



UCSC

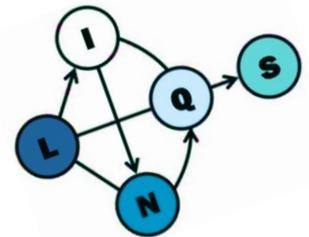
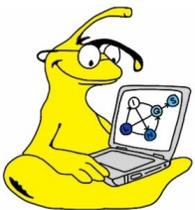


# Collective Graph Identification

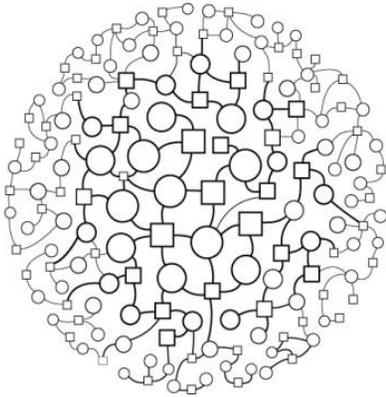
Lise Getoor

University of California, Santa Cruz

Workshop on Incomplete Networked Data  
March 23, 2016







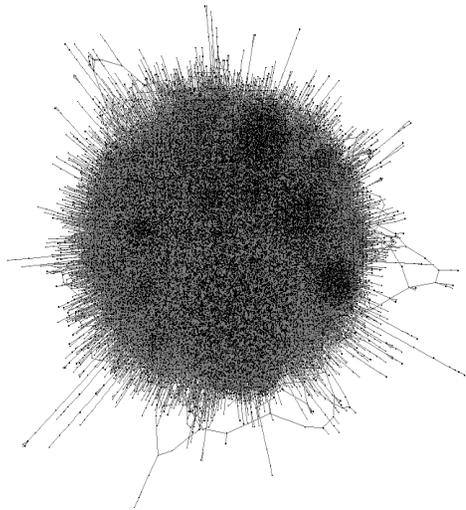
+

**Graph  
Analytics**

=

**Insights!**





+

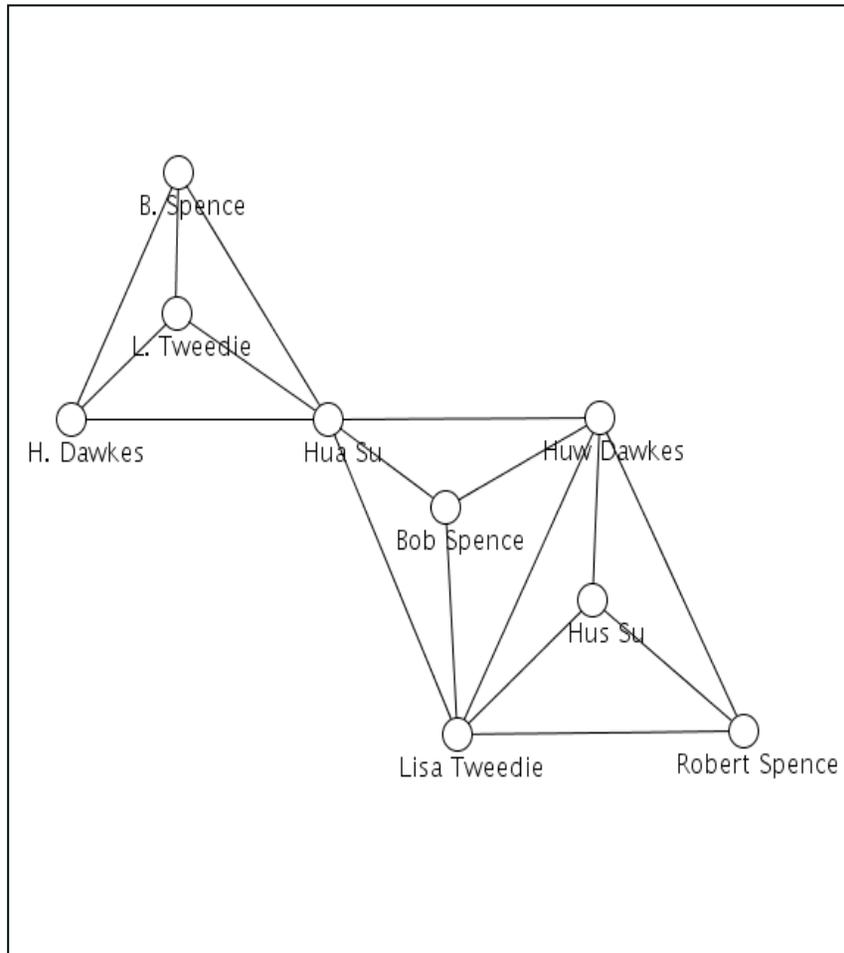
Graph  
Analytics

=

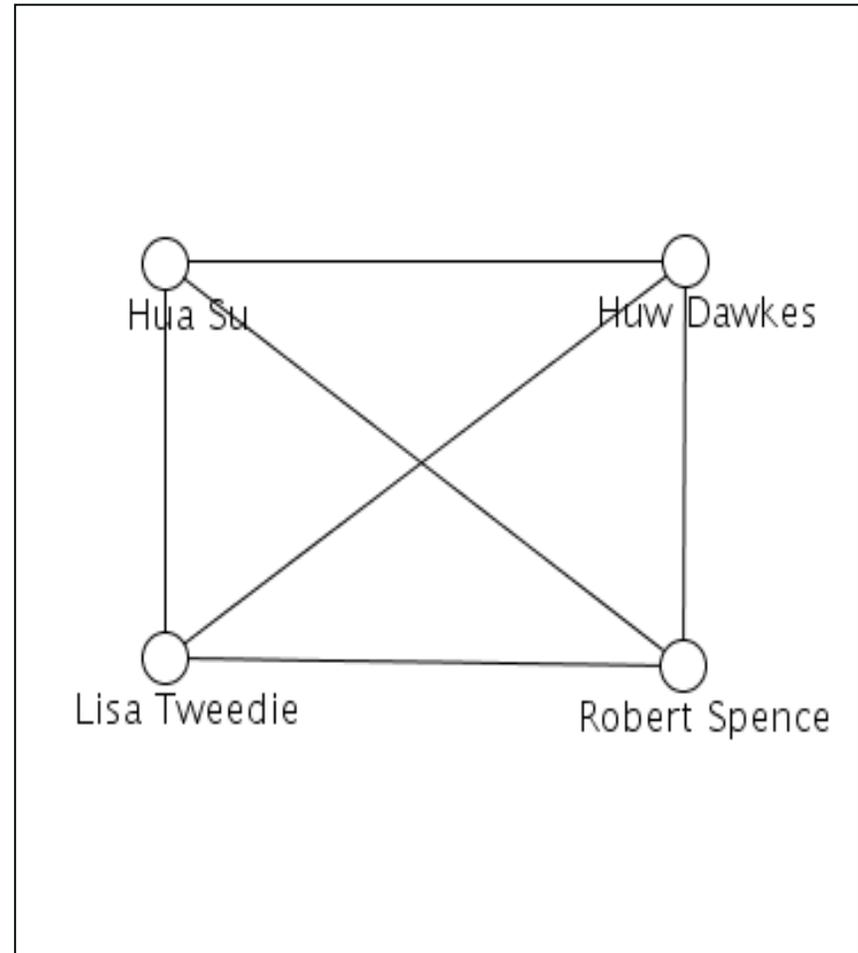
1+1=3



# Co-Author Graph



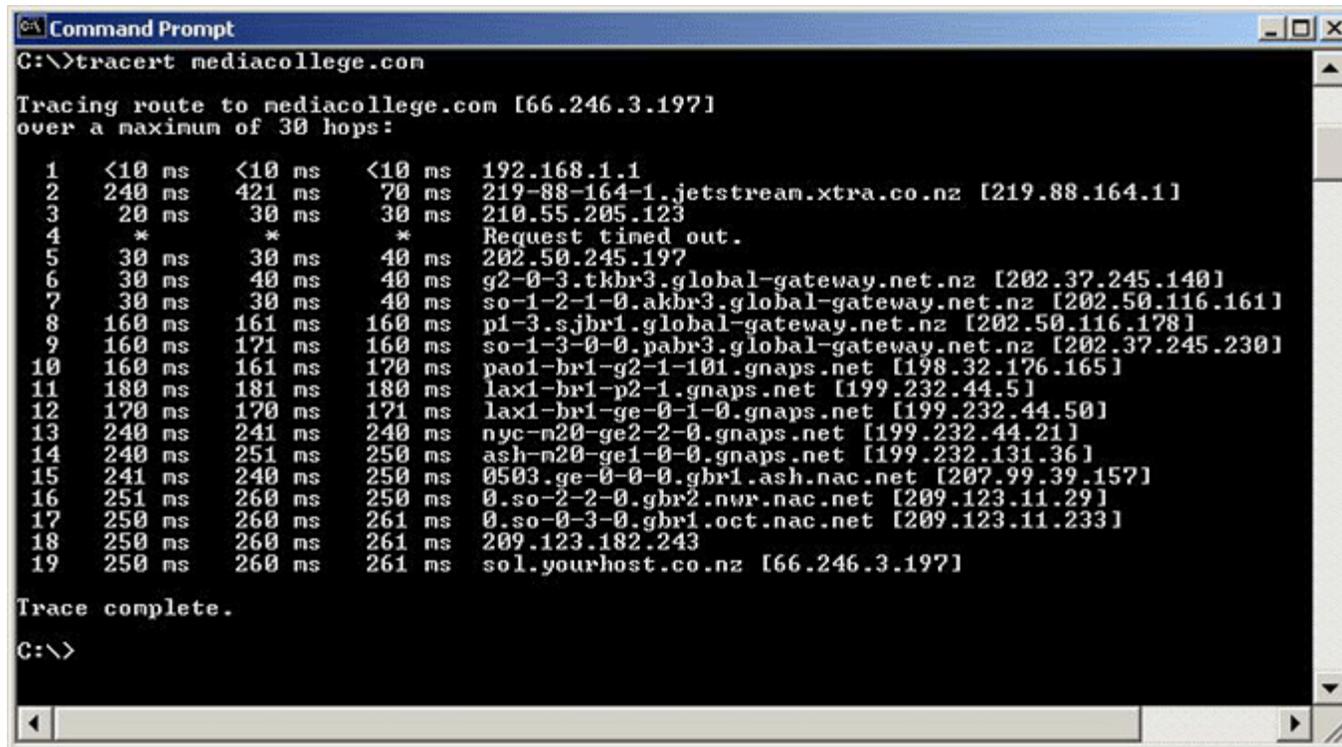
before



after

# Motivation: ER and Network Analysis

- o Measuring the topology of the internet ... using traceroute



```
C:\>tracert mediacollege.com

Tracing route to mediacollege.com [66.246.3.197]
over a maximum of 30 hops:

  0  <10 ms  <10 ms  <10 ms  192.168.1.1
  1  240 ms  421 ms  70 ms  219-88-164-1.jetstream.xtra.co.nz [219.88.164.1]
  2  20 ms   30 ms   30 ms  210.55.205.123
  3  *       *       *      Request timed out.
  4  30 ms   30 ms   40 ms  202.50.245.197
  5  30 ms   40 ms   40 ms  g2-0-3.tkbr3.global-gateway.net.nz [202.37.245.140]
  6  30 ms   30 ms   40 ms  so-1-2-1-0.akbr3.global-gateway.net.nz [202.50.116.161]
  7  160 ms  161 ms  160 ms  p1-3.sjbr1.global-gateway.net.nz [202.50.116.178]
  8  160 ms  171 ms  160 ms  so-1-3-0-0.pabr3.global-gateway.net.nz [202.37.245.230]
  9  160 ms  161 ms  170 ms  pao1-br1-g2-1-101.gnaps.net [198.32.176.165]
 10  180 ms  181 ms  180 ms  lax1-br1-p2-1.gnaps.net [199.232.44.5]
 11  170 ms  170 ms  171 ms  lax1-br1-ge-0-1-0.gnaps.net [199.232.44.50]
 12  240 ms  241 ms  240 ms  nyc-n20-ge2-2-0.gnaps.net [199.232.44.21]
 13  240 ms  251 ms  250 ms  ash-n20-ge1-0-0.gnaps.net [199.232.131.36]
 14  241 ms  240 ms  250 ms  0503.ge-0-0-0.gbr1.ash.nac.net [207.99.39.157]
 15  251 ms  260 ms  250 ms  0.so-2-2-0.gbr2.nwr.nac.net [209.123.11.29]
 16  250 ms  260 ms  261 ms  0.so-0-3-0.gbr1.oct.nac.net [209.123.11.233]
 17  250 ms  260 ms  261 ms  209.123.182.243
 18  250 ms  260 ms  261 ms  sol.yourhost.co.nz [66.246.3.197]
 19  250 ms  260 ms  261 ms

Trace complete.

C:\>
```

