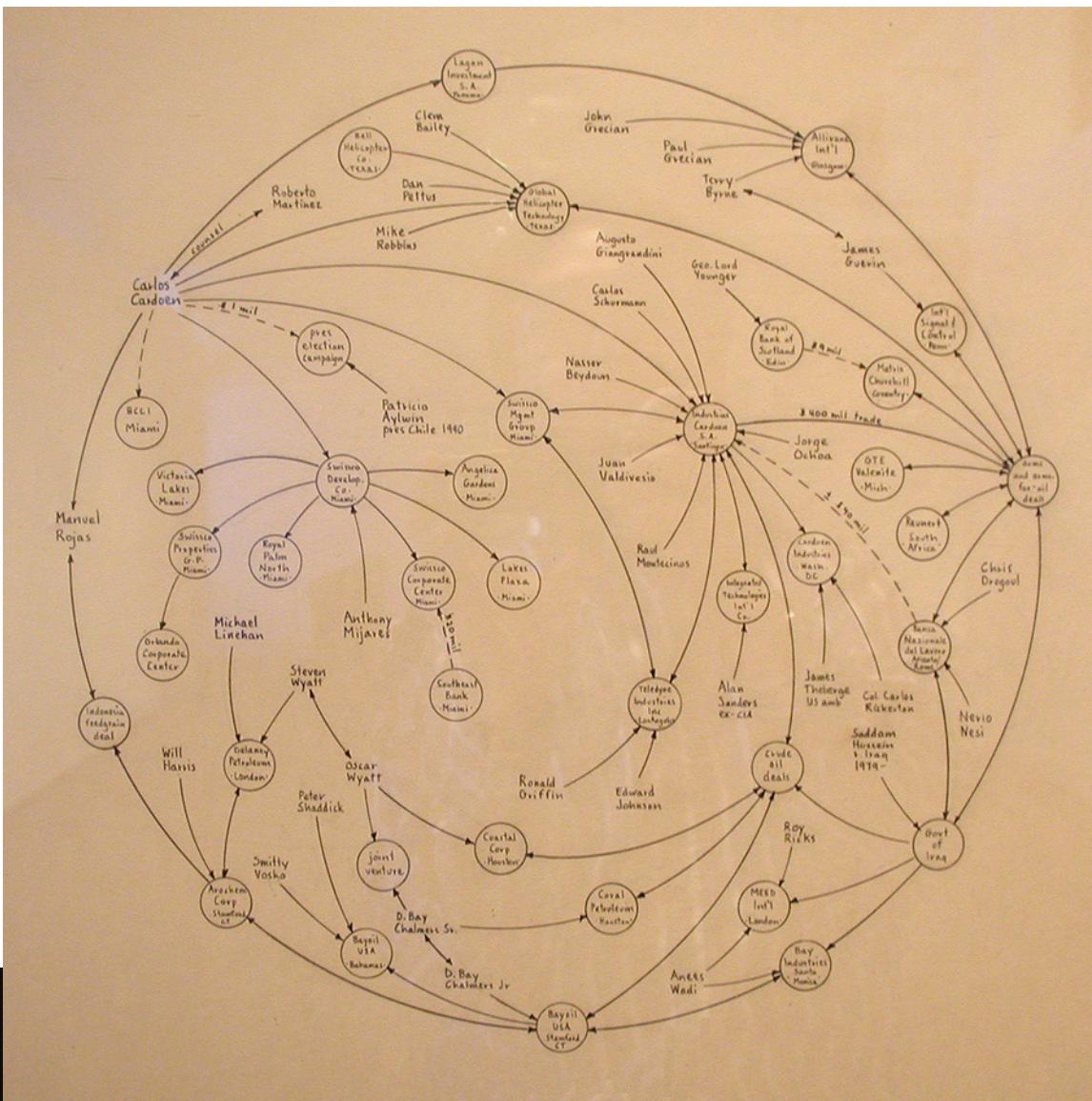
Rumi Chunara, PhD

CS3943/9223

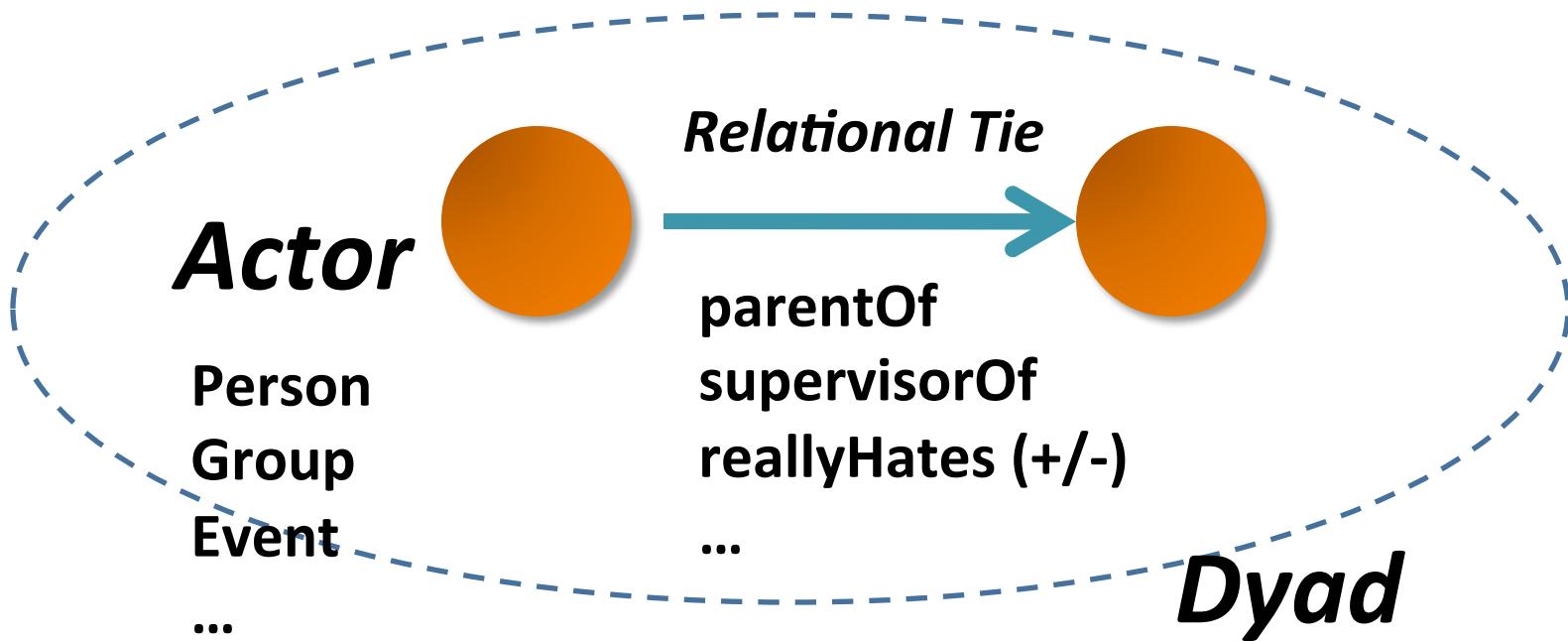








Vocabulary Lesson

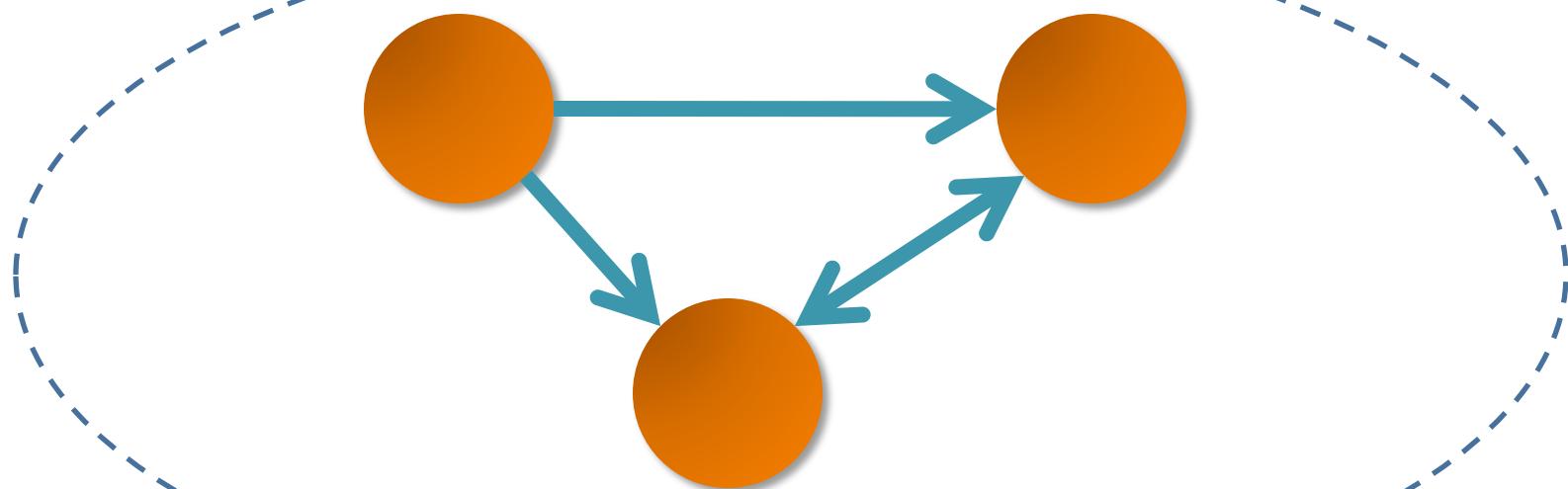


Relation: collection of ties of a specific type
(every parentOf tie)