# IP Aliasing Problem [Willinger et al. 2009]

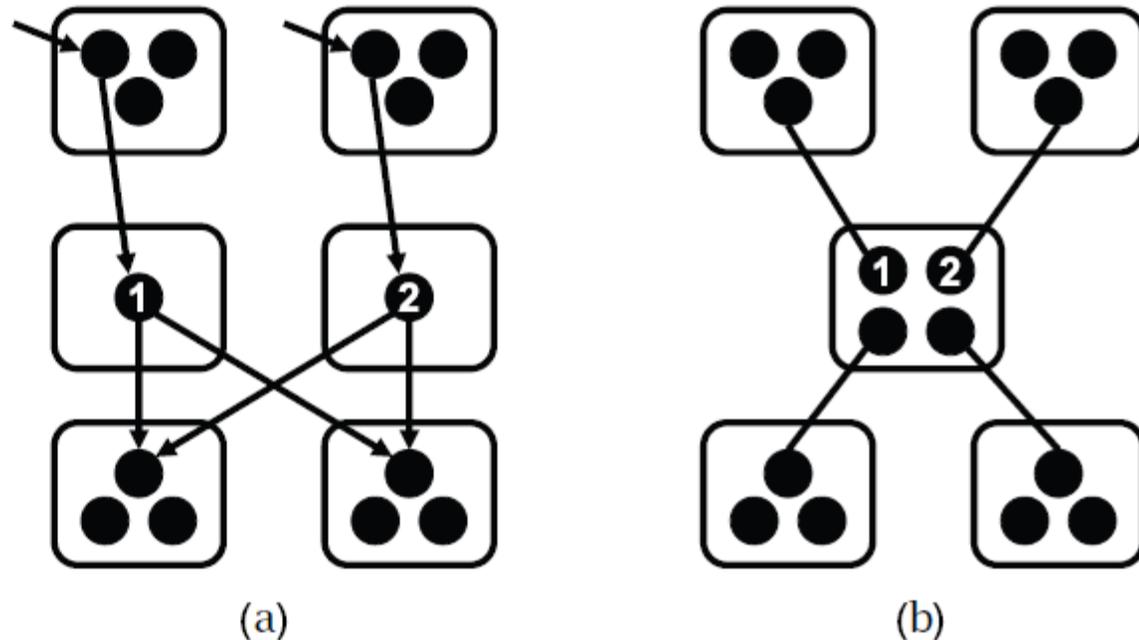


Figure 2. The IP alias resolution problem. Paraphrasing Fig. 4 of [50], traceroute does not list routers (boxes) along paths but IP addresses of input interfaces (circles), and alias resolution refers to the correct mapping of interfaces to routers to reveal the actual topology. In the case where interfaces 1 and 2 are aliases, (b) depicts the actual topology while (a) yields an “inflated” topology with more routers and links than the real one.

# IP Aliasing Problem [Willinger et al. 2009]

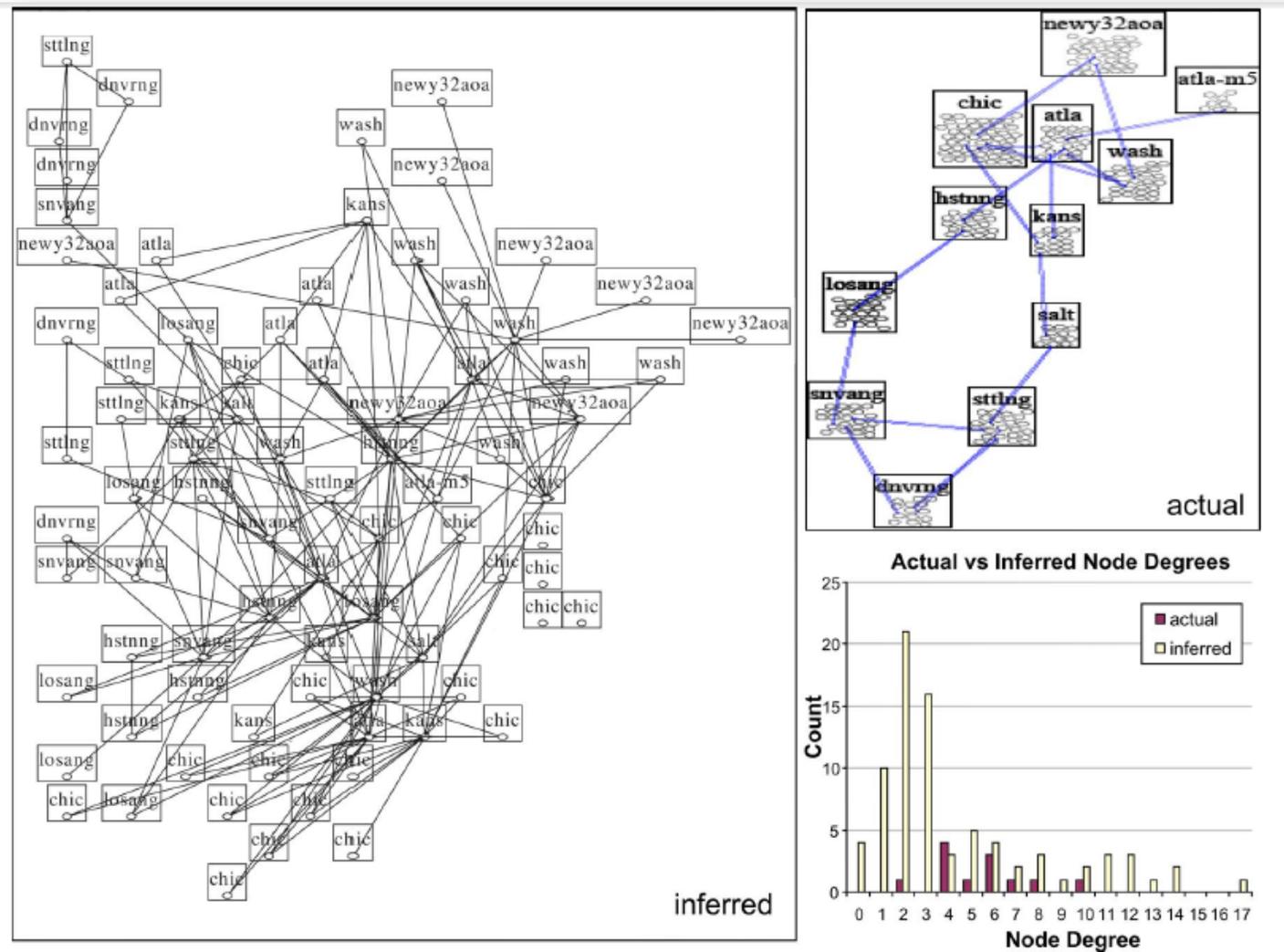
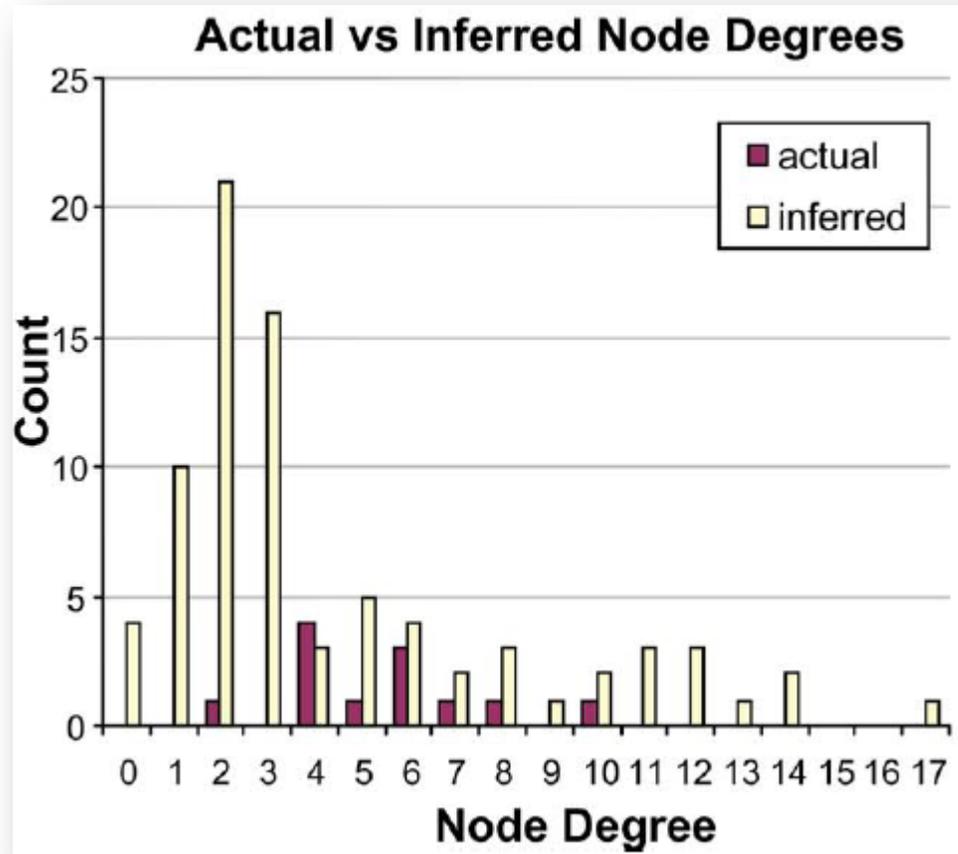


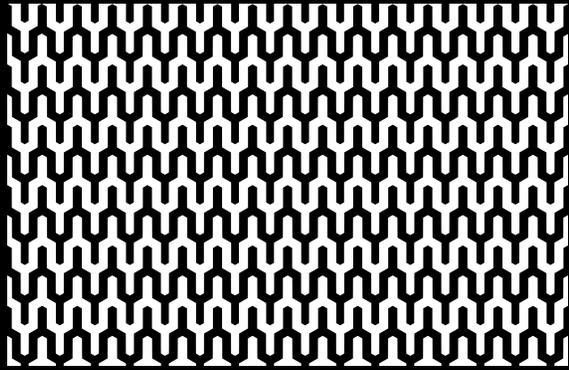
Figure 3. The IP alias resolution problem in practice. This is re-produced from [48] and shows a comparison between the Abilene/Internet2 topology inferred by Rocketfuel (left) and the actual topology (top right). Rectangles represent routers with interior ovals denoting interfaces. The histograms of the corresponding node degrees are shown in the bottom right plot. © 2008 ACM,

# IP Aliasing Problem [Willinger et al. 2009]



Lesson: Make sure you are working on the right graph before performing analytics!

How do you get the right graph?  
Infer it from the data!



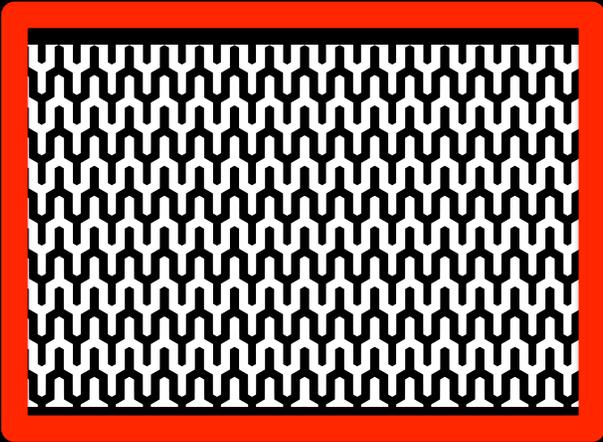
Patterns



Key Ideas



Tools



Patterns



Key Ideas



Tools

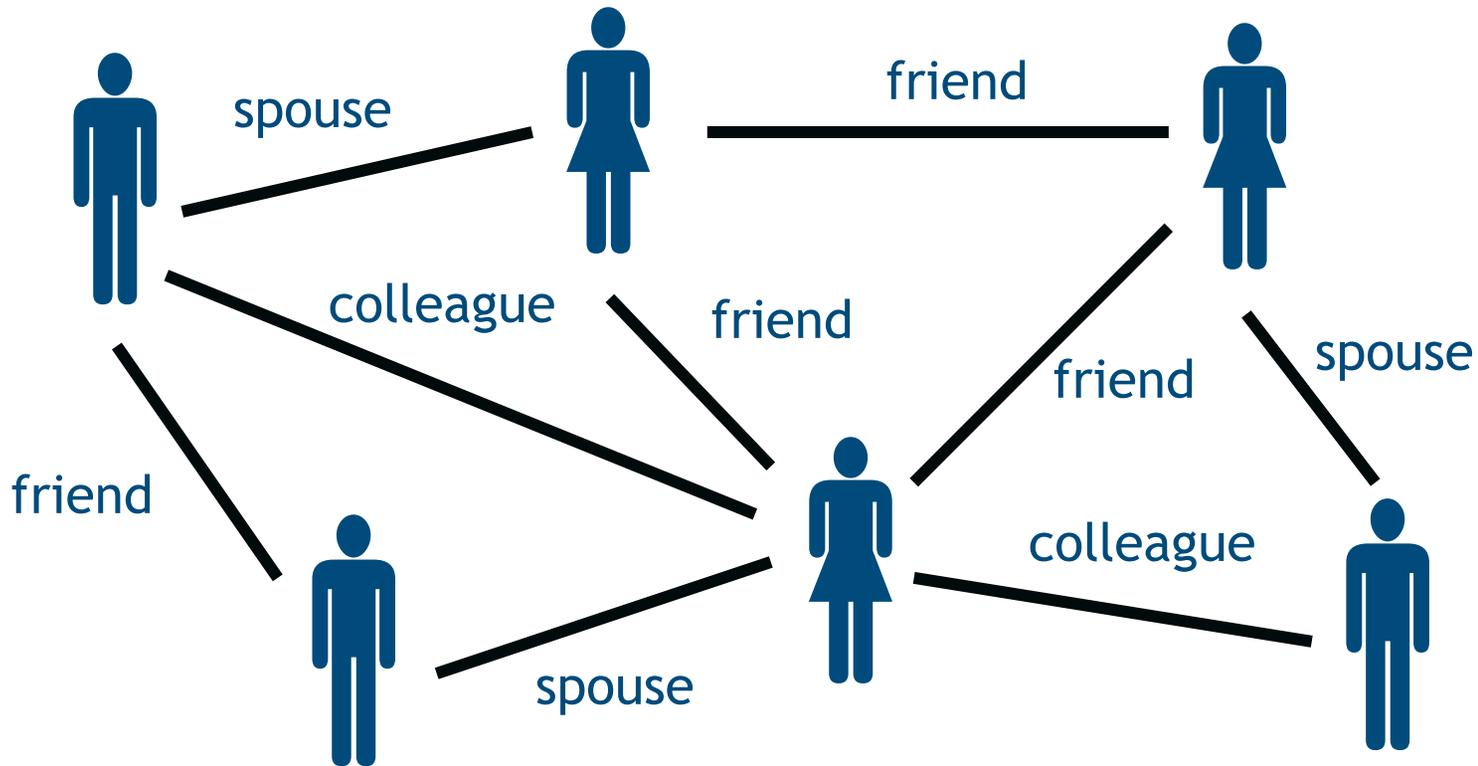
# Graph Inference Patterns

- Entity Resolution
- Collective Classification
- Link Prediction

**Entity Resolution:** determining which nodes refer to same underlying entity

**Collective Classification:** inferring  
the labels of nodes in a graph

# Collective Classification



Question:

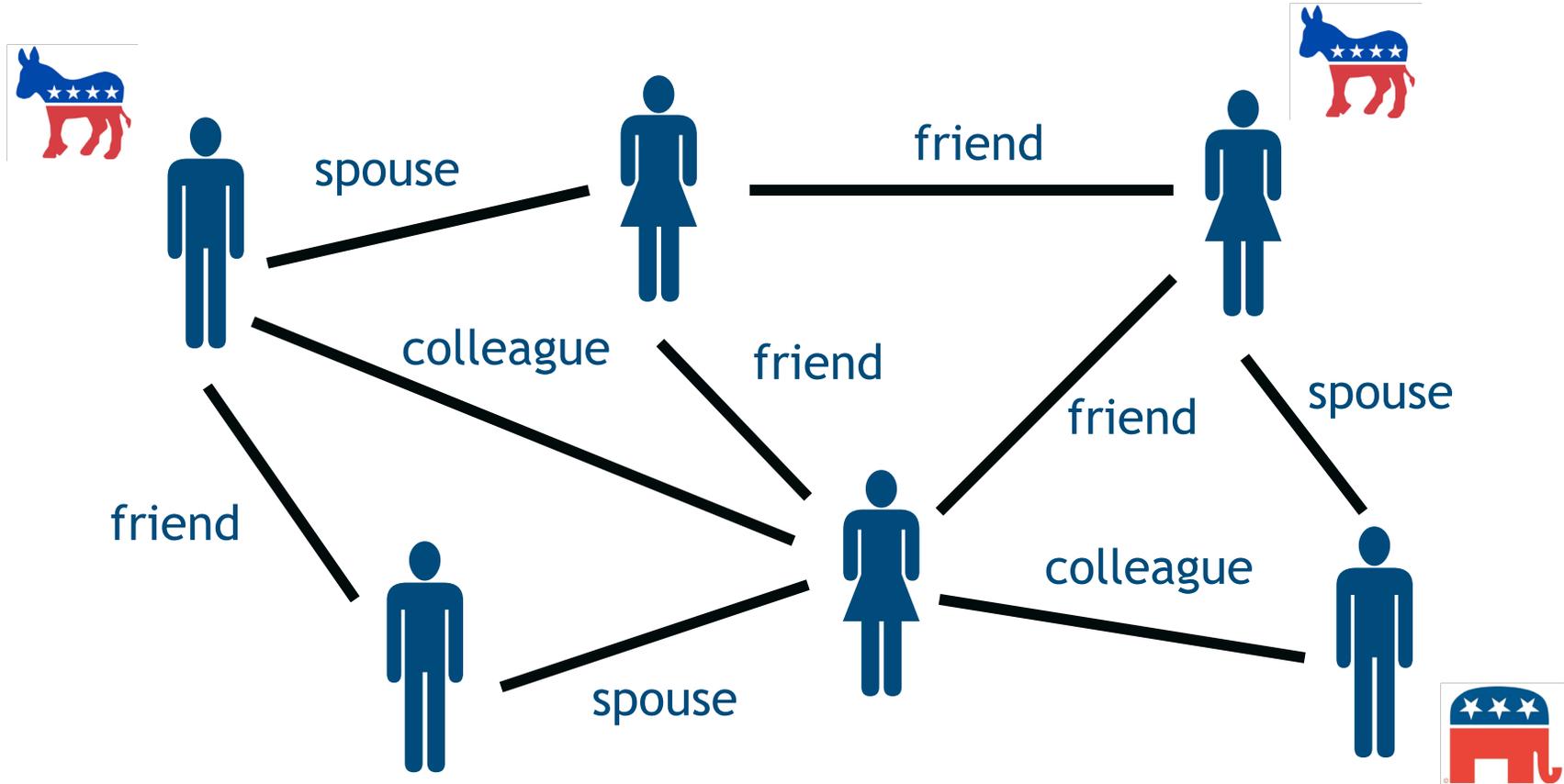


or

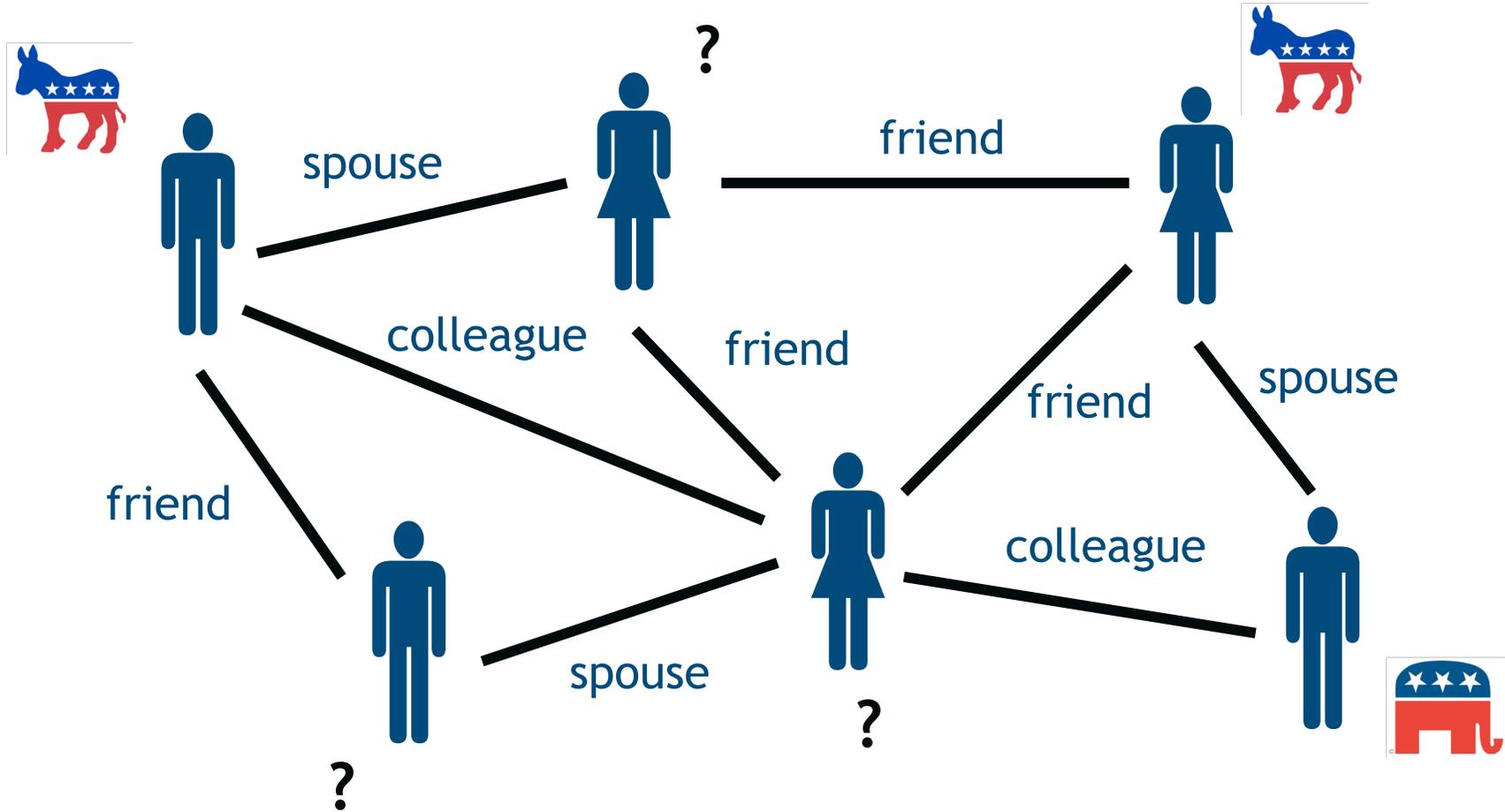


?