Vocabulary Lesson

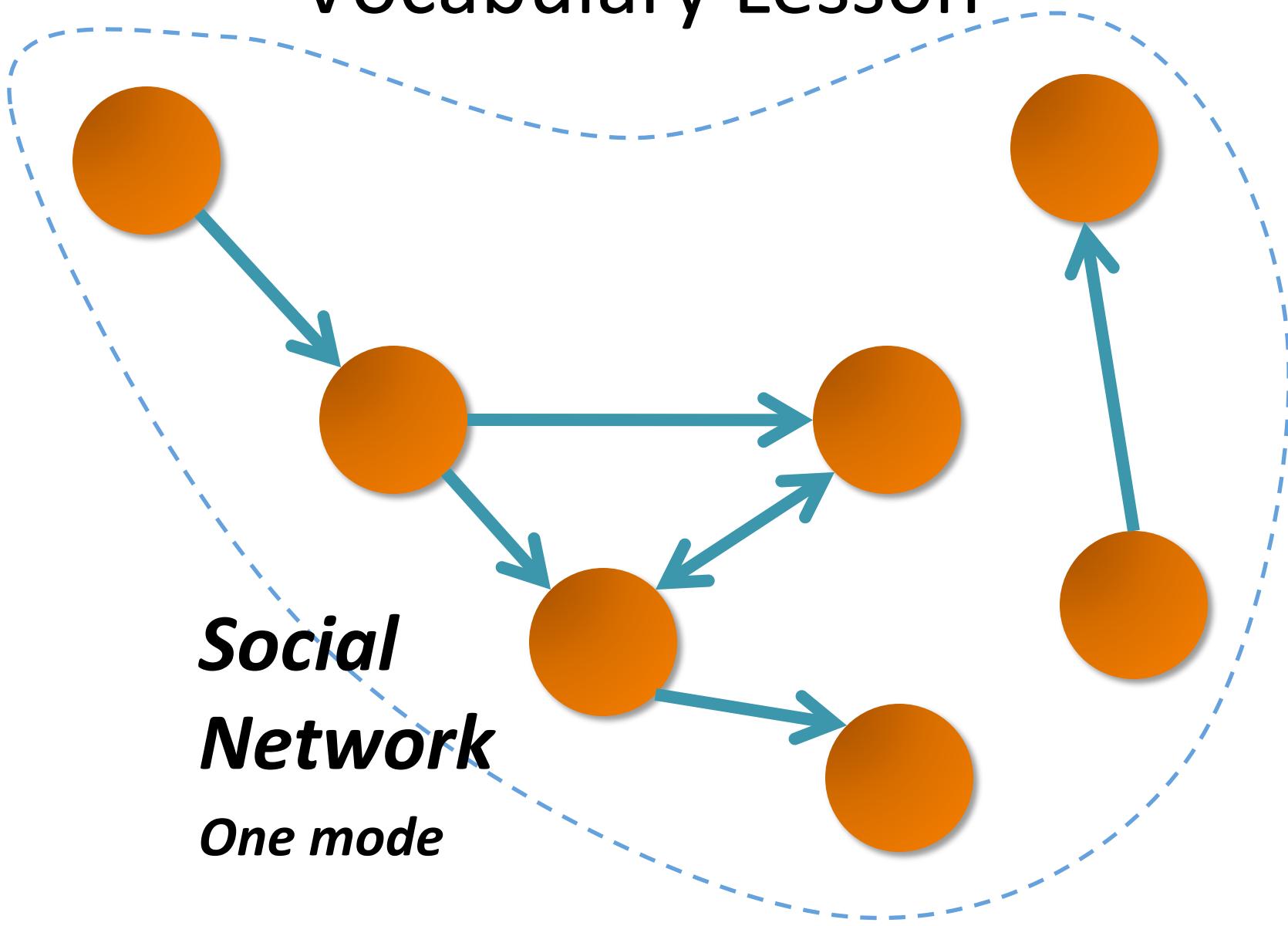
If A likes B and B likes C then A likes C (transitivity)
If A likes B and C likes B then A likes C

...

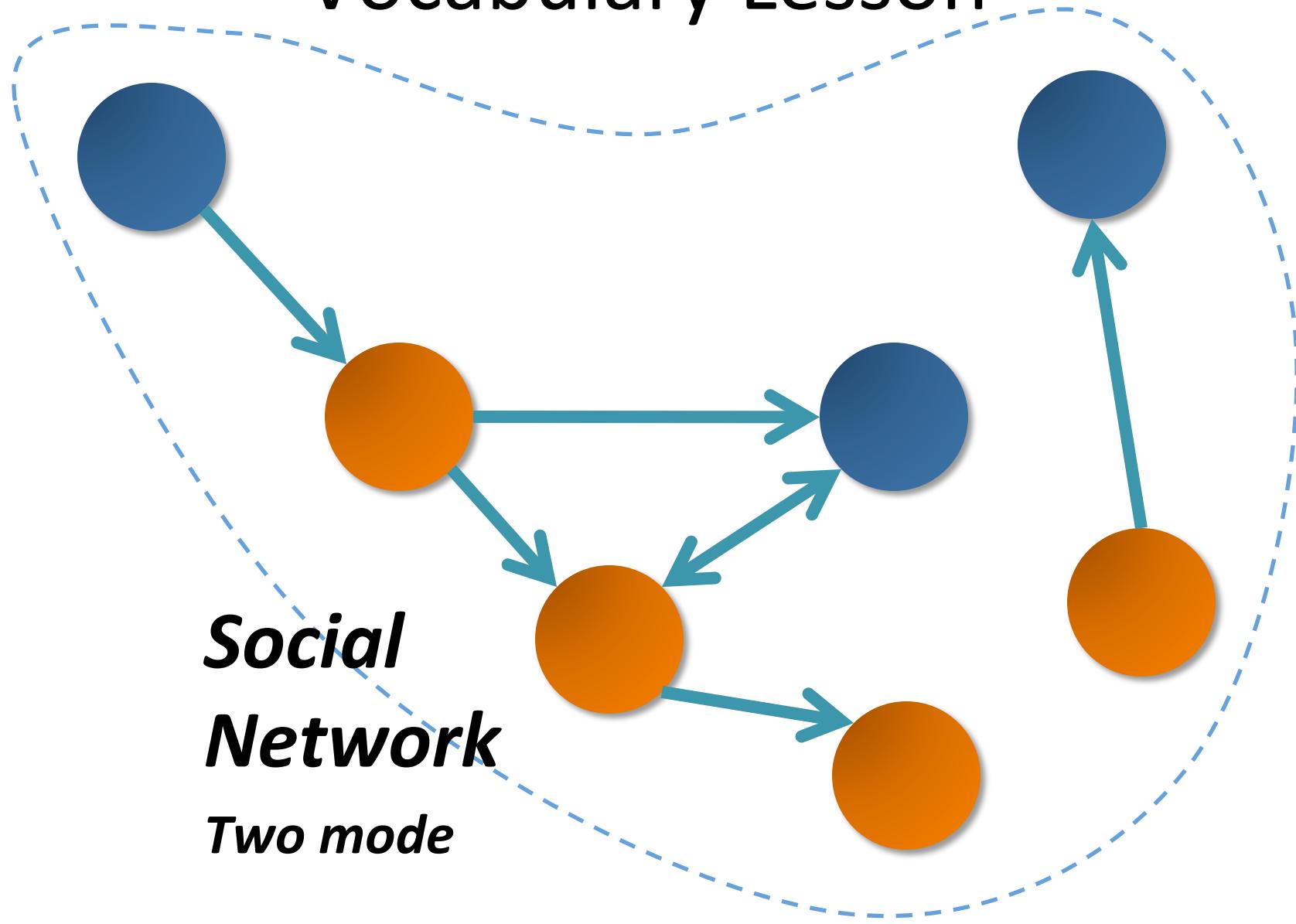


Triad

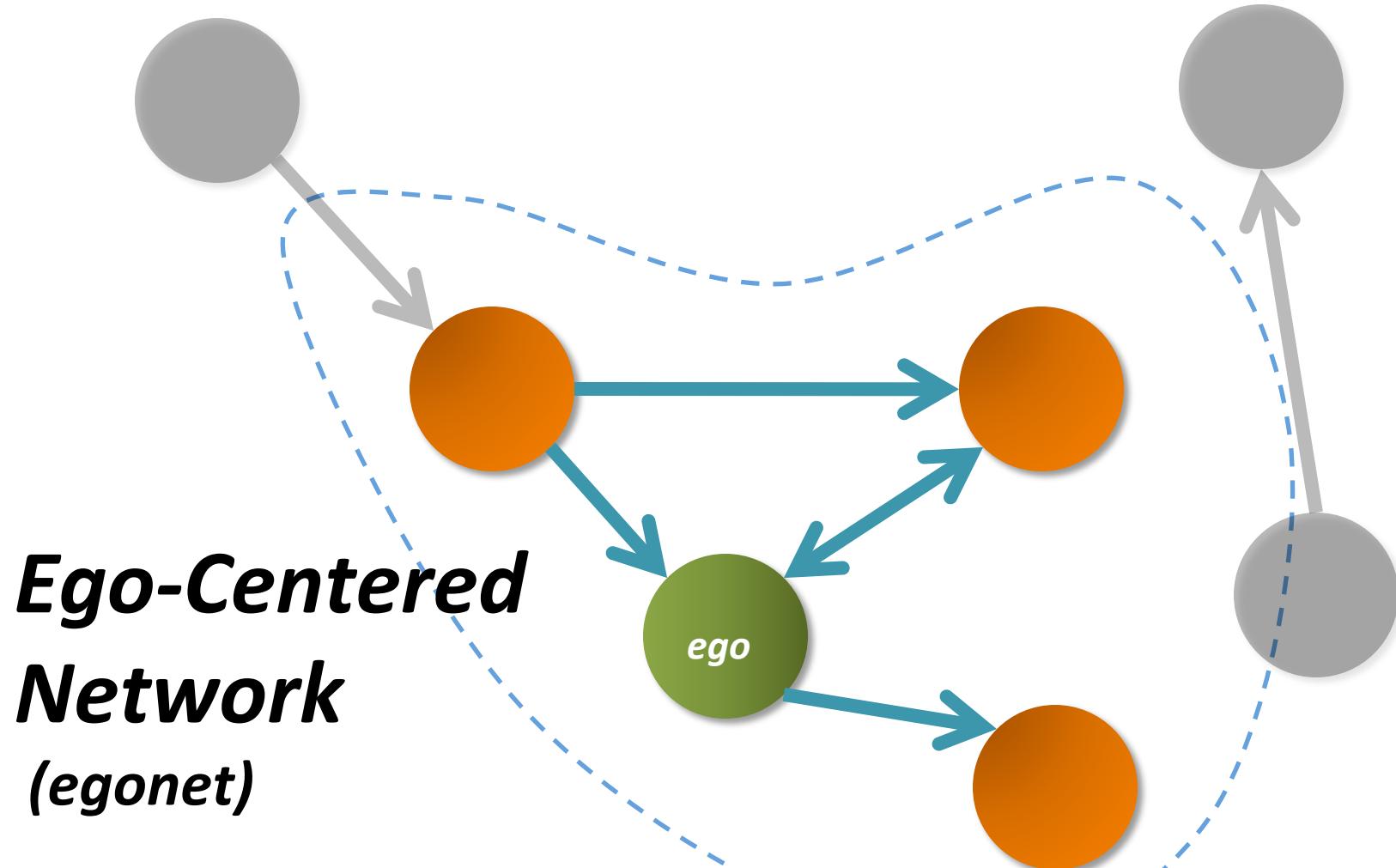
Vocabulary Lesson



Vocabulary Lesson



Vocabulary Lesson

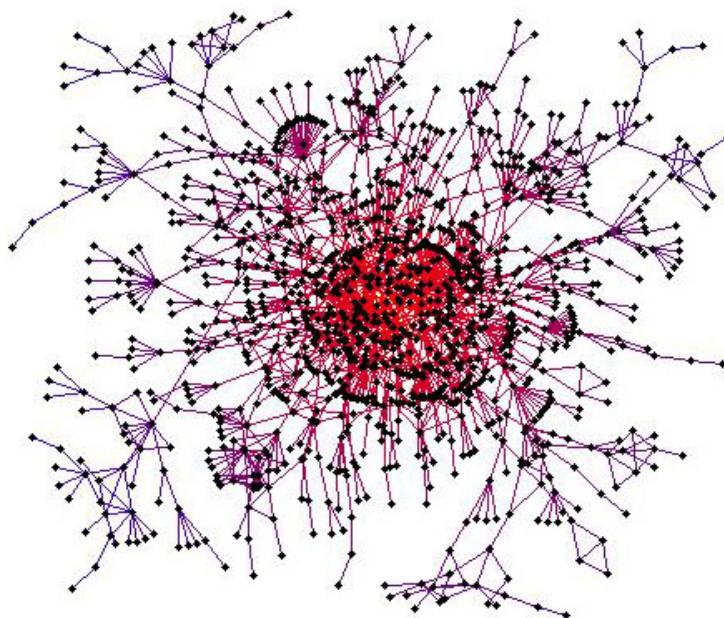


Describing Networks

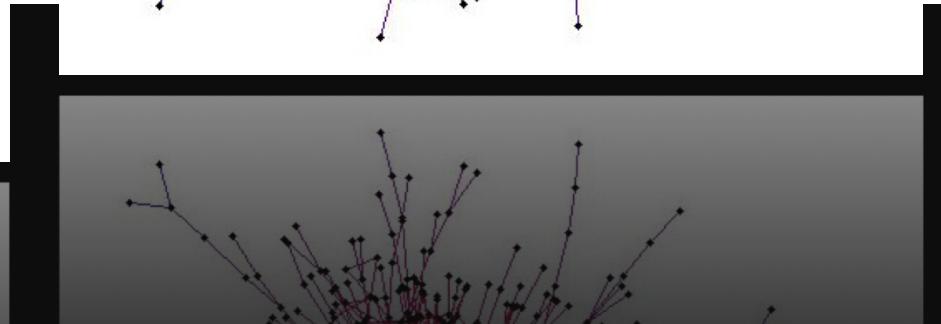
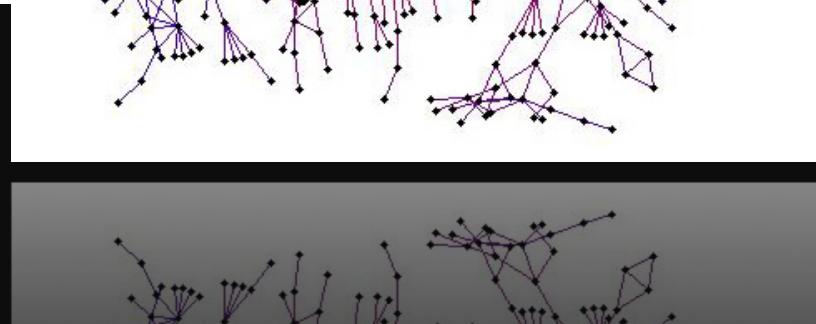
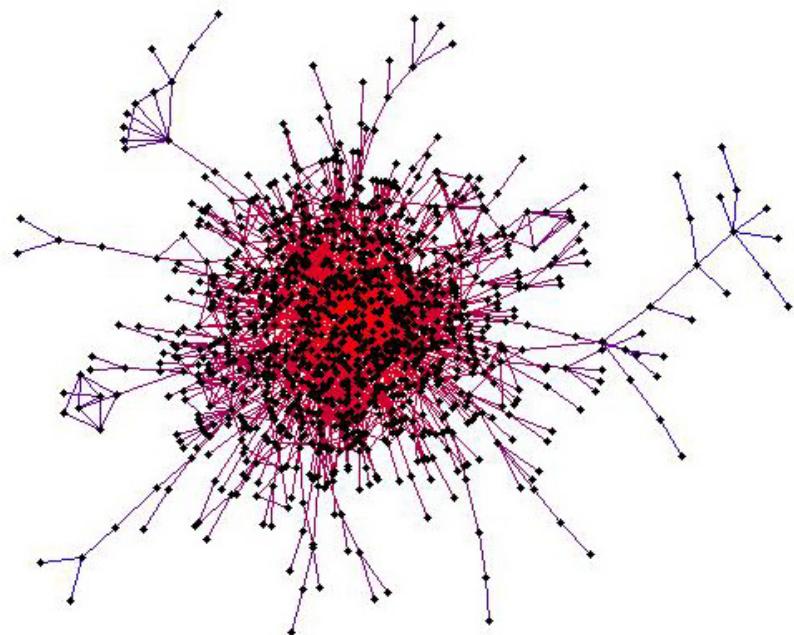
- Graph theoretic
 - Nodes/edges
- Sociometric
 - Sociomatrix (2D matrix representation)
 - Sociogram (the adjacency matrix)
- Algebraic
 - $n_i \rightarrow n_j$
- Basically complimentary

Describing Networks

Stanford



MIT



Describing Networks

- *Geodesic*
 - `shortest_path(n,m)`
- *Diameter*
 - `max(geodesic(n,m))` n,m actors in graph
- *Density*
 - Number of existing edges / All possible edges
- *Degree distribution*

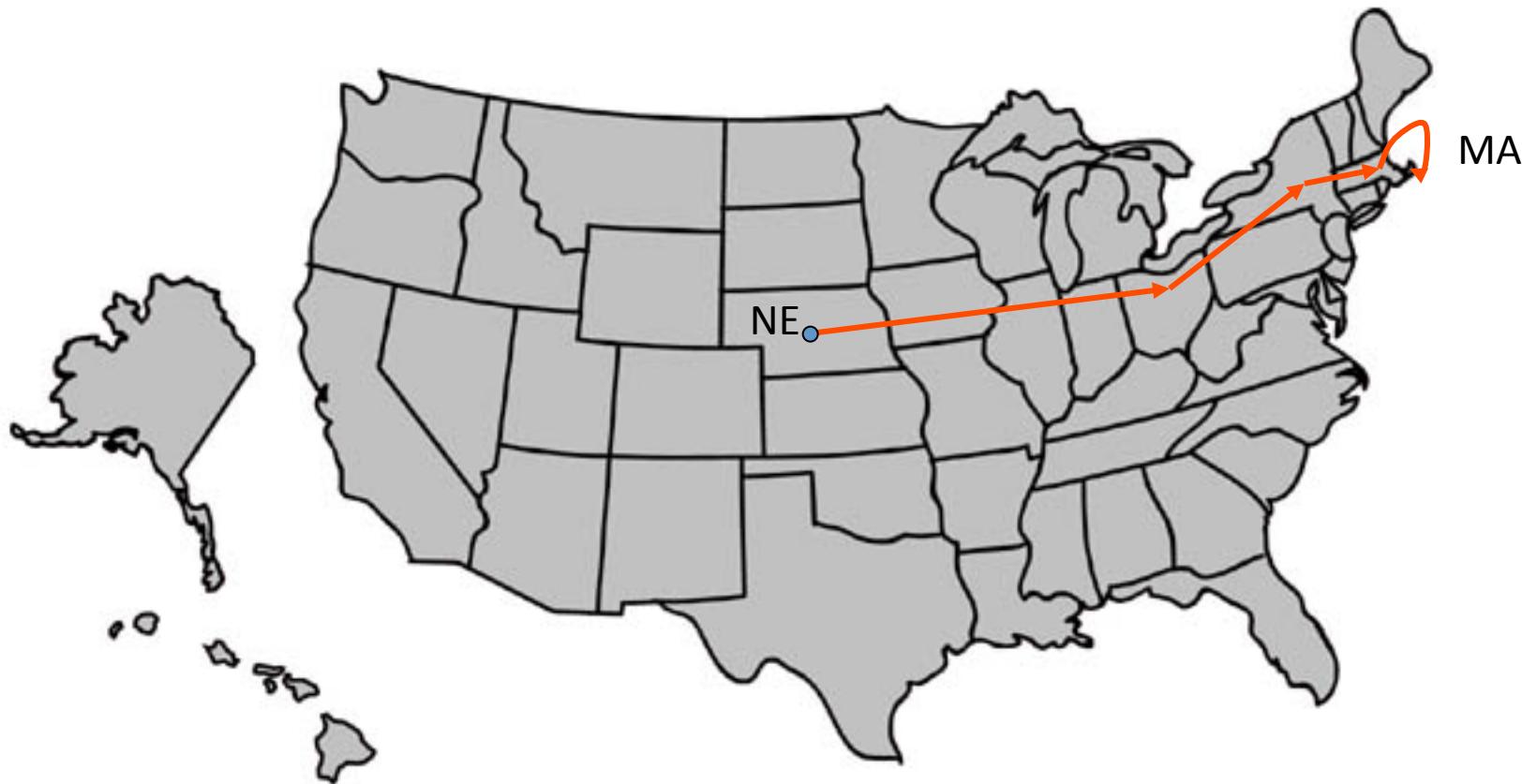
Types of Networks/Models

- A few quick examples
 - Erdős–Rényi
 - $G(n,M)$: randomly draw M edges between n nodes
 - Does not really model the real world
 - Average connectivity on nodes conserved

Types of Networks/Models

- A few quick examples
 - Erdős–Rényi
 - Small World
 - Watts-Strogatz
 - Kleinberg lattice model

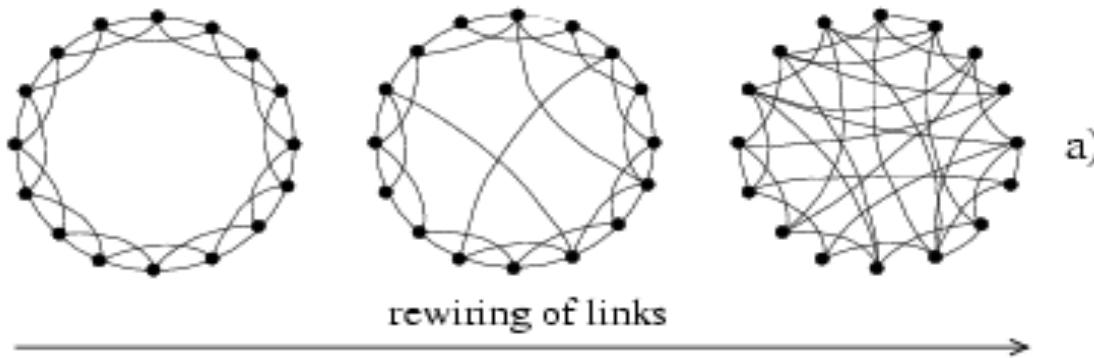
Small world experiments then



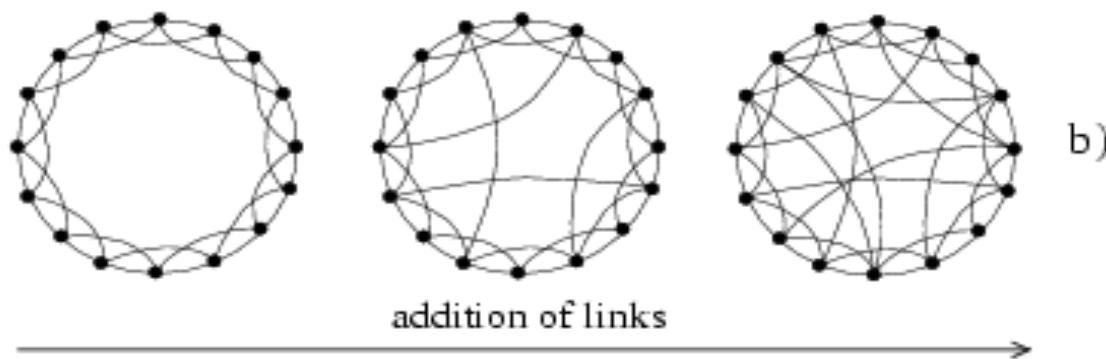
Milgram's experiment (1960's):

Given a target individual and a particular property, pass the message to a person you correspond with who is “closest” to the target.