# Collective Classification



# Collective Classification



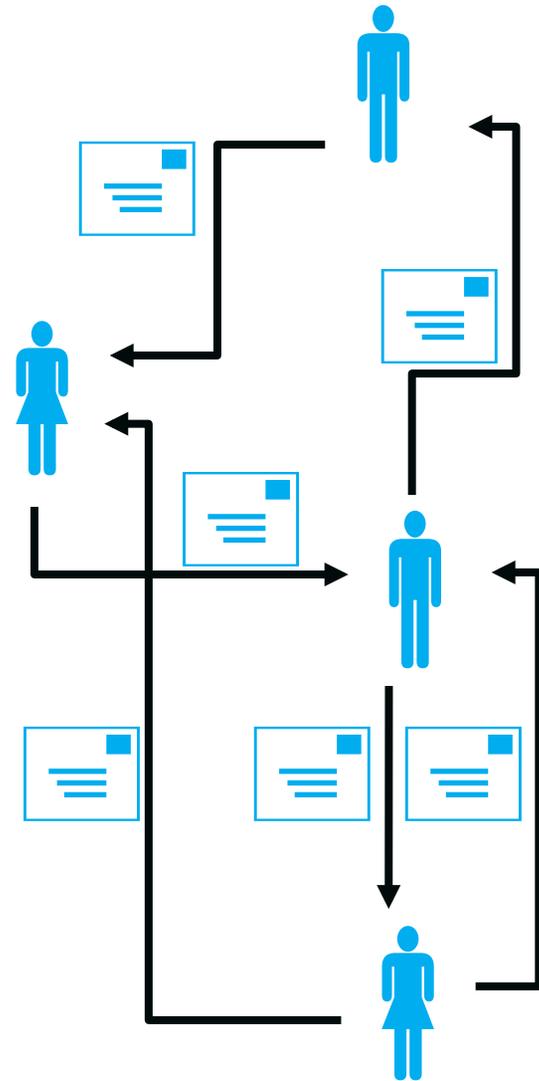
## Common Graph Inference Patterns

- Collective Classification
- Link Prediction
- Entity Resolution

**Link Prediction:** inferring the existence of edges in a graph

# Link Prediction

- Entities
  - People, Emails
- Observed relationships
  - communications, co-location
- Predict relationships
  - Supervisor, subordinate, colleague



# Graph Inference Patterns

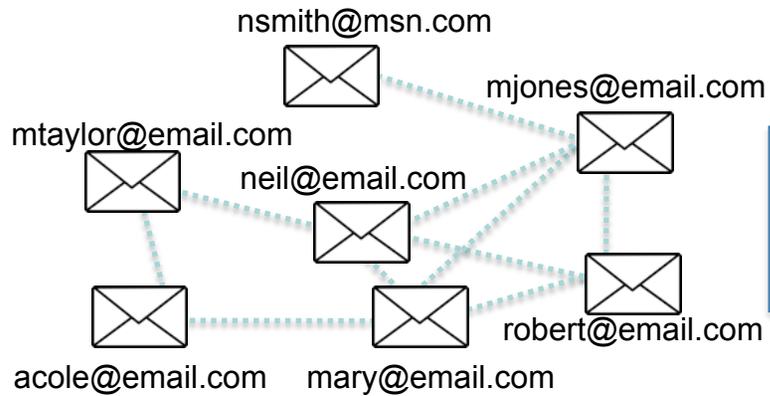
- Entity Resolution
- Collective Classification
- Link Prediction

**My favorite  
Graph Inference Pattern**

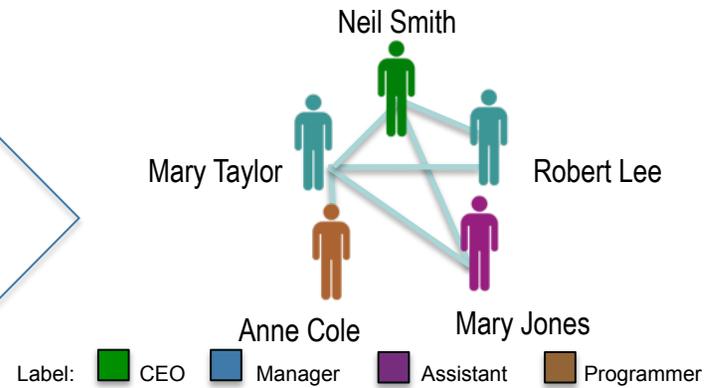
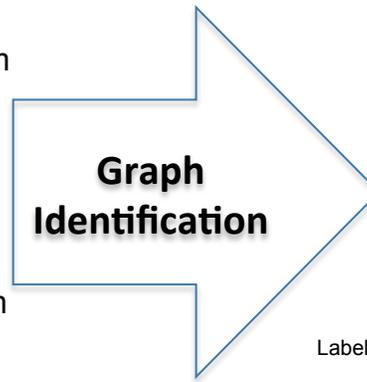
# Graph Identification

- Goal:
  - Given an **input graph** infer an **output graph**
- Three major components:
  - **Entity Resolution (ER)**: Infer the set of nodes
  - **Link Prediction (LP)**: Infer the set of edges
  - **Collective Classification (CC)**: Infer the node labels
- Challenge: The components are intra and inter-dependent

# Graph Identification

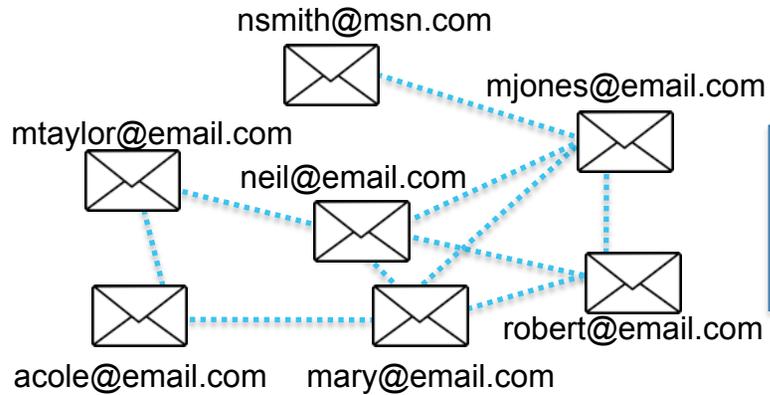


Input Graph: Email Communication Network

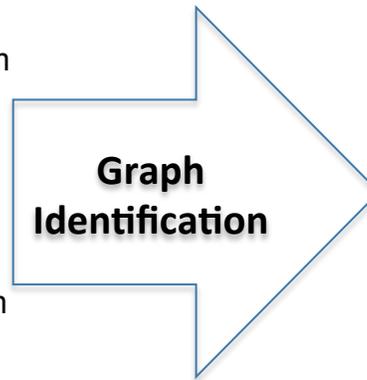


Output Graph: Social Network

# Graph Identification



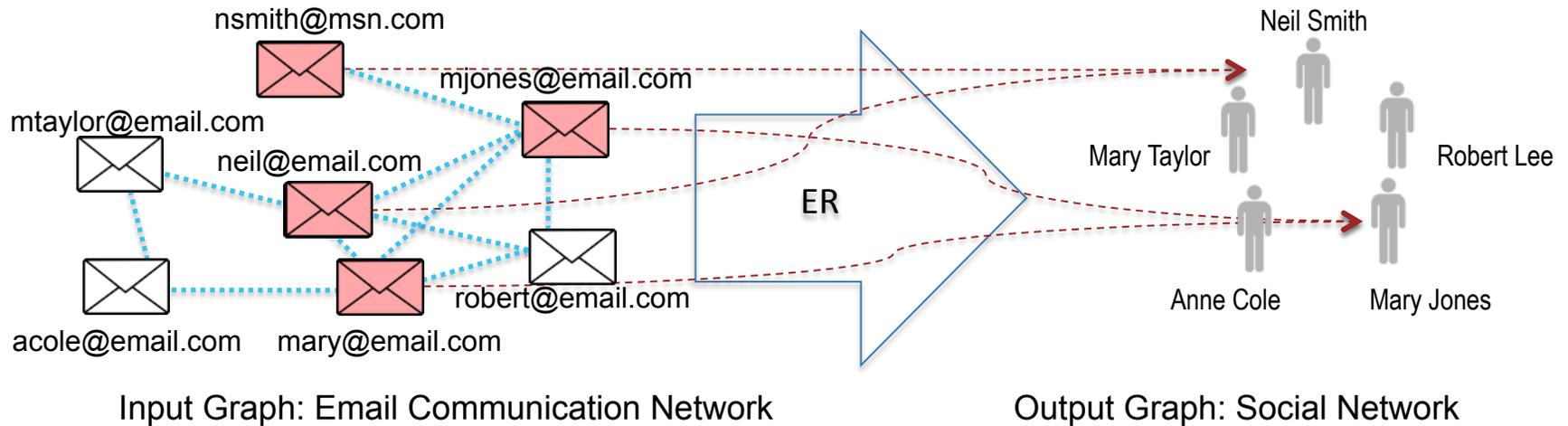
Input Graph: Email Communication Network



Output Graph: Social Network

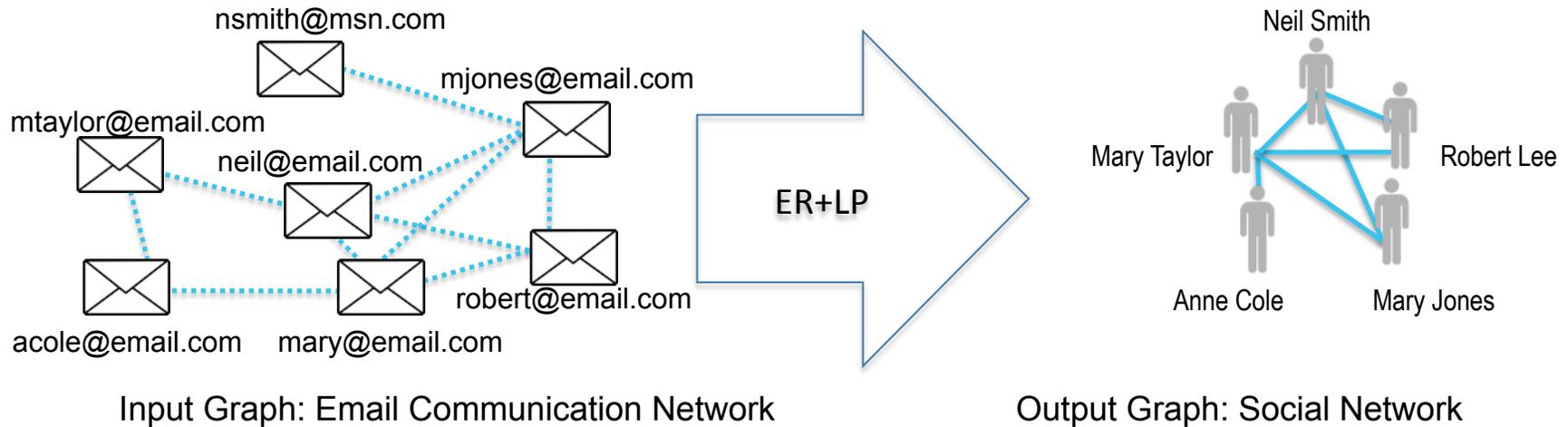
- What's involved?

# Graph Identification



- What's involved?
  - Entity Resolution (ER): Map input graph nodes to output graph nodes

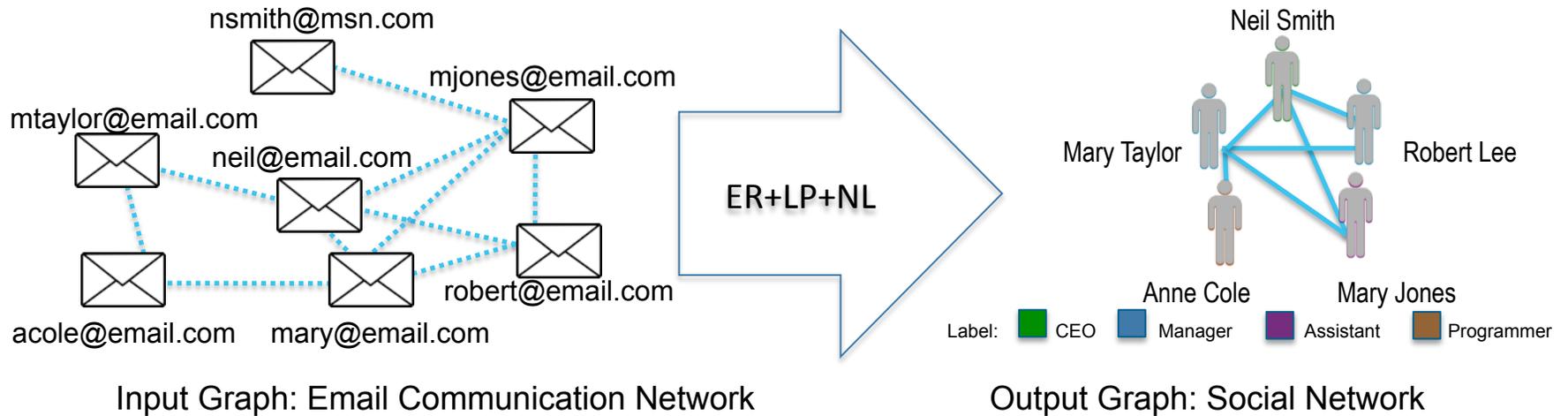
# Graph Identification



- What's involved?