Watts-Strogatz Ring Lattice Rewiring



a)
Select a fraction p of
edges
Reposition one of their
endpoints

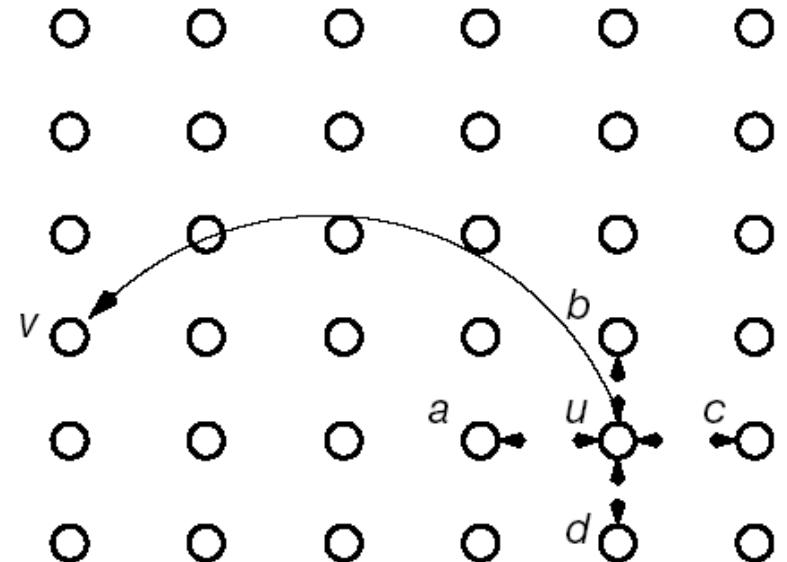


b)
Add a fraction p of
additional edges leaving
underlying lattice intact

- As in many network generating algorithms
 - Disallow self-edges
 - Disallow multiple edges

Kleinberg Lattice Model

nodes are placed on a lattice and connect to nearest neighbors

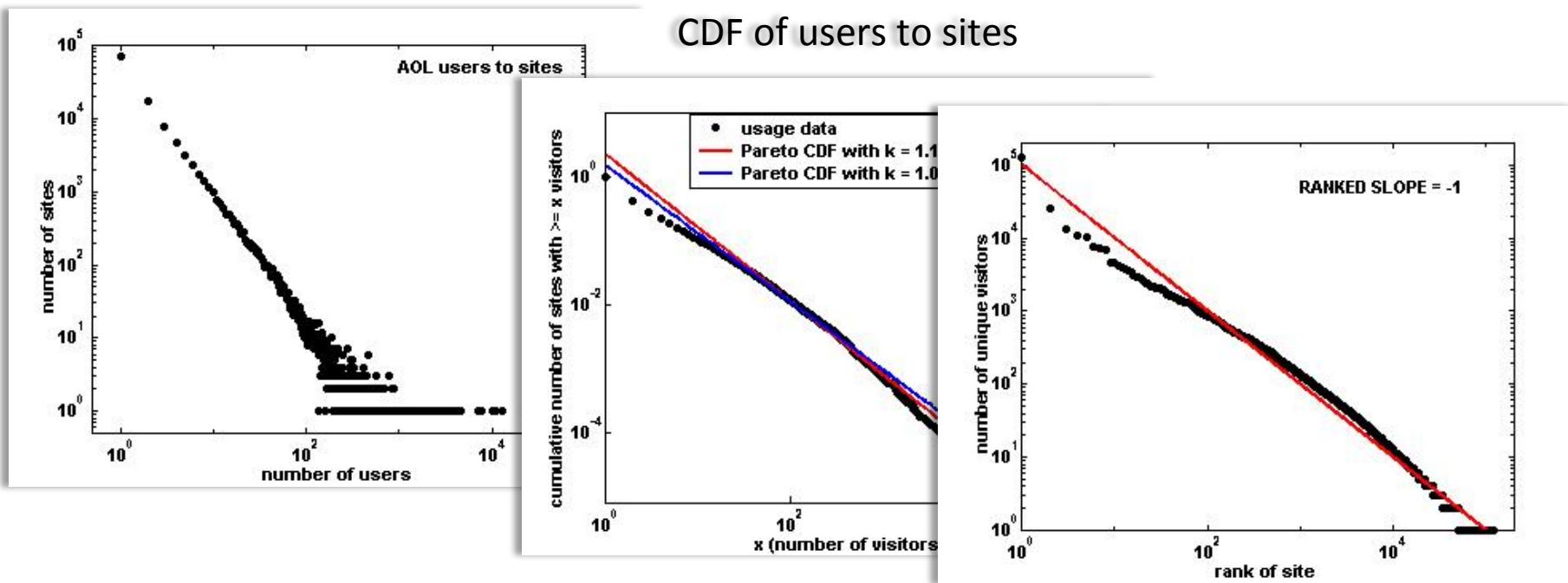


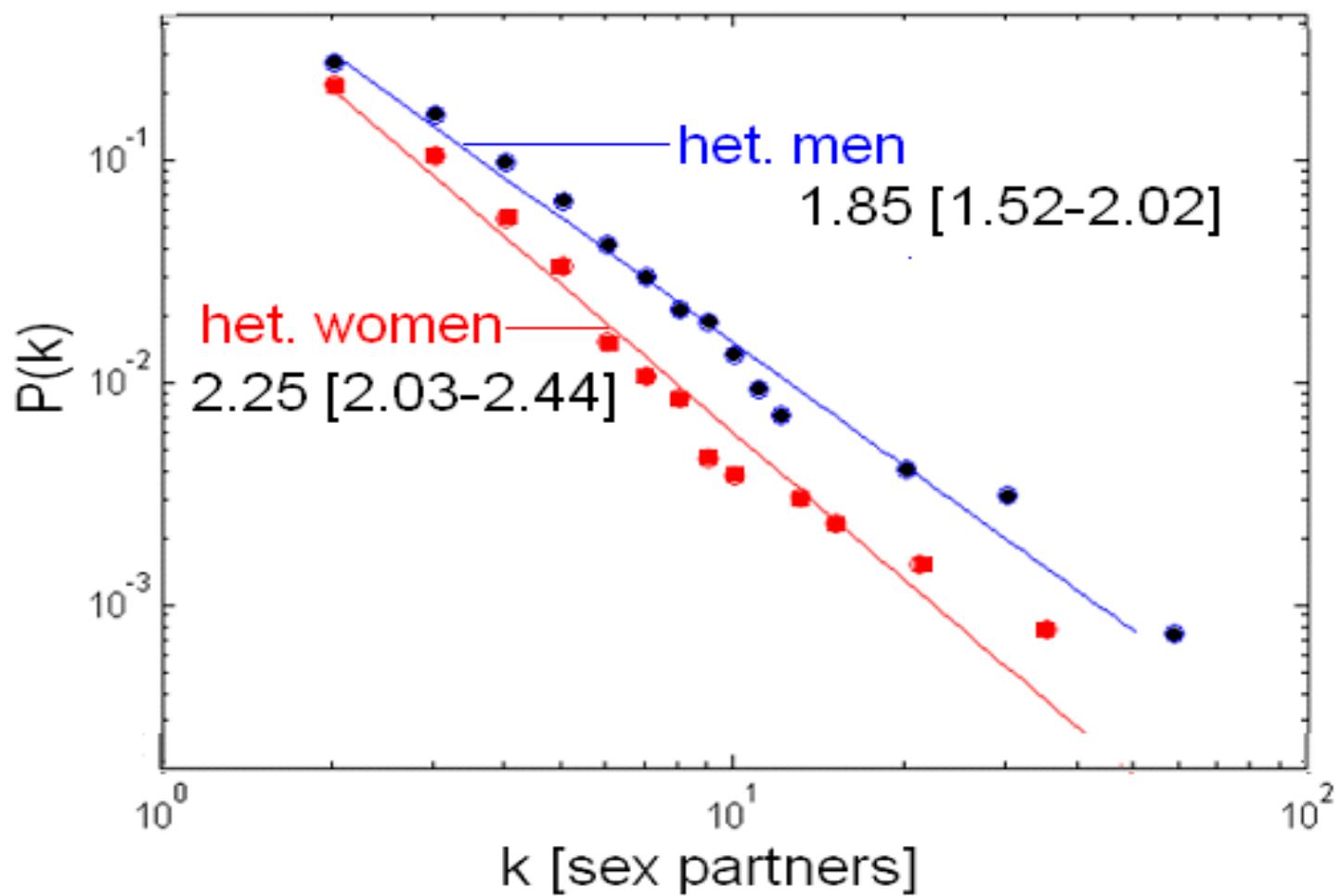
additional links placed with $p_{uv} \sim d_{uv}^{-r}$

Kleinberg, 'The Small World Phenomenon, An Algorithmic Perspective'
(Nature 2000)

A little more on degree distribution

- Power-laws, zipf, etc.





A little more on degree distribution

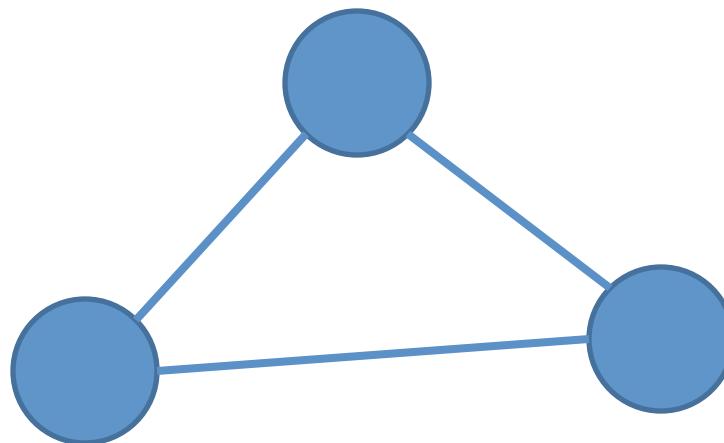
- Pareto/Power-law
 - Pareto: CDF $P[X > x] \sim x^{-k}$
 - Power-law: PDF $P[X = x] \sim x^{-(k+1)} = x^{-a}$
 - Some recent debate (Aaron Clauset)
 - <http://arxiv.org/abs/0706.1062>
- Zipf
 - Frequency versus rank $y \sim r^{-b}$ (small b)
- More info:
 - Zipf, Power-laws, and Pareto – a ranking tutorial (<http://www.hpl.hp.com/research/idl/papers/ranking/ranking.html>)

Types of Networks/Models

- A few quick examples
 - Erdős–Rényi
 - Small World
 - Watts-Strogatz
 - Kleinberg lattice model
 - Preferential Attachment
 - Generally attributed to Barabási & Albert

Basic BA-model

- Very simple algorithm to implement
 - start with an initial set of m_0 fully connected nodes
 - e.g. $m_0 = 3$



- now add new vertices one by one, each one with exactly m edges
 - each new edge connects to an existing vertex in proportion to the number of edges that vertex already has
→ ***preferential attachment***

Properties of the BA graph

- The distribution is scale free with exponent $\alpha = 3$
 $P(k) = 2 m^2/k^3$
- The graph is connected
 - Every new vertex is born with a link or several links (depending on whether $m = 1$ or $m > 1$)
 - It then connects to an ‘older’ vertex, which itself connected to another vertex when it was introduced
 - And we started from a connected core
- The older are richer
 - Nodes accumulate links as time goes on, which gives older nodes an advantage since newer nodes are going to attach preferentially – and older nodes have a higher degree to tempt them with than some new kid on the block

Common Tasks

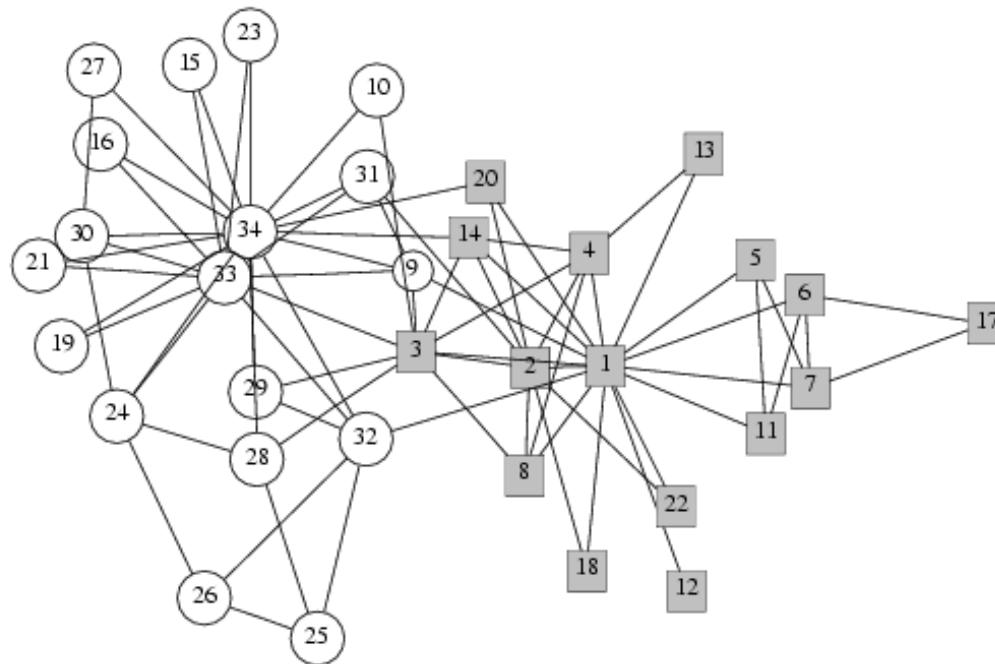
- Measuring “importance”
 - Centrality, prestige
- Link prediction
- Diffusion modeling
 - Epidemiological
- Clustering
 - Blockmodeling, Girvan-Newman
- Structure analysis
 - Motifs, Isomorphisms, etc.
- Visualization/Privacy/etc.

Data Collection / Cleaning



Analysis

(a)



Past

Data Collection / Cleaning

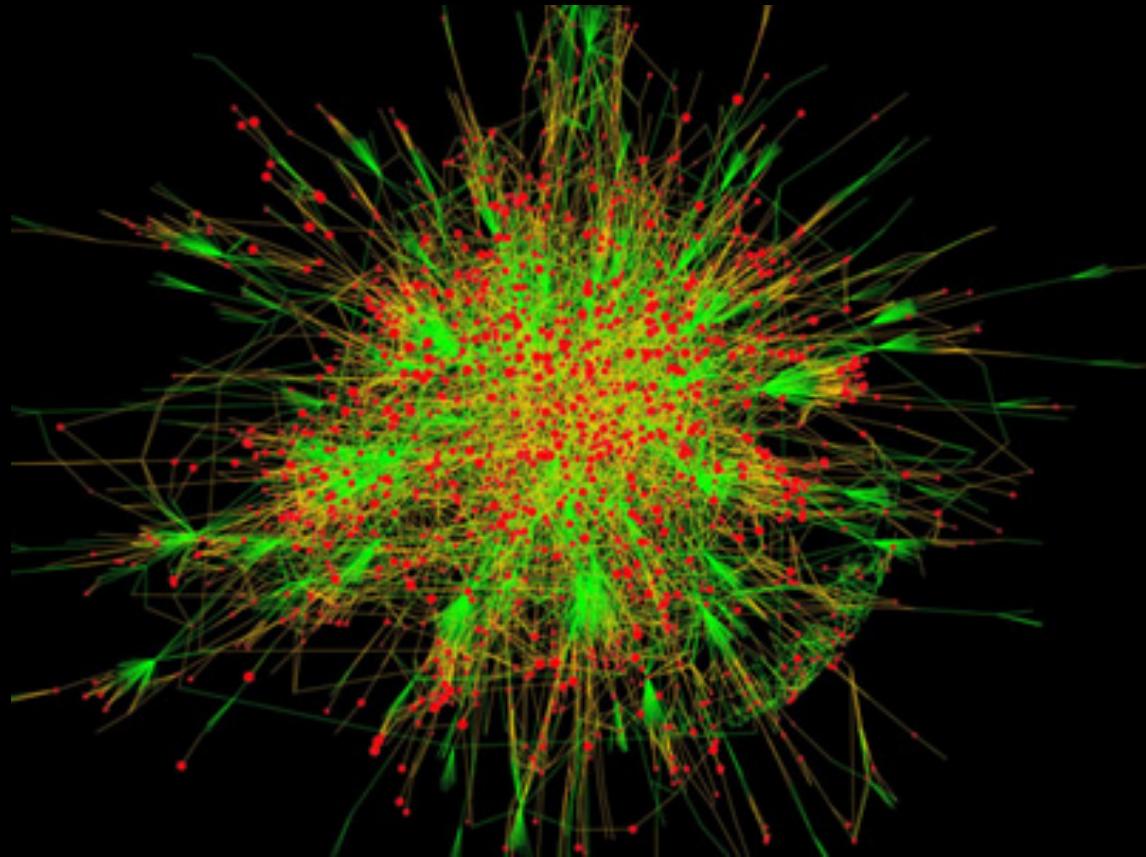
Small datasets
Pretty explicit connections



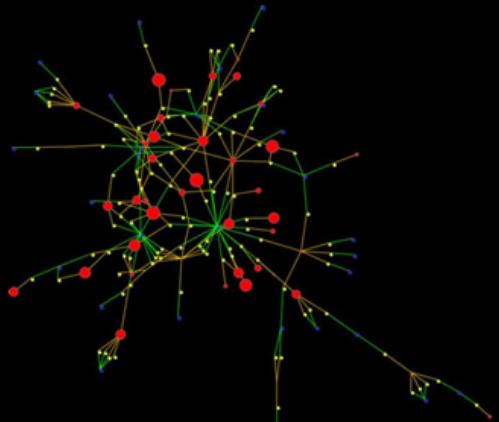
Analysis

Understand the properties

Past



Present



Data Collection / Cleaning



Large datasets
Entity resolution
Implicit connections

Analysis

Understand the properties

Present

Common Tasks

- Measuring “importance”
 - Centrality, prestige (incoming links)
- Link prediction
- Diffusion modeling
 - Epidemiological
- Clustering
 - Blockmodeling, Girvan-Newman
- Structure analysis
 - Motifs, Isomorphisms, etc.
- Visualization/Privacy/etc.

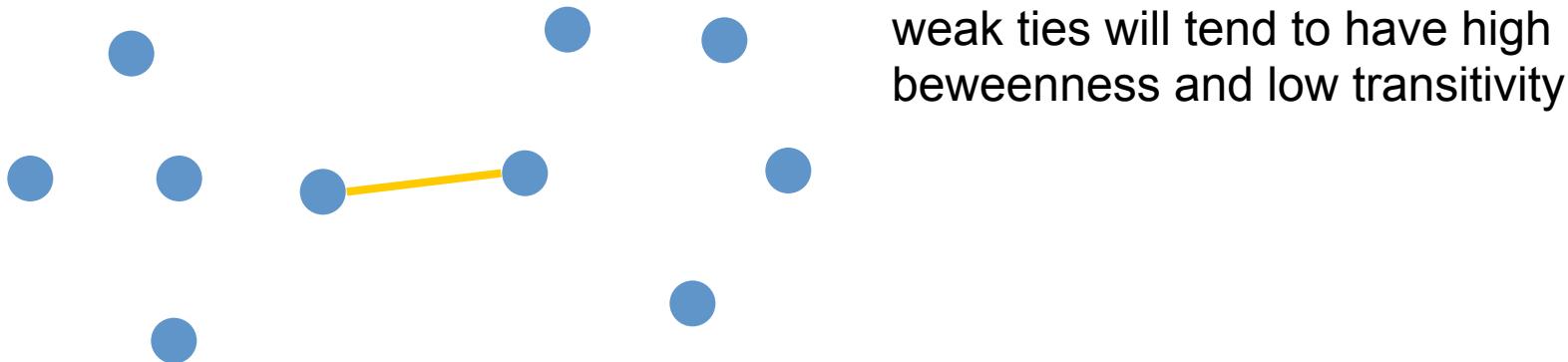
Centrality Measures

- Degree centrality
 - Edges per node (the more, the more important the node)
- Closeness centrality
 - How close the node is to every other node
- Betweenness centrality
 - How many shortest paths go through the edge node (communication metaphor)
- Information centrality
 - All paths to other nodes weighted by path length
- Bibliometric + Internet style
 - PageRank

Tie Strength

- **Strength of Weak Ties (Granovetter)**

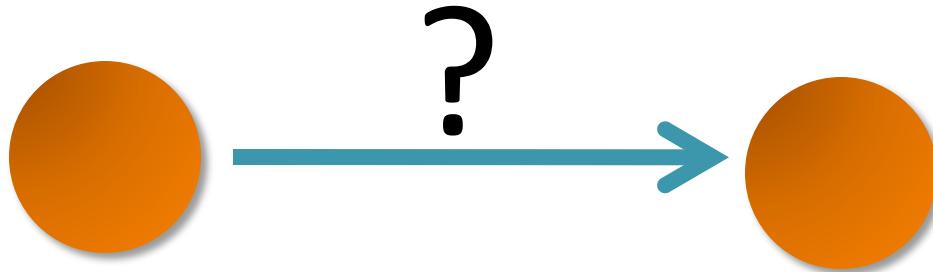
- Granovetter: How often did you see the contact that helped you find the job prior to the job search
 - 16.7 % often (at least once a week)
 - 55.6% occasionally (more than once a year but less than twice a week)
 - 27.8% rarely – once a year or less
- Weak ties will tend to have different *information* than we and our close contacts do



Common Tasks

- Measuring “importance”
 - Centrality, prestige (incoming links)
- Link prediction
- Diffusion modeling
 - Epidemiological
- Clustering
 - Blockmodeling, Girvan-Newman
- Structure analysis
 - Motifs, Isomorphisms, etc.
- Visualization/Privacy/etc.