- Entity Resolution (ER): Map input graph nodes to output graph nodes
- Link Prediction (LP): Predict existence of edges in output graph

# Graph Identification

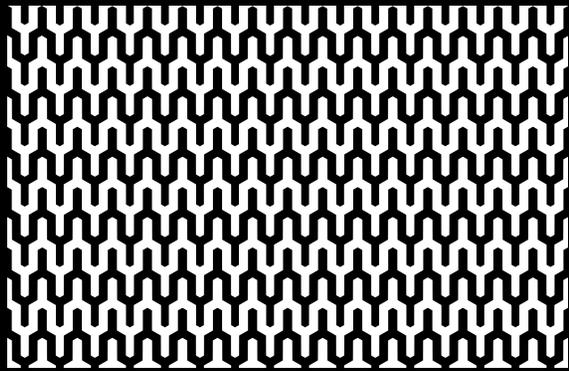


## •What's involved?

- Entity Resolution (ER): Map input graph nodes to output graph nodes
- Link Prediction (LP): Predict existence of edges in output graph
- Node Labeling (NL): Infer the labels of nodes in the output graph

# Graph Identification

- Goal:
  - Given an **input graph** infer an **output graph**
- Three major components:
  - **Entity Resolution (ER)**: Infer the set of nodes
  - **Link Prediction (LP)**: Infer the set of edges
  - **Collective Classification (CC)**: Infer the node labels
- Challenge: The components are intra and inter-dependent



Patterns



Key Ideas

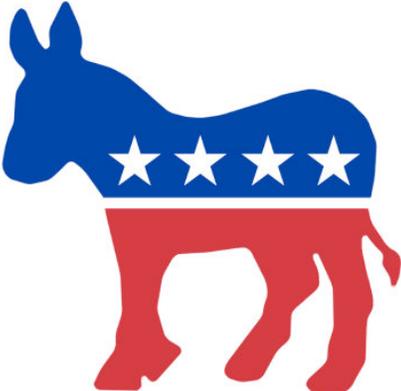
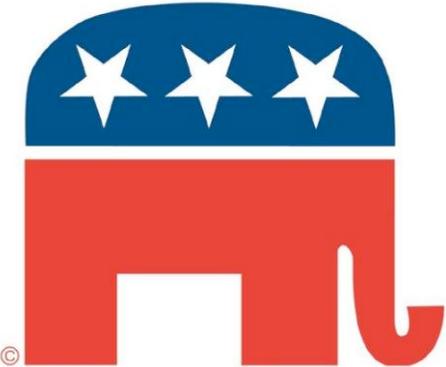


Tools

*richly structured*

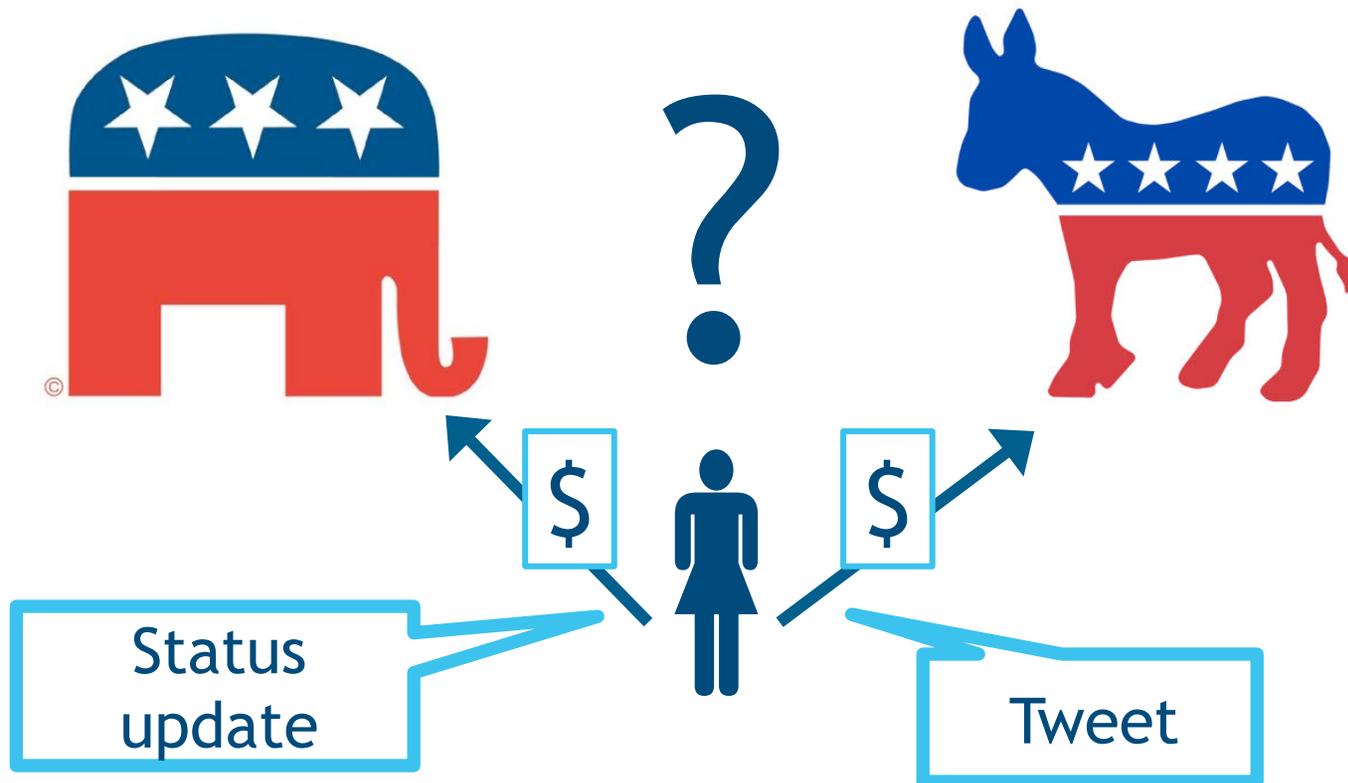
Collective reasoning over ~~complex~~  
graphs