Link Prediction

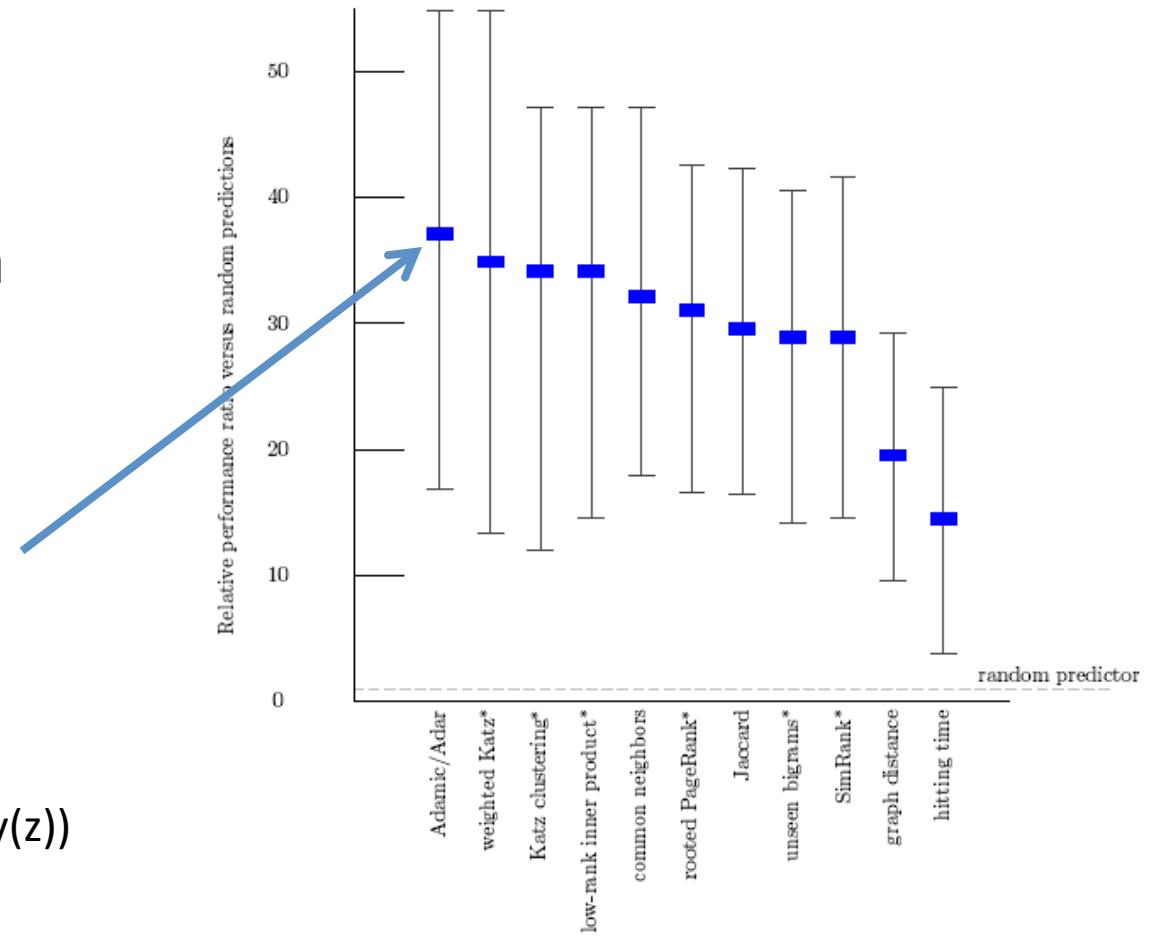


Link Prediction in Social Net Data

- *We know things about structure*
 - *Homophily* = like likes like or bird of a feather flock together or similar people group together
 - *Mutuality*
 - *Triad closure*
- Various measures that try to use this

Link Prediction

- Simple metrics
 - Only take into account graph properties

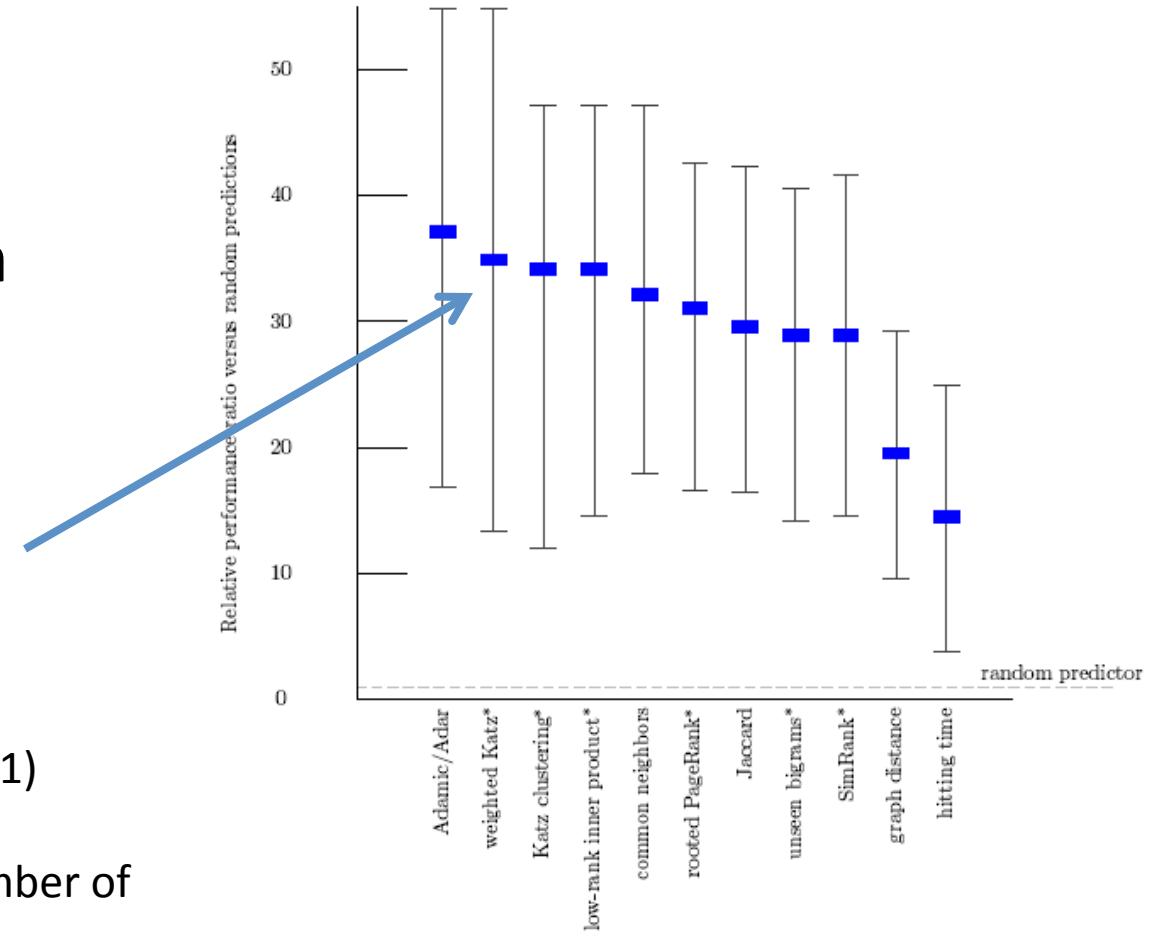


$\Gamma(x)$ = neighbors of x

Originally: $1 / \log(\text{frequency}(z))$

Link Prediction

- Simple metrics
 - Only take into account graph properties



Paths of length ℓ (generally 1)
from x to y
weighted variant is the number of
times the two collaborated

Link Prediction in Relational Data

- *We know things about structure*
 - *Homophily* = like likes like or bird of a feather flock together or similar people group together
 - *Mutuality*
 - *Triad closure*
- Slightly more interesting problem if we have relational data on actors and ties
 - Move beyond structure

Relationship & Link Prediction



Common Tasks

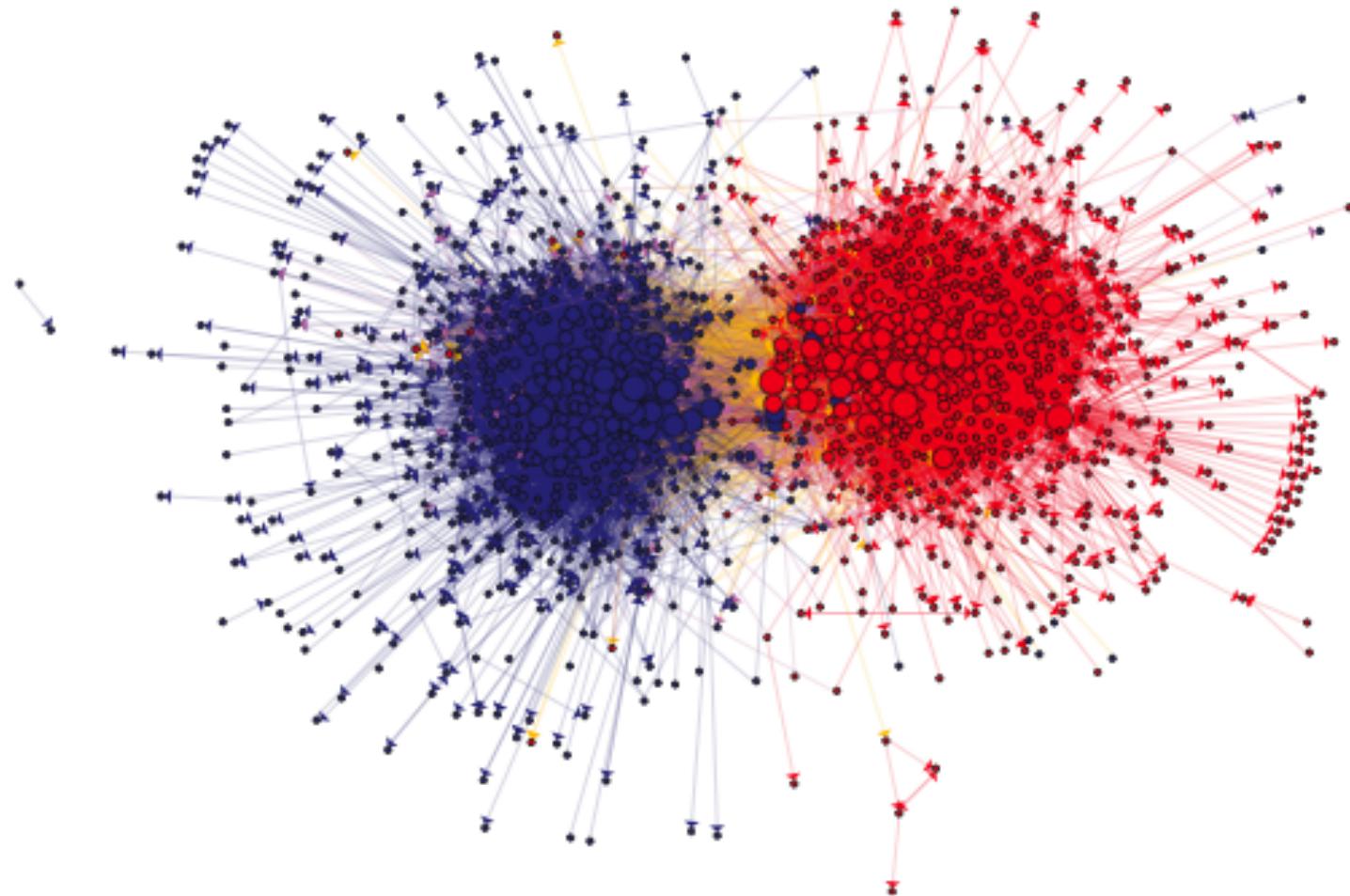
- Measuring “importance”
 - Centrality, prestige (incoming links)
- Link prediction
- Diffusion modeling
 - Epidemiological
- Clustering
 - Blockmodeling, Girvan-Newman
- Structure analysis
 - Motifs, Isomorphisms, etc.
- Visualization/Privacy/etc.

Epidemiological

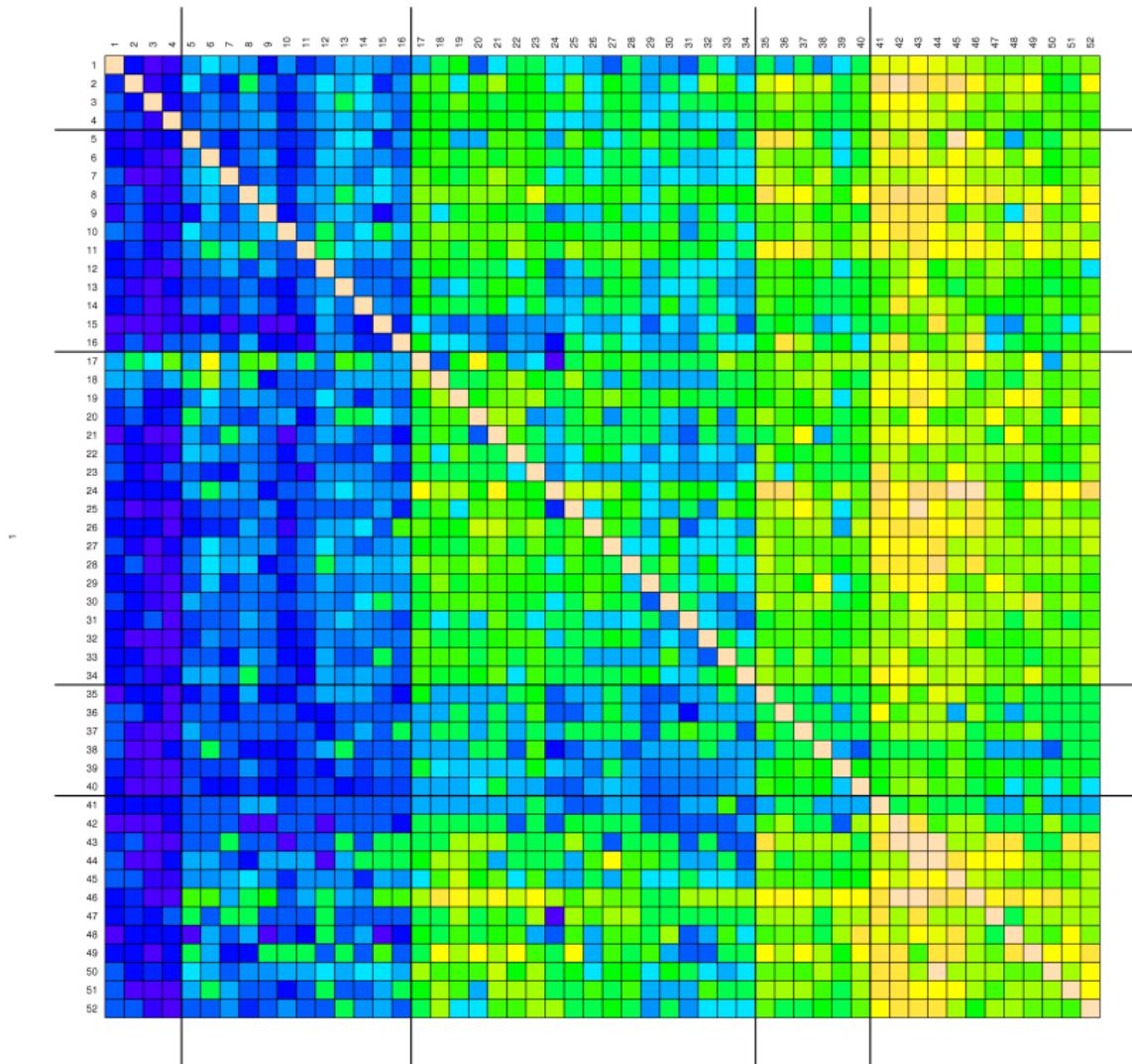
- Viruses
 - Biological, computational
 - STDs, needle sharing, etc.
 - Mark Handcock at UW
- Blog networks
 - Applying SIR models (Info Diffusion Through Blogspace, Gruhl et al.)
 - Induce transmission graph, cascade models, simulation
 - Link prediction (Tracking Information Epidemics in Blogspace, Adar et al.)
 - Find repeated “likely” infections
 - Outbreak detection (Cost-effective Outbreak Detection in Networks, Leskovec et al.)
 - Submodularity

Common Tasks

- Measuring “importance”
 - Centrality, prestige (incoming links)
- Link prediction
- Diffusion modeling
 - Epidemiological
- Clustering/Community Detection
 - Blockmodeling, Girvan-Newman
- Structure analysis
 - Motifs, Isomorphisms, etc.
- Visualization/Privacy/etc.



Blockmodel of U.S. Philosophy Departments. Note that row/column numbers do not correspond to PGR rankings.



Domingo

Carlos

Alejandro

Eduardo

Frank

Hal

Karl

Bob

Ike

Gill

Lanny

Mike

John

Xavier

Utrecht

Norm

Russ

Quint

Wendle

Ozzie

Ted

Sam

Vern

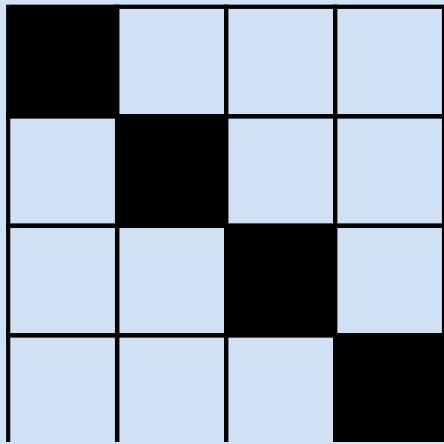
Paul

Blockmodels

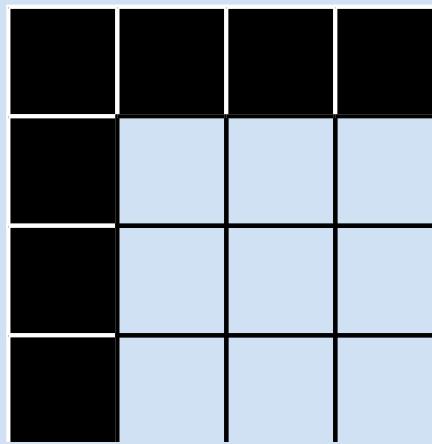
- Actors are portioned into *positions*
 - Rearrange rows/columns
- The sociomatrix is then reduced to a smaller *image*
- Hierarchical clustering
 - Various distance metrics
 - Euclidean, CONvergence of CORrelation (CONCOR)
 - Various “fit” metrics

Image matrix

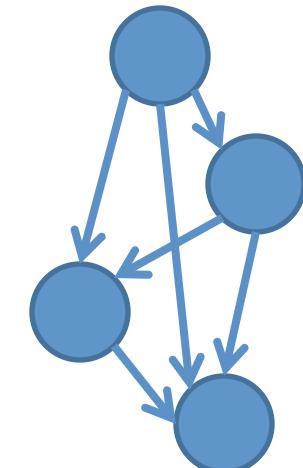
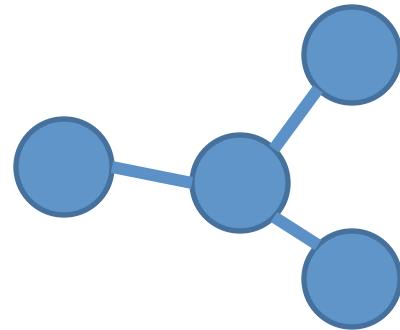
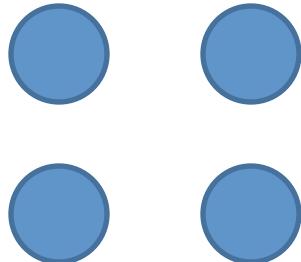
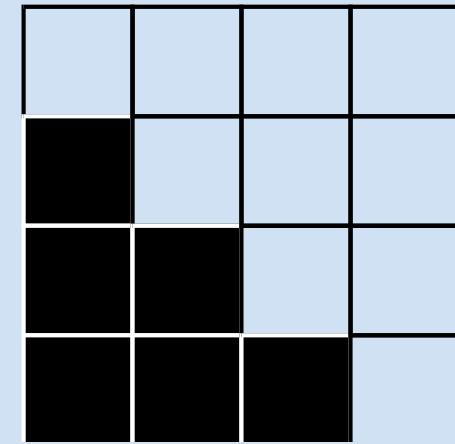
Cohesion



Center-periphery



Ranking



Girvan-Newman Algorithm

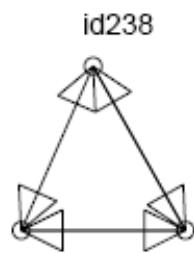
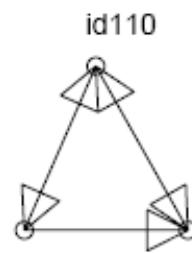
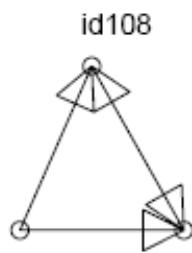
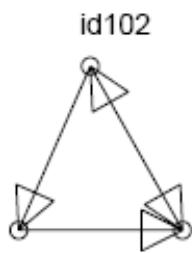
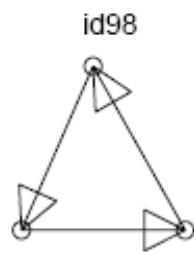
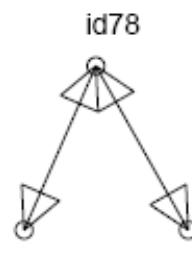
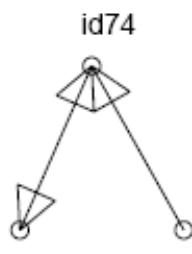
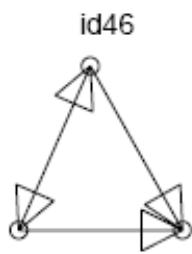
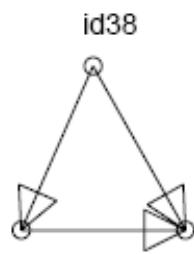
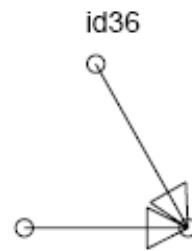
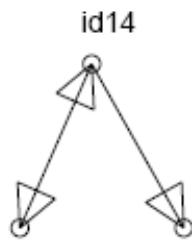
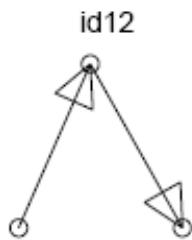
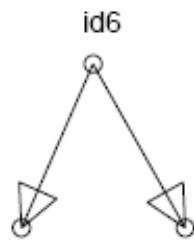
- Split on shortest paths (“weak ties”)
1. Calculate betweenness on all edges
 2. Remove highest betweenness edge
 3. Recalculate
 4. Goto 1