# Collective Classification



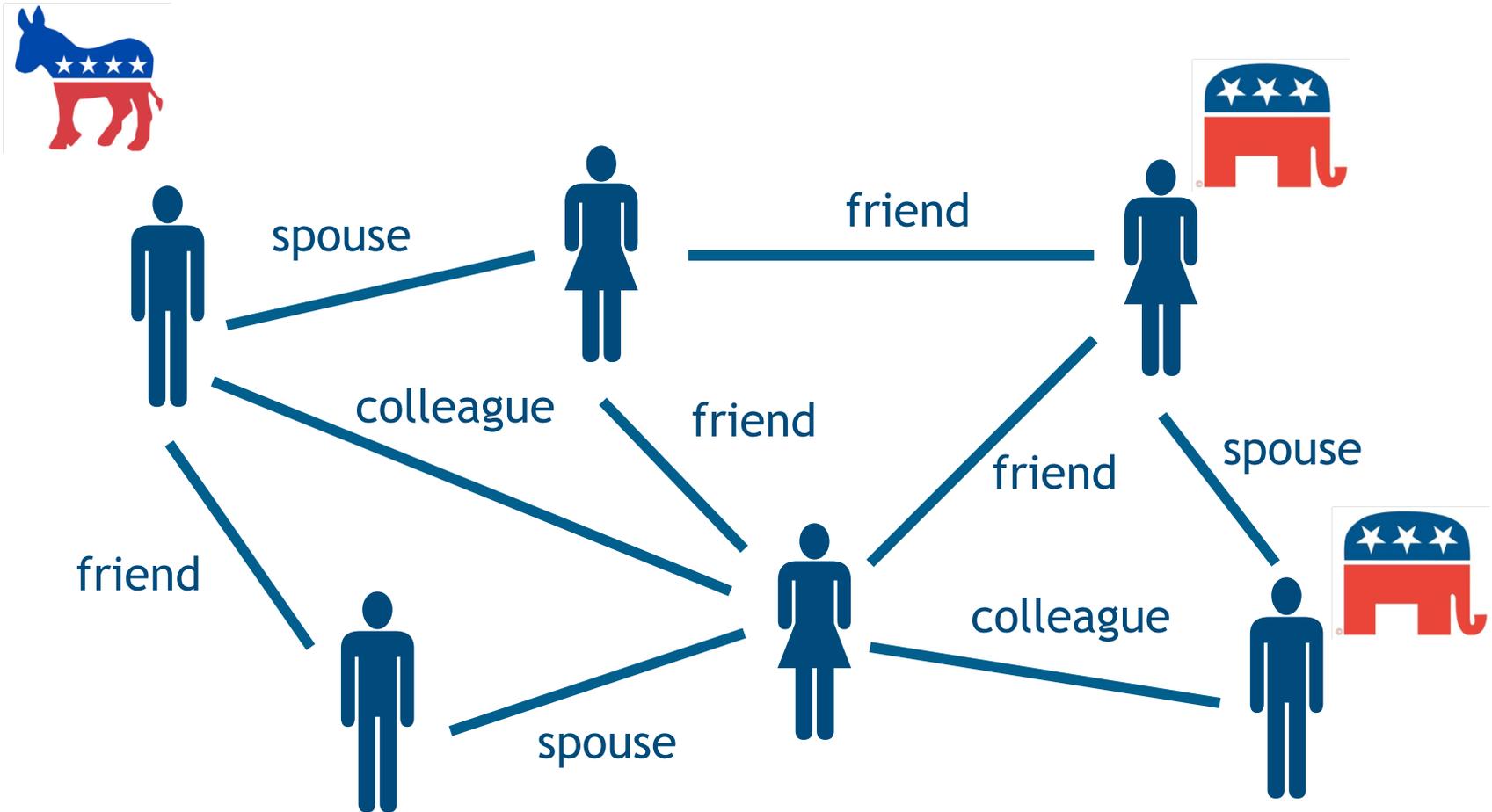
# Collective Classification

Donates(A, 🇺🇸🐘) => Votes(A, 🇺🇸🐘): 5.0



Mentions(A, "Affordable Health") => Votes(A, 🇺🇸🐘): 0.3

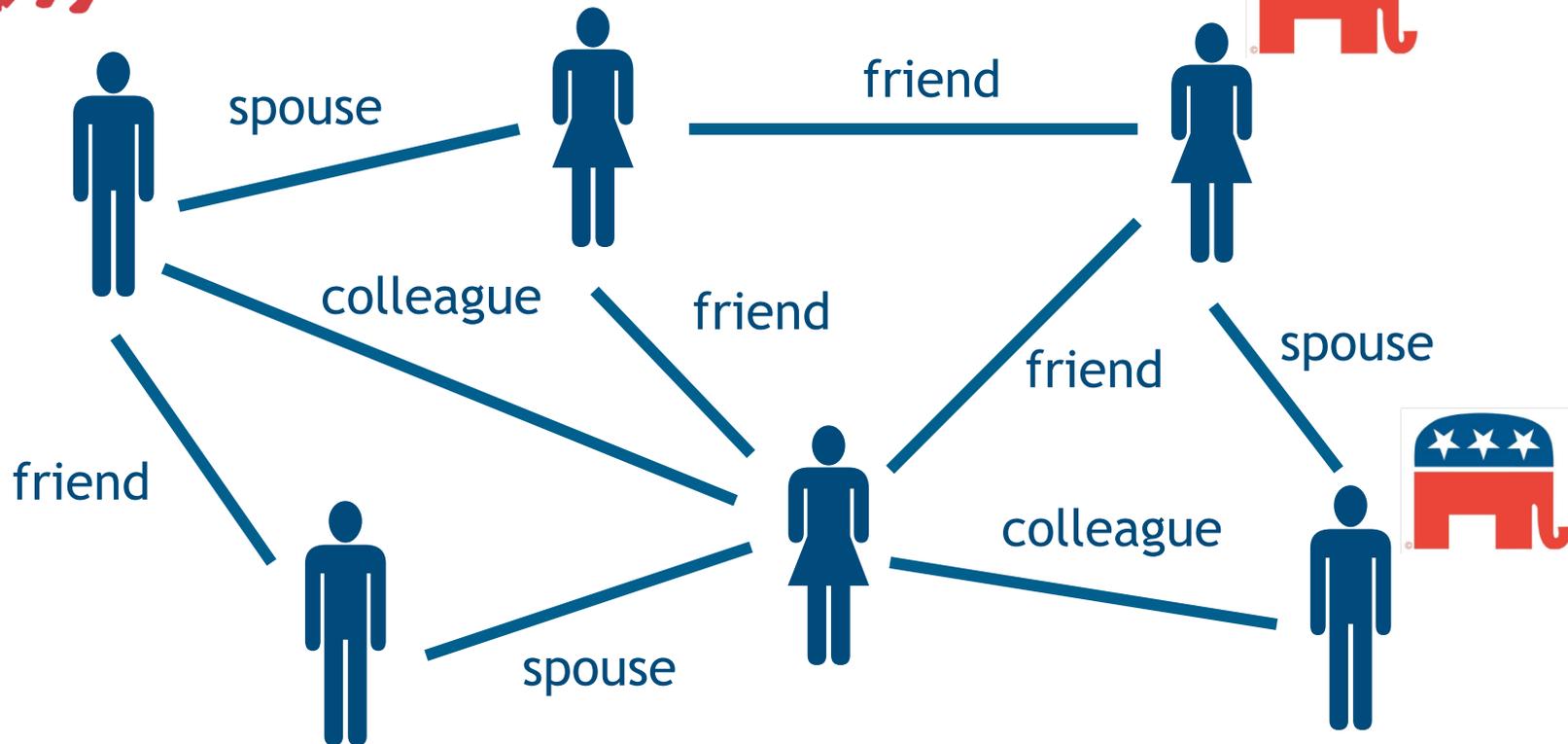
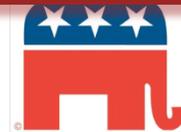
# Collective Classification



# Collective Classification

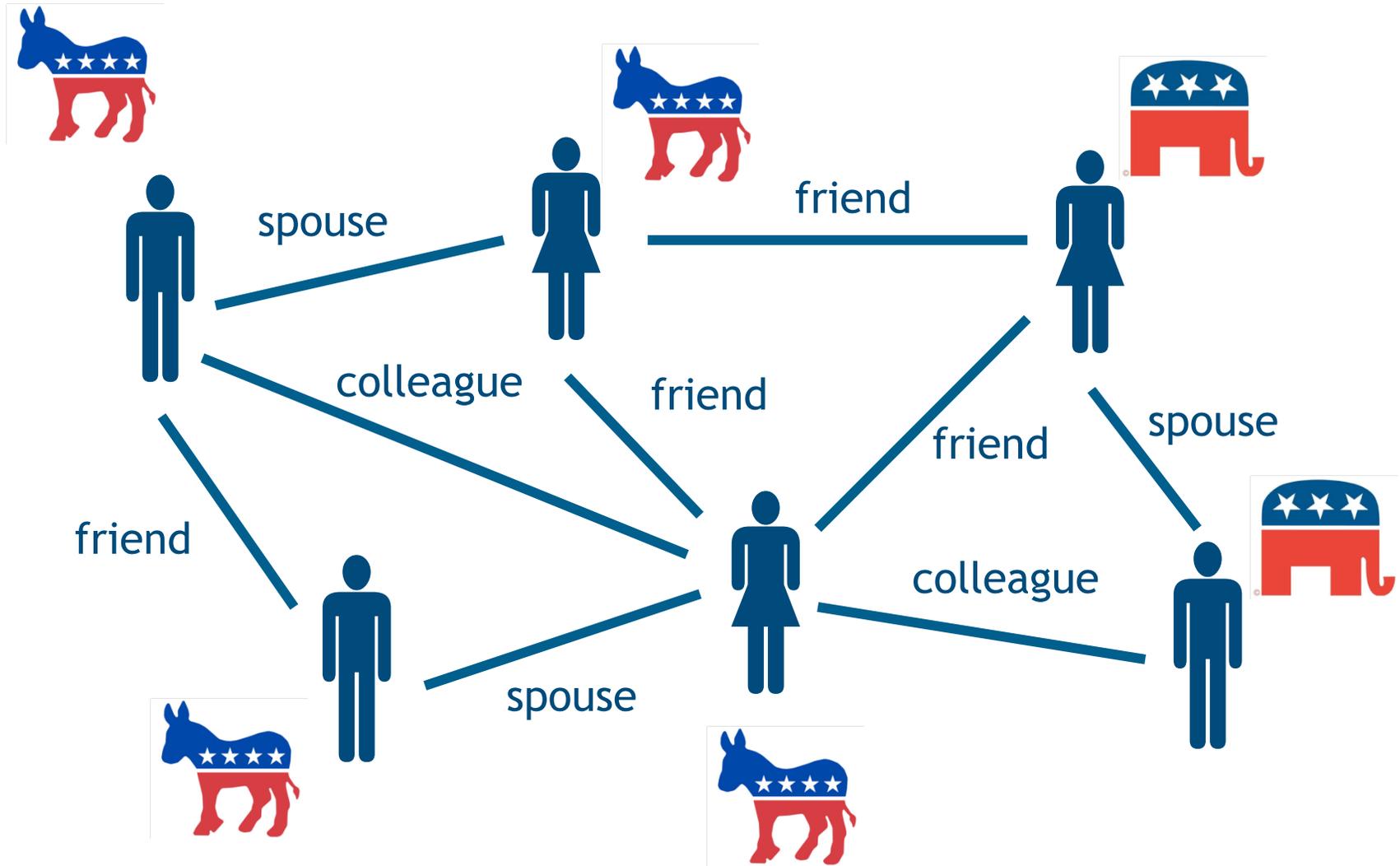


$\text{vote}(A,P) \wedge \text{friend}(B,A) \rightarrow \text{vote}(B,P) : 0.3$



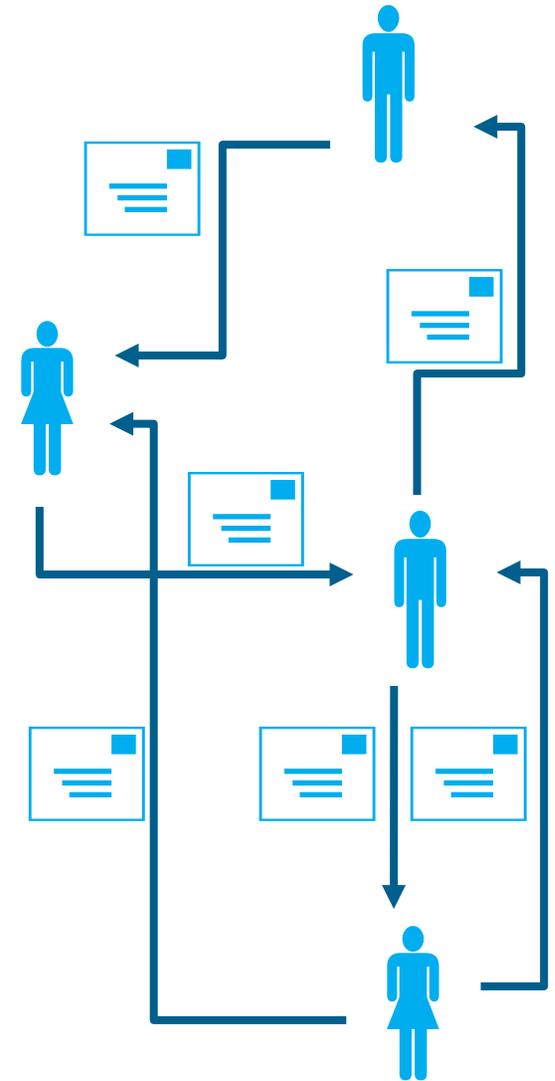
$\text{vote}(A,P) \wedge \text{spouse}(B,A) \rightarrow \text{vote}(B,P) : 0.8$

# Collective Classification



# Link Prediction

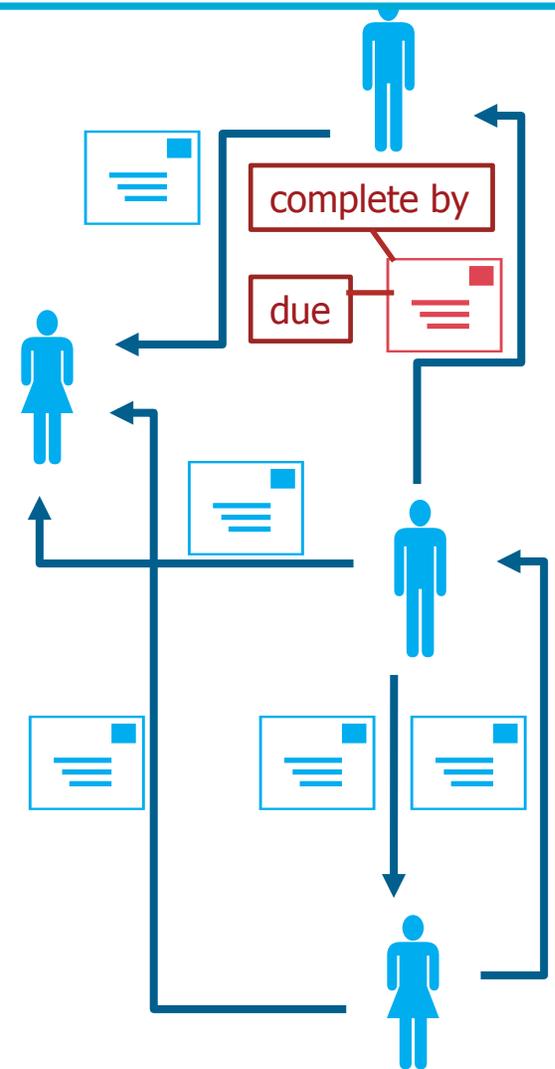
- People, emails, words, communication, relations
- Use model to express dependencies
  - “If email content suggests type X, it is of type X”
  - “If A sends deadline emails to B, then A is the supervisor of B”
  - “If A is the supervisor of B, and A is the supervisor of C, then B and C are colleagues”



# Link Prediction

HasWord(E, "due") => Type(E, deadline) : 0.6

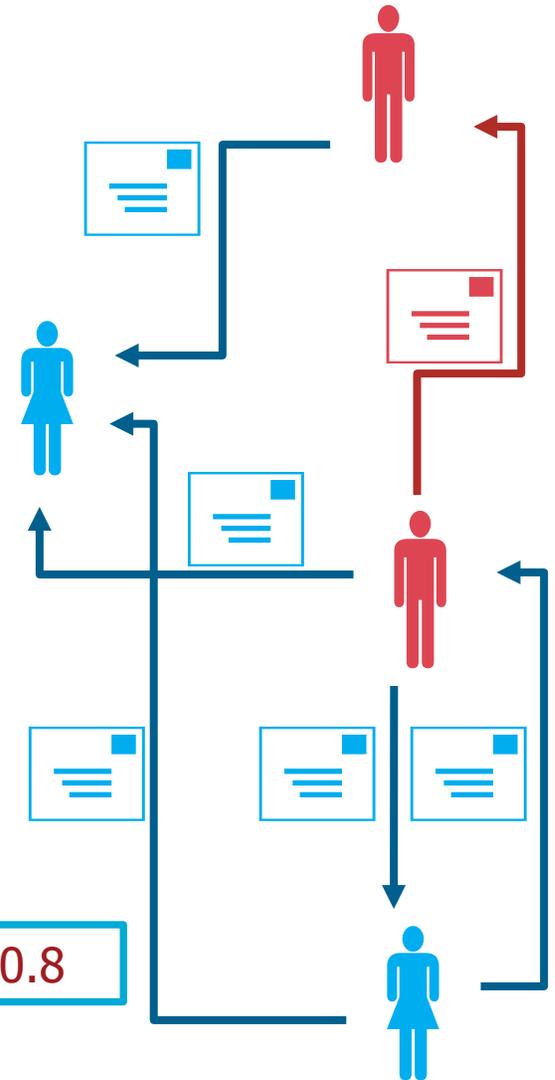
- People, emails, words, communication, relations
- Use model to express dependencies
  - “If email content suggests type X, it is of type X”
  - “If A sends deadline emails to B, then A is the supervisor of B”
  - “If A is the supervisor of B, and A is the supervisor of C, then B and C are colleagues”



# Link Prediction

- People, emails, words, communication, relations
- Use model to express dependencies
  - “If email content suggests type X, it is of type X”
  - “If A sends deadline emails to B, then A is the supervisor of B”
  - “If A is the supervisor of B, and A is the supervisor of C, then B and C are colleagues”

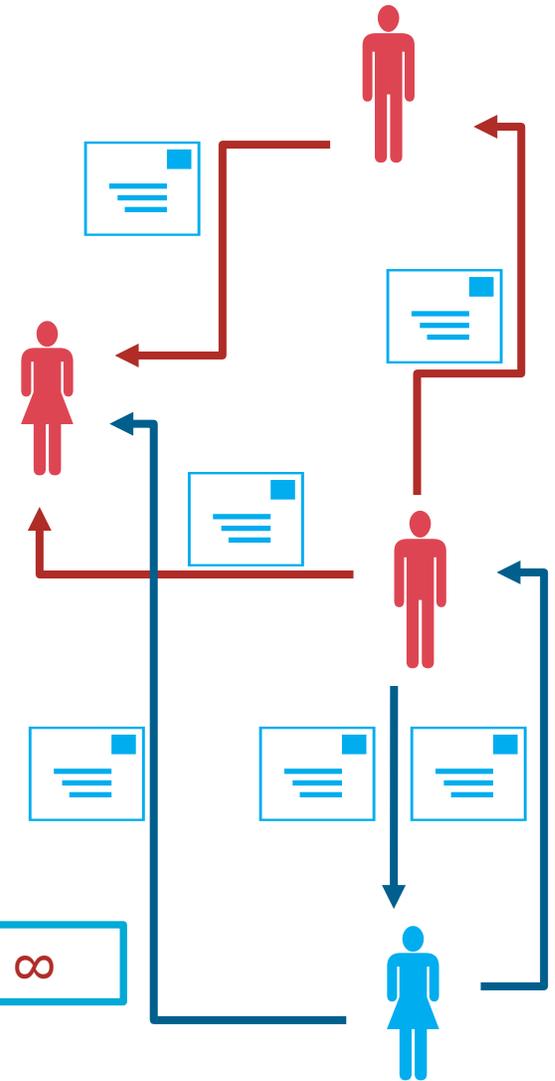
$\text{Sends}(A,B,E) \wedge \text{Type}(E,\text{deadline}) \Rightarrow \text{Supervisor}(A,B) : 0.8$



# Link Prediction

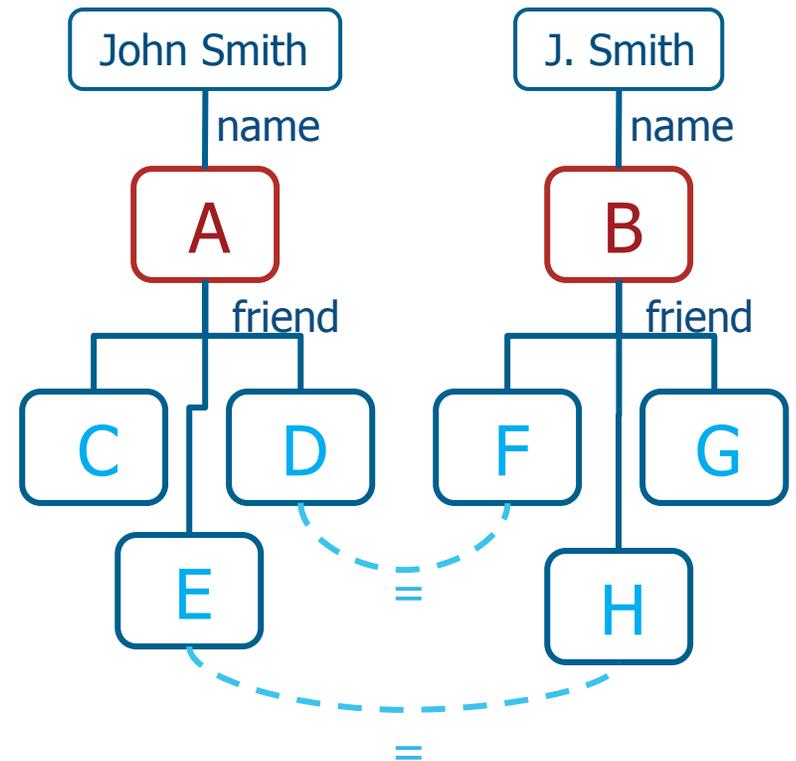
- People, emails, words, communication, relations
- Use model to express dependencies
  - “If email content suggests type X, it is of type X”
  - “If A sends deadline emails to B, then A is the supervisor of B”
  - “If A is the supervisor of B, and A is the supervisor of C, then B and C are colleagues”

$\text{Supervisor}(A,B) \wedge \text{Supervisor}(A,C) \Rightarrow \text{Colleague}(B,C) : \infty$



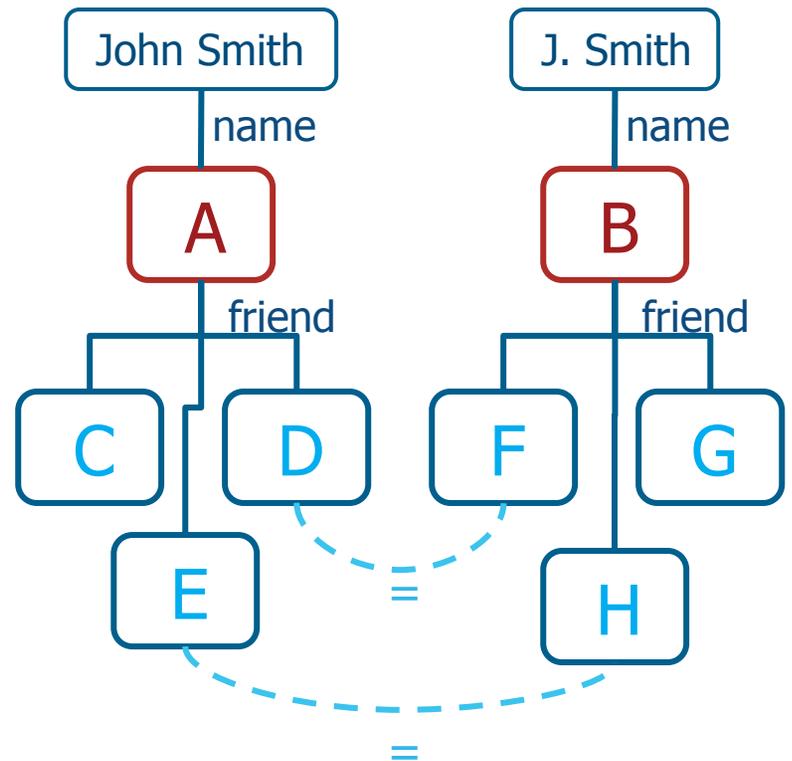
# Entity Resolution

- Entities
  - People References
- Attributes
  - Name
- Relationships
  - Friendship
- Goal: Identify references that denote the same person



# Entity Resolution

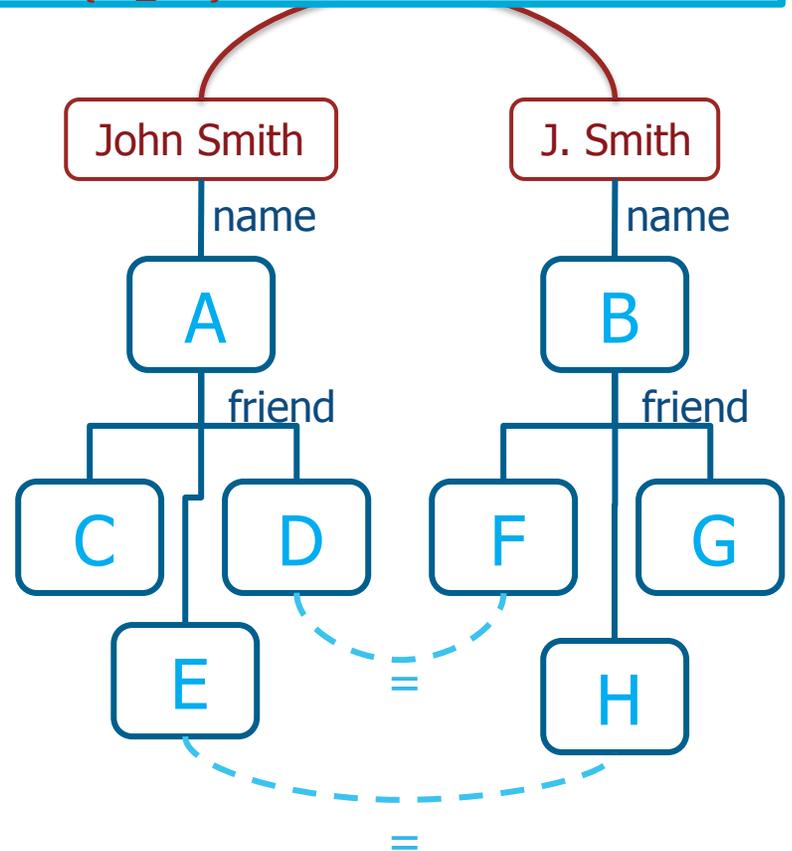
- References, names, friendships
- Use model to express dependencies
  - “ If two people have similar names, they are probably the same’ ’
  - “ If two people have similar friends, they are probably the same’ ’
  - “ If  $A=B$  and  $B=C$ , then  $A$  and  $C$  must also denote the same person’ ’



# Entity Resolution

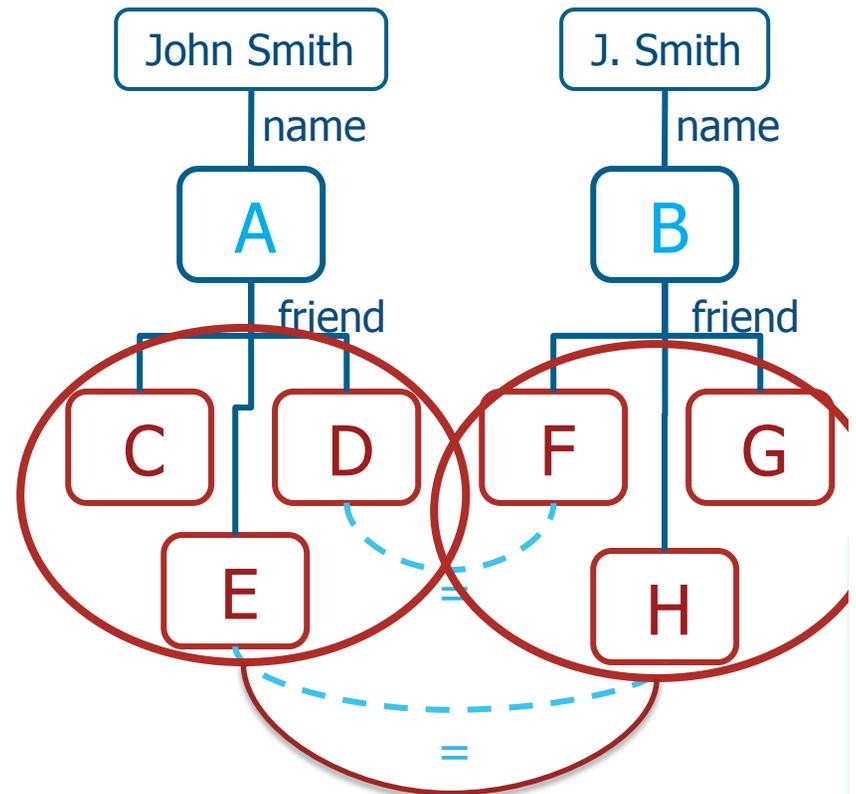
$$A.name \approx_{\{str\_sim\}} B.name \Rightarrow A \approx B : 0.8$$

- References, names, friendships
- Use model to express dependencies
  - “ If two people have similar names, they are probably the same”
  - “ If two people have similar friends, they are probably the same”
  - “ If  $A=B$  and  $B=C$ , then  $A$  and  $C$  must also denote the same person”



# Entity Resolution

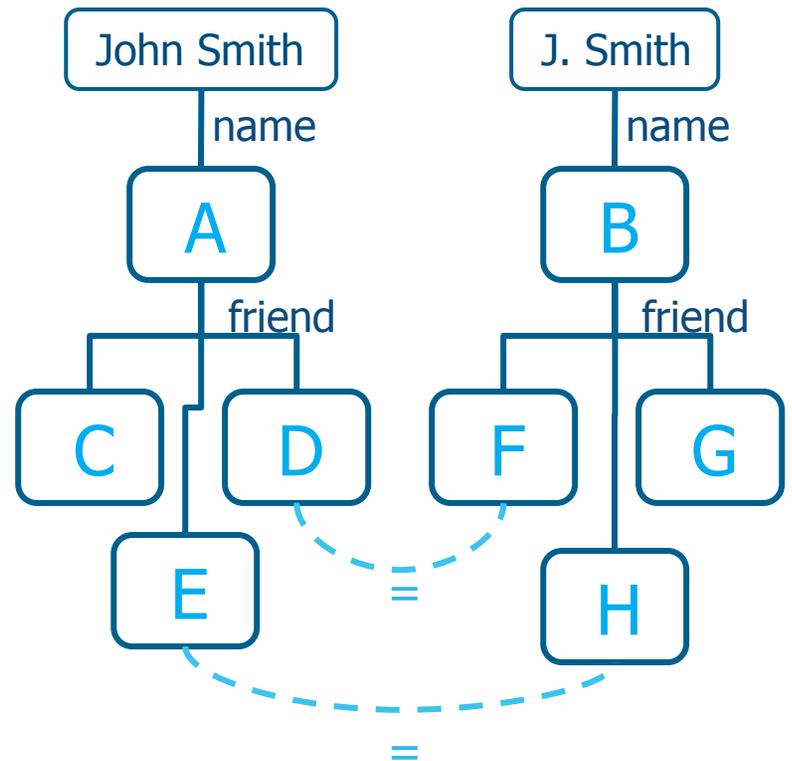
- References, names, friendships
- Use model to express dependencies
  - “If two people have similar names, they are probably the same”
  - “If two people have similar friends, they are probably the same”
  - “If  $A=B$  and  $B=C$ , then  $A$  and  $C$  must also denote the same person”



$$\{A.friends\} \approx_{\{ \}} \{B.friends\} \Rightarrow A \approx B : 0.6$$

# Entity Resolution

- References, names, friendships
- Use model to express dependencies
  - “ If two people have similar names, they are probably the same’ ’
  - “ If two people have similar friends, they are probably the same’ ’
  - “ If  $A=B$  and  $B=C$ , then  $A$  and  $C$  must also denote the same person’ ’



$$A \approx B \wedge B \approx C \Rightarrow A \approx C : \infty$$

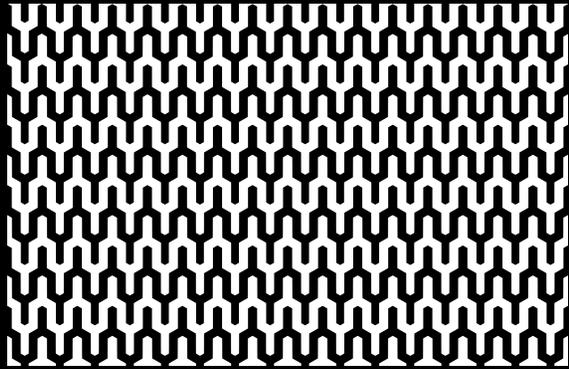
# Challenges

- **Collective Classification:** labeling nodes in graph
  - irregular structure, not a chain, not a grid
  - Challenge: One large partially labeled cluster

- **Link prediction:** predicting missing edges
  - Dependencies between nodes
  - Don't want to consider all possible edges
  - Challenge: extremely skewed probabilities

**Key Idea: Predictions/Outputs depend on each other, joint reasoning is required!**

- **Community detection:** determine nodes that refer to the same entities in a graph
  - Dependencies between clusters
  - Challenge: enforcing constraints, e.g. transitive closure



Patterns



Key Ideas



Tools



# Probabilistic Soft Logic



Stephen Bach



Matthias Broecheler



Alex Memory



Lily Mihalkova



Stanley Kok



Angelika Kimmig



Bert Huang



Ben London



Arti Ramesh



Jay Pujara



Shobeir Fakhraei



Hui Miao

# Probabilistic Soft Logic (PSL)

**Declarative language** based on logics to express collective probabilistic inference problems

- Predicate = relationship or property
- Atom = **(continuous)** random variable
- Rule = capture dependency or **constraint**
- Set = define **aggregates**

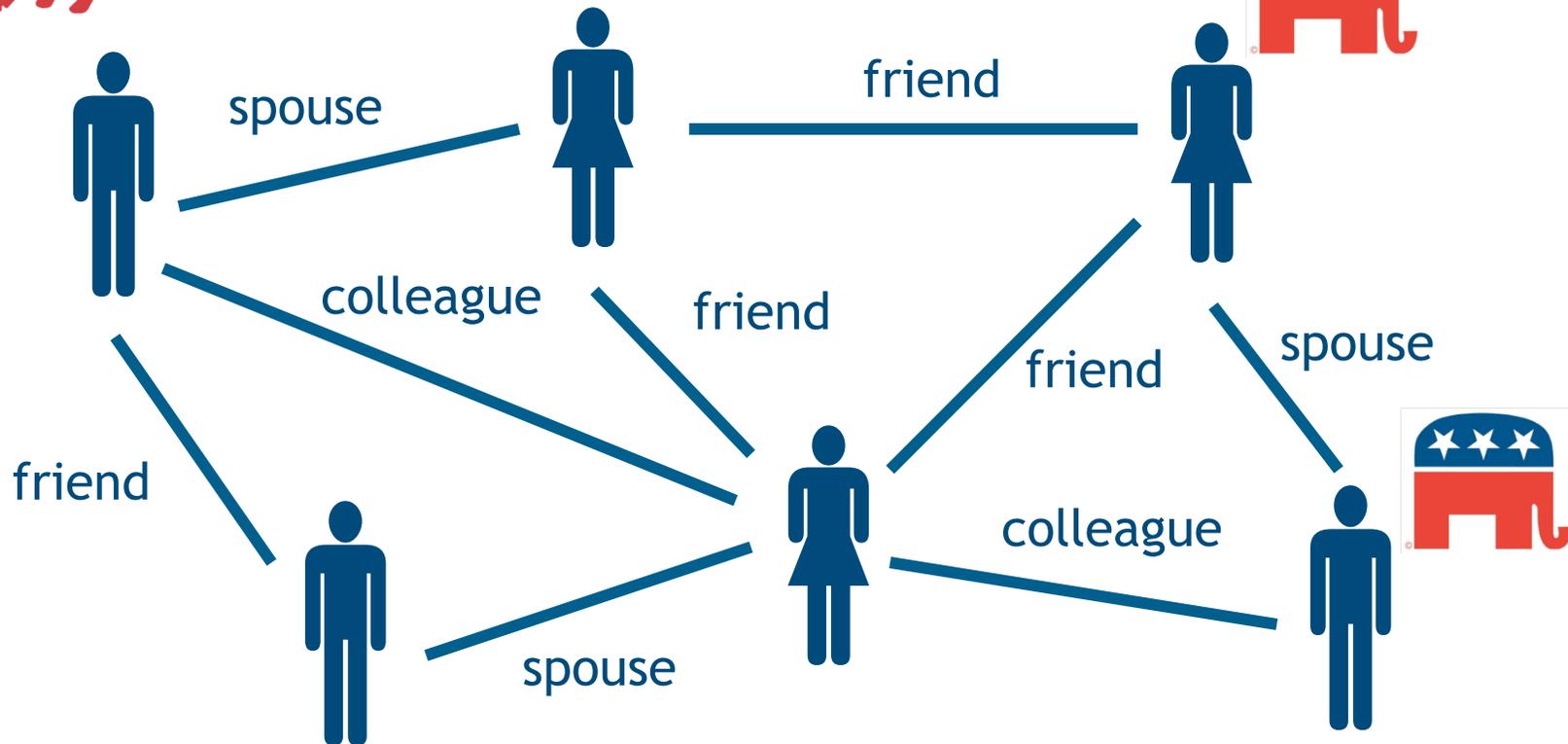
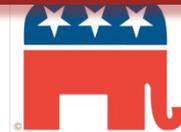
PSL Program = Rules + Input DB

*Reference: Hinge-Loss Markov Random Fields and Probabilistic Soft Logic, Stephen H. Bach, Matthias Broecheler, Bert Huang, Lise Getoor, arXiv 2015*

# Collective Classification



$\text{vote}(A,P) \wedge \text{friend}(B,A) \rightarrow \text{vote}(B,P) : 0.3$



$\text{vote}(A,P) \wedge \text{spouse}(B,A) \rightarrow \text{vote}(B,P) : 0.8$

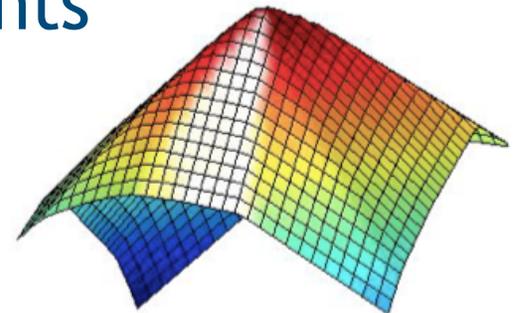
# PSL Foundations

- PSL makes large-scale reasoning scalable by mapping logical rules to convex functions
- Three principles justify this mapping:
  - LP programs for MAX SAT with approximation guarantees [Goemans and Williamson, '94]
  - Pseudomarginal LP relaxations of Boolean Markov random fields [Wainwright, et al., '02]
  - Łukasiewicz logic, a logic for reasoning about continuous values [Klir and Yuan, '95]

# Hinge-loss Markov Random Fields

$$P(\mathbf{Y} | \mathbf{X}) = \frac{1}{Z} \exp \left[ - \sum_{j=1}^m w_j \max\{\ell_j(\mathbf{Y}, \mathbf{X}), 0\}^{p_j} \right]$$

- Continuous variables in  $[0,1]$
- Potentials are hinge-loss functions
- Subject to arbitrary linear constraints
- Log-concave!



# PSL in a Slide

- MAP Inference in PSL translates into convex optimization problem -> **inference is really fast!**
- Inference further enhanced with state-of-the-art optimization and distributed processing paradigms such as ADMM & GraphLab -> **inference even faster!**
- **Outperforms discrete MRFs** in terms of speed, and (very) often accuracy
- **PSL is *flexible***: Applied to image segmentation, activity recognition, stance-detection, sentiment analysis, document classification, drug target prediction, latent social groups and trust, engagement modeling, ontology alignment, and looking for more!



Discussion

Lesson: Make sure you are working on the right graph before performing analytics!

How do you get the right graph?

Use **graph identification** to infer it from the data!

# Closing Comments

- Make sure you're working on the right graph before performing analytics!
- Combining sampling (incomplete data) with inference
  - Active inference and surveying
  - Latent variable learning
- Important flipside to graph identification: ***Privacy***
  - How to ensure that the graph can't be re-identified?
- Research challenges and compelling applications abound!

# Thank You!



[psl.umiacs.umd.edu](http://psl.umiacs.umd.edu)

Contact information:  
[getoor@ucsc.edu](mailto:getoor@ucsc.edu)