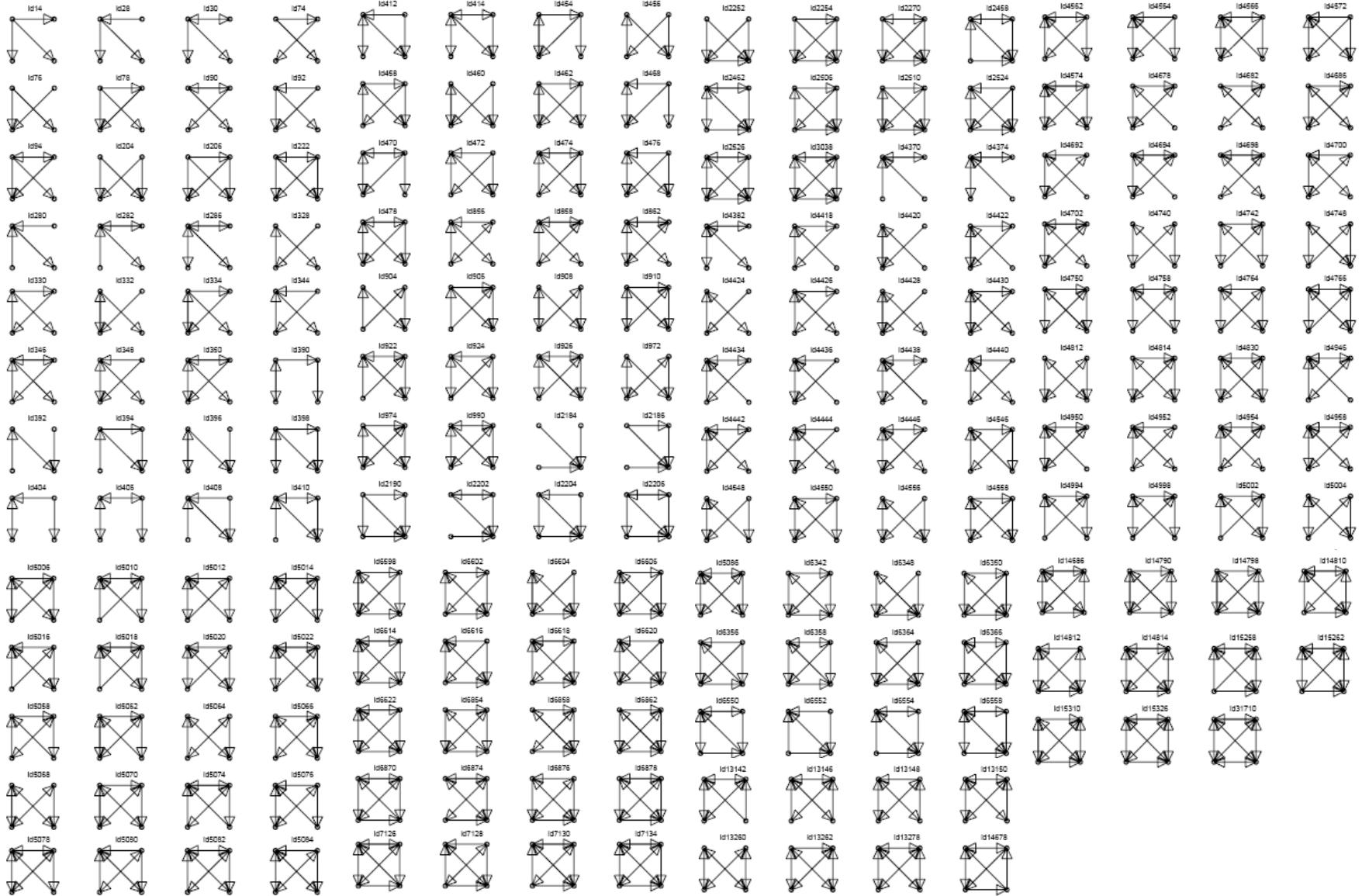
Other solutions

- Min-cut based
- “Voltage” based
- Hierarchical schemes

Common Tasks

- Measuring “importance”
 - Centrality, prestige (incoming links)
- Link prediction
- Diffusion modeling
 - Epidemiological
- Clustering
 - Blockmodeling, Girvan-Newman
- Structure analysis
 - Motifs, Isomorphisms, etc.
- Visualization/Privacy/etc.





Network motif detection

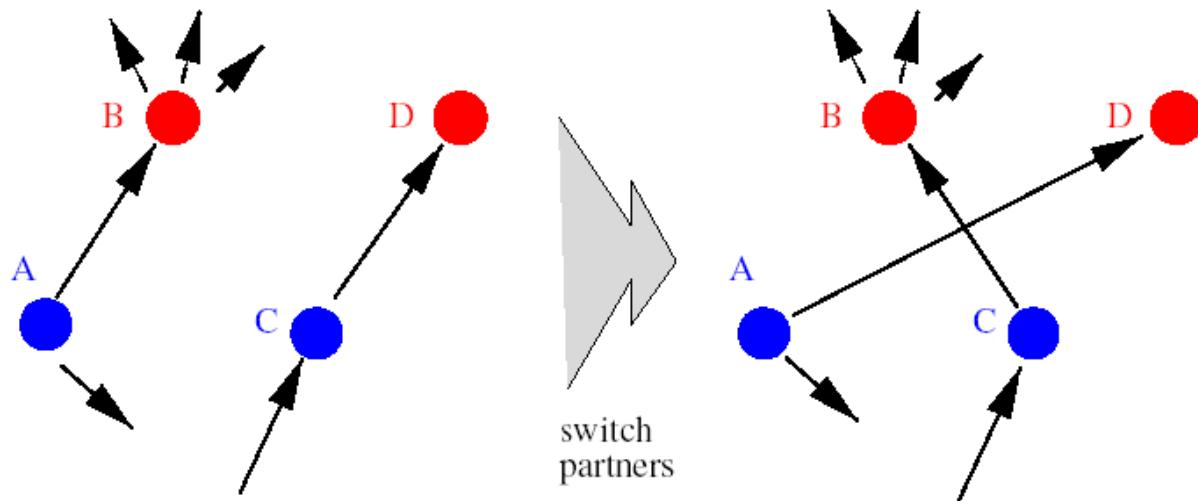
- How many more motifs of a certain type exist over a random network
- Started in biological networks
 - <http://www.weizmann.ac.il/mcb/UriAlon/>

Basic idea

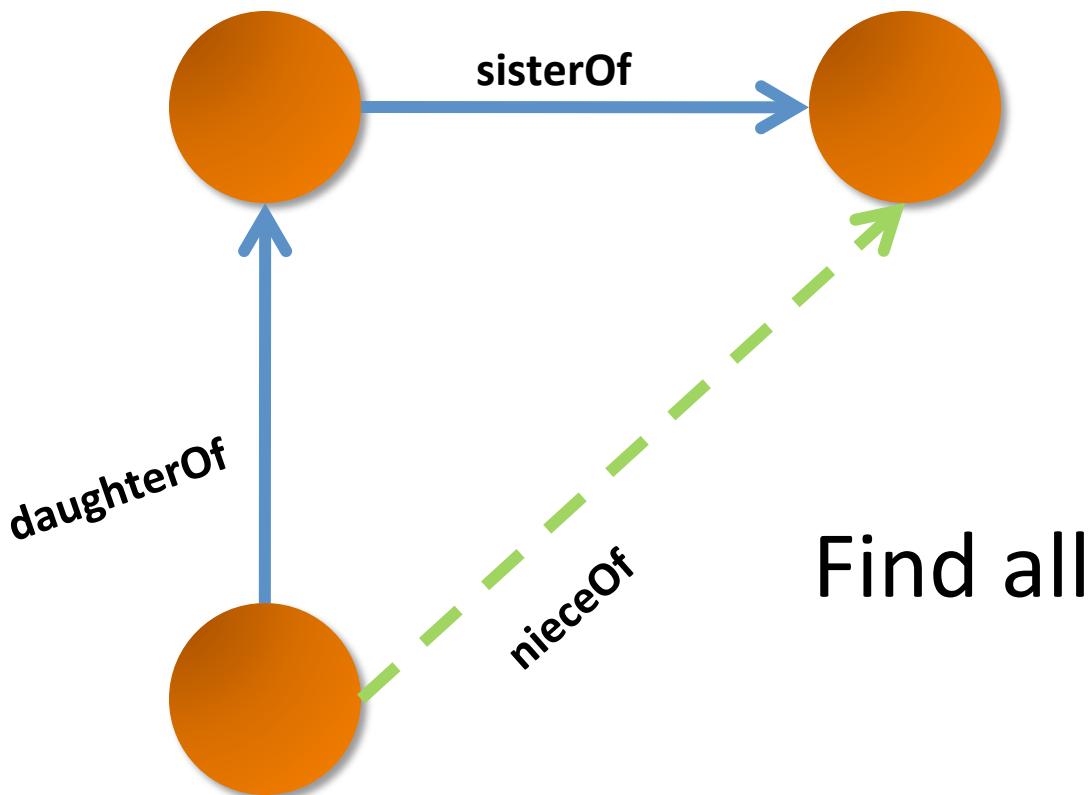
- construct many random graphs with the same number of nodes and edges (same node degree distribution?)
- count the number of motifs in those graphs
- calculate the Z score: the probability that the given number of motifs in the real world network could have occurred by chance

Generating random graphs

- Many models don't preserve the desired features
- Have to be careful how we generate



Other Structural Analysis



Common Tasks

- Measuring “importance”
 - Centrality, prestige (incoming links)
- Link prediction
- Diffusion modeling
 - Epidemiological
- Clustering
 - Blockmodeling, Girvan-Newman
- Structure analysis
 - Motifs, Isomorphisms, etc.
- Visualization/Privacy/etc.

Privacy

- Emerging interest in anonymizing networks
 - Lars Backstrom (WWW'07) demonstrated one of the first attacks
- How to remove labels while preserving graph properties?
 - While ensuring that labels cannot be reapplied

Software

- Pajek
 - <http://vlado.fmf.uni-lj.si/pub/networks/pajek/>
- UCINET
 - <http://www.analytictech.com/>
- KrackPlot
 - <http://www.andrew.cmu.edu/user/krack/krackplot.shtml>
- GUESS
 - <http://www.graphexploration.org>
- Etc.

Books/Journals/Conferences

- *Social Networks/Phs. Rev*
- Social Network Analysis (Wasserman + Faust)
- The Development of Social Network Analysis (Freeman)
- Linked (Barabsi)
- Six Degrees (Watts)
- Sunbelt/ICWSM/KDD/CIKM/NIPS

Assortativity

- Social networks are assortative:
 - the gregarious people associate with other gregarious people
 - the loners associate with other loners
- The Internet is disassortative:

