# KDD'14 Tutorial
# Network Mining and Analysis for Social Applications

**Feida Zhu**
Singapore Management University
**Huan Sun, Xifeng Yan**
University of California at Santa Barbara

August 24, 2014

# Our perspective: social applications

**Retweet**

**Rate**

**Adopt**

**Friendship Network**

**Collective behavior**

## User
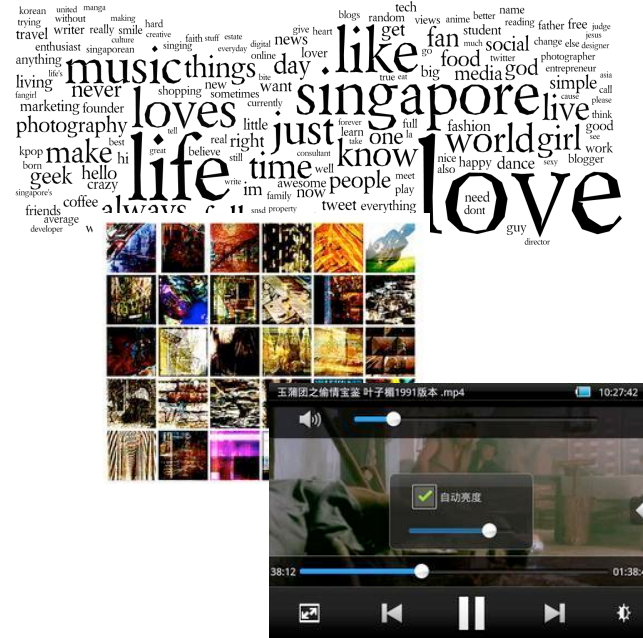
- Follow
- Befriend
- …

## Feedback

- Retweet
- Reply
- Rate
- Adopt
- …

## Content

- Usage of hashtags
- Usage of lexicons
- …

# Our perspective: social applications

- The essence of "Social Data"--- what makes it different from traditional relational table data
  - o Interactive behavior
    - ➢ Behavioral patterns, anomalies, collaborations, relationship, alignment, etc.
  - o Network Structure
    - ➢ Structural properties, measures, patterns, correlation, evolution, etc.

- This tutorial examines network mining and analysis from these two aspects as motivated by real social applications.

# Tutorial Coverage

- Network Structure
  - o Part I: Network Correlation and Patterns (Huan)
  - o Part II: Frequent Network Patterns (Feida)

- Interactive Behavior
  - o Part III: Collaboration Patterns (Huan)

- Network Structure + Interactive Behavior
  - o Part IV: Relationship Mining (Feida)
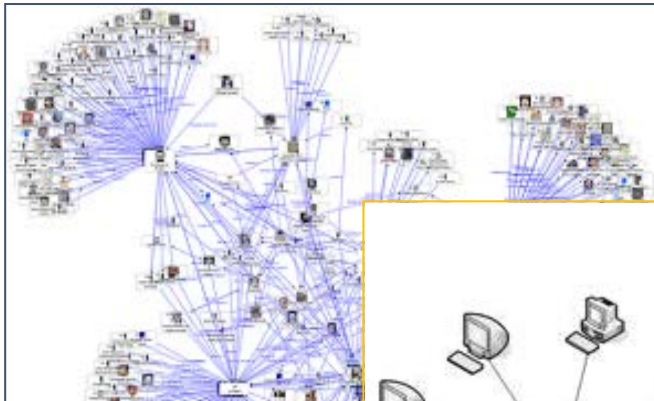  - o Part V: User Identity Linkage (Feida)

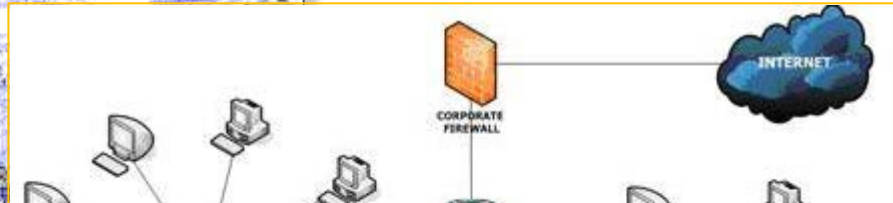# Network Mining and Analysis
# for Social Applications

# Part I
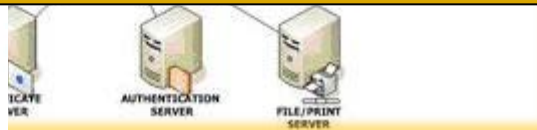## Network Correlation and Patterns

# Networks with Rich Attributes



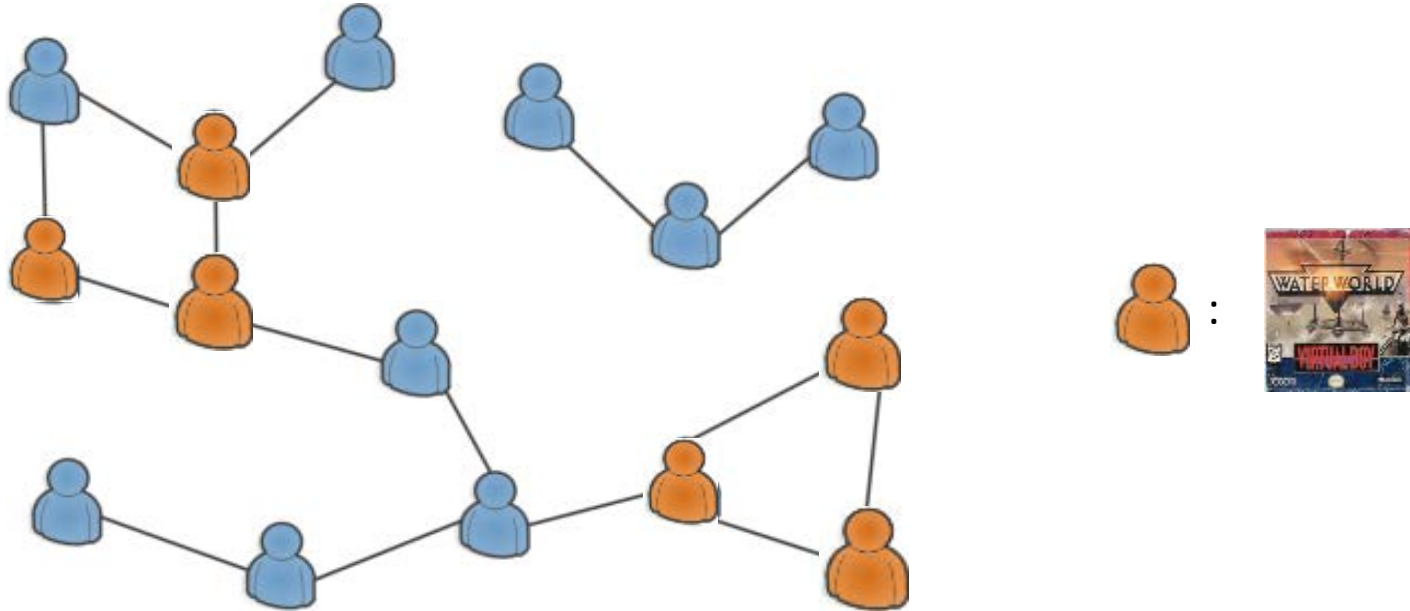Social networks

Research Collaboration Networks

## New Perspectives of Patterns and Correlations

# Example of Correlations



Correlation between the occurrence of an event and the network structure

Further study:
1. Is this correlation due to influence in the network? (network structure→ event occurrence)

2. Conversely, is it the preferences over the same product facilitate the formation of the link between users, i.e., their friendship development? (event occurrence → network structure)

# Pattern Kaleidoscope

❑ Proximity Pattern

❑ Attribute-Structure Correlations

❑ Cohesive Pattern

❑ Itemset-sharing Pattern

❑ Graph Topological pattern

❑ Graph Iceberg

❑ Graph Anomaly

❑ Frequent Network Pattern
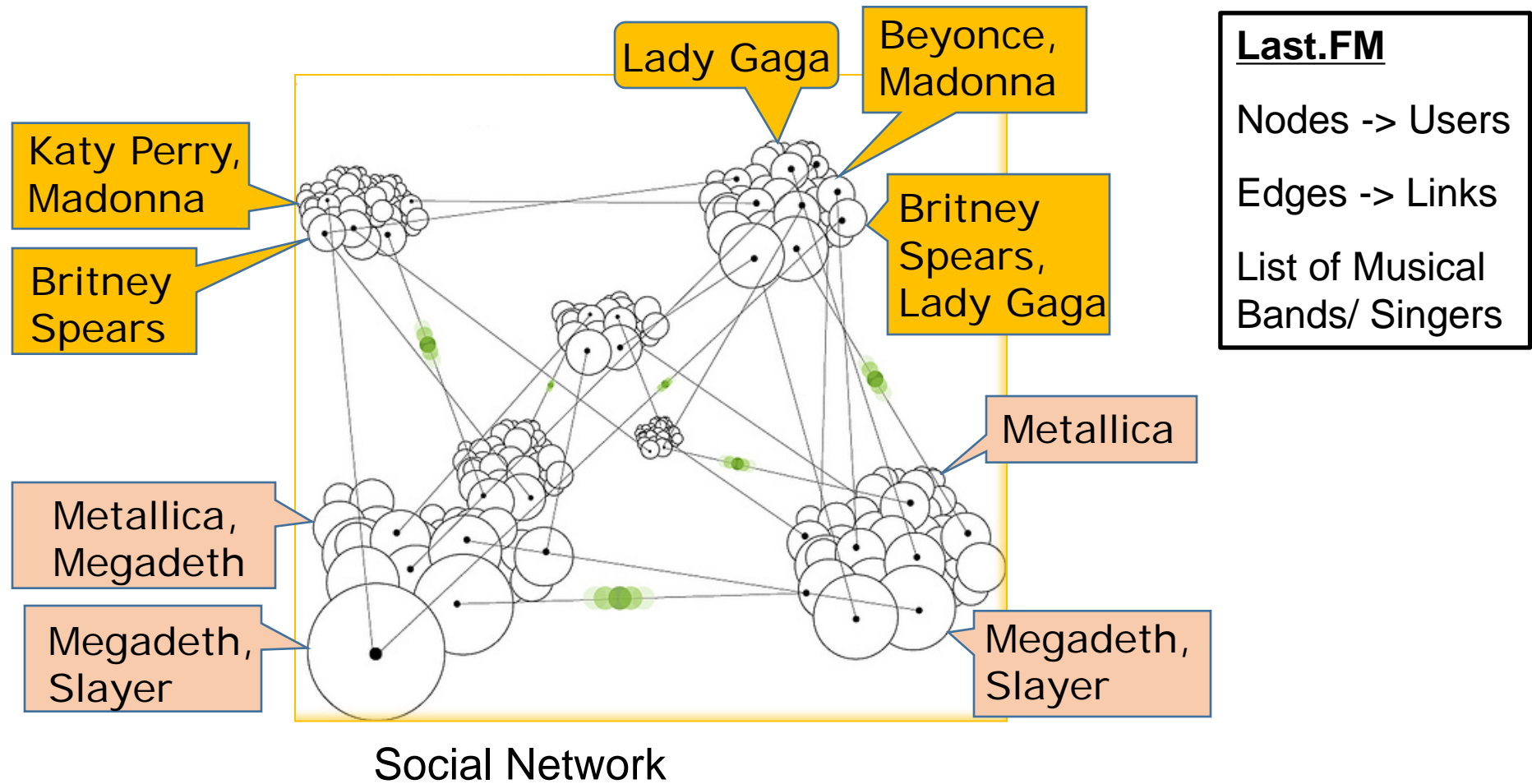
# Pattern Kaleidoscope

❑ Proximity Pattern ✔

❑ Attribute-Structure Correlations ✔

❑ Cohesive Pattern

❑ Itemset-sharing Pattern **Brief introduction**

❑ Graph Topological Pattern

❑ Graph Iceberg ✔

❑ Graph Anomaly ➡ **Akoglu et al., Tutorial at WSDM'13**

❑ Frequent Network Pattern ✔ (Part II)

# What are the **related Musical Bands/ Singers** that **co-occur frequently in neighborhood?**



Social Network

# What are **Related Computer Attacks** that **Co-occur Frequently in Neighborhood?**



Intrusion Network

**Computers in LAN**

Computers in Same LAN Attacked by Similar Intrusions

TFTP_Put, Ping_Flood
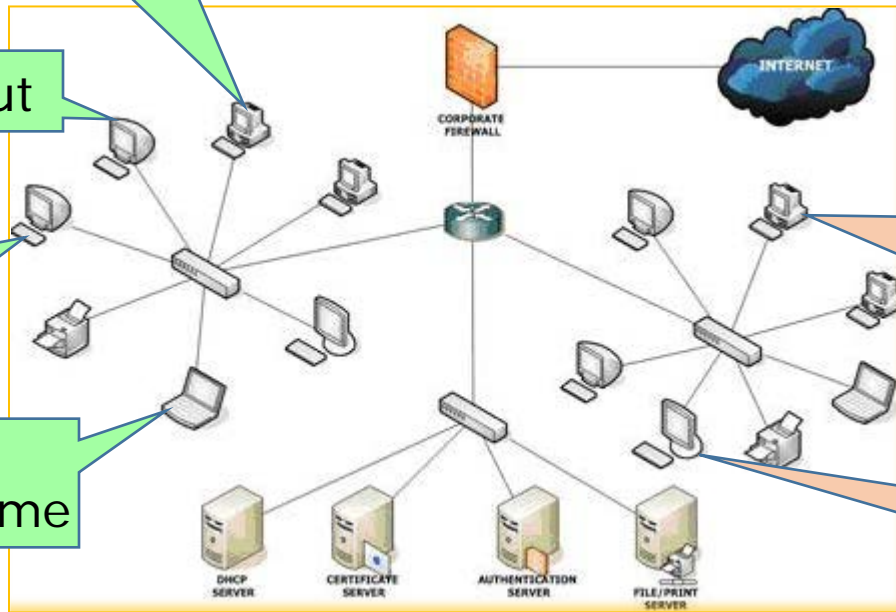
TFTP_Put

TFTP_Put, ICMP_Flood

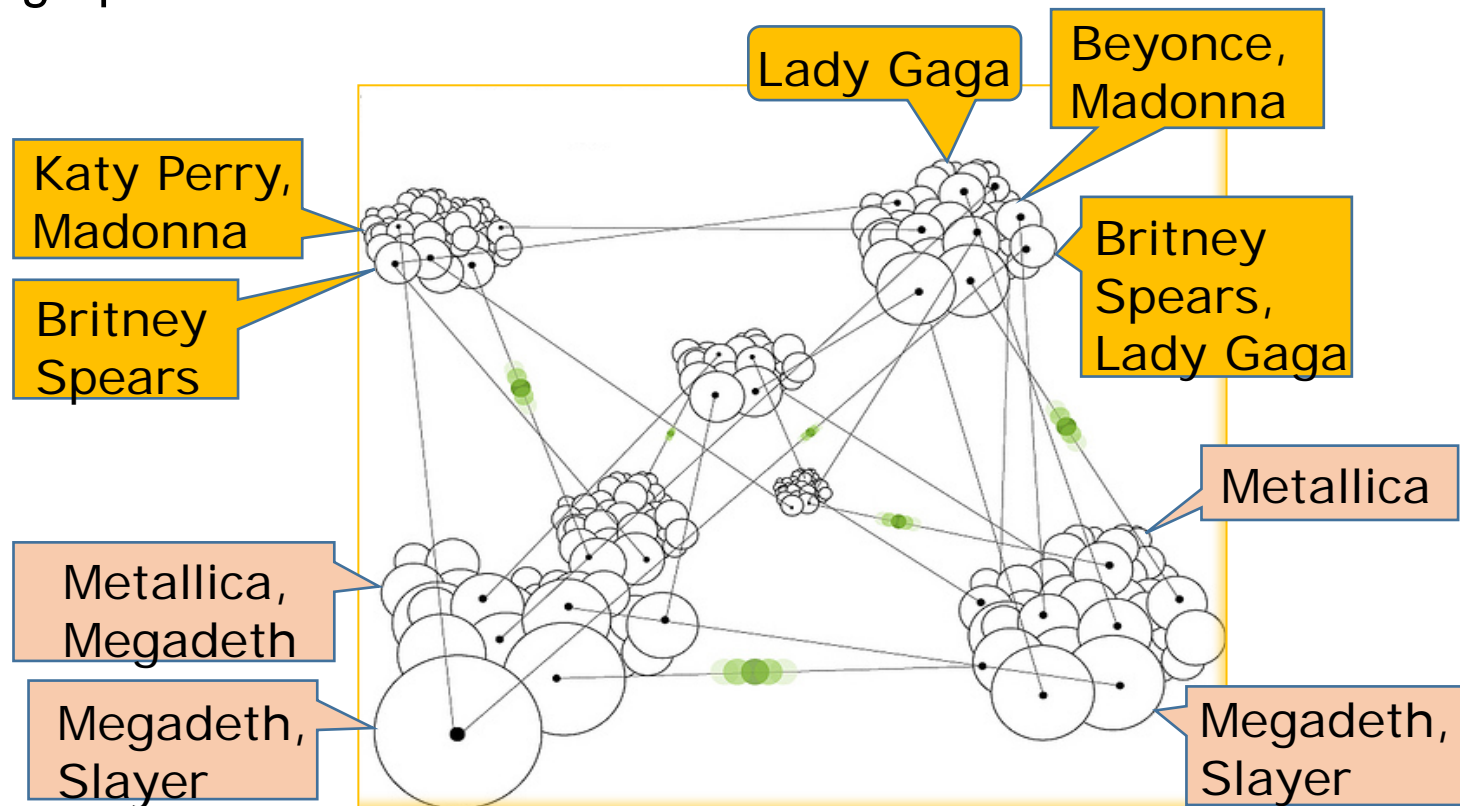Audit_TFTP _Get_Filename

SQL _SSRP _Slammer _Worm

SQL _SSRP _ StackBo
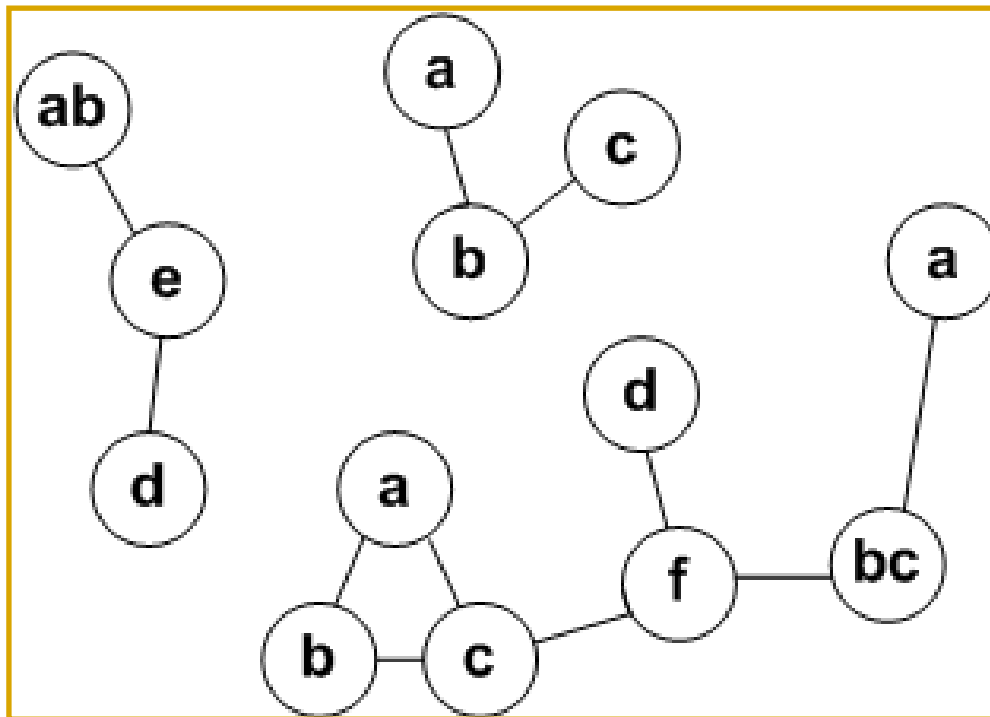
# Proximity Pattern Mining

## ❑ Definition

A subset of labels that repeatedly appear in tightly connected subgraphs in G.

# Proximity Pattern Mining

## ❑ **Definition**

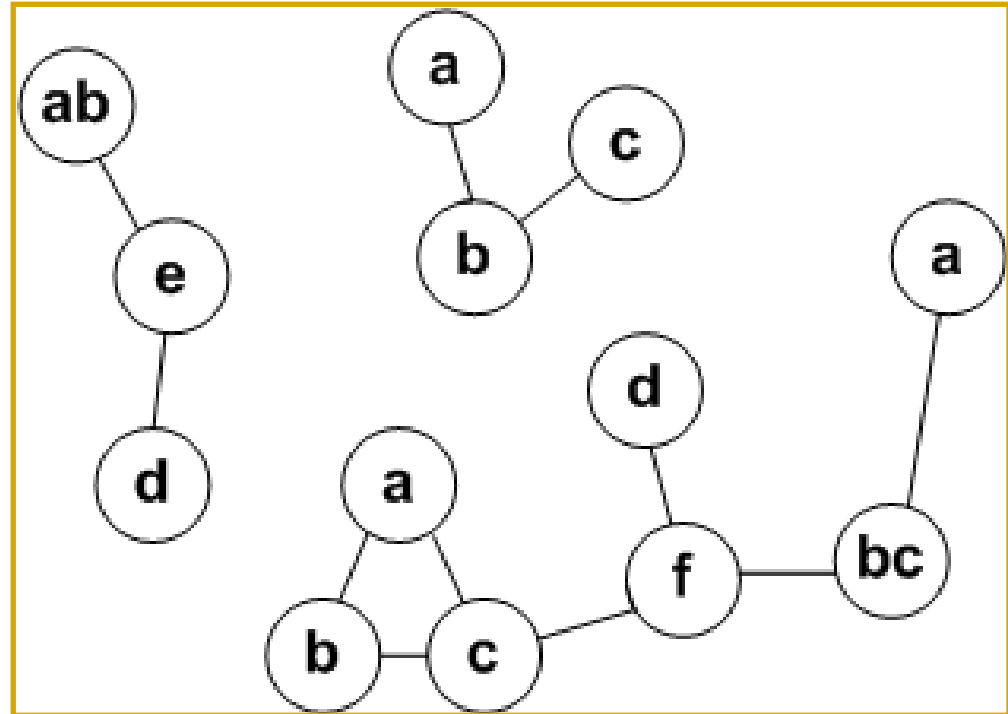A subset of labels that repeatedly appear in tightly connected subgraphs in G.



a, b – **YES**
a, b, c – **YES**
d, e, f - **NO**

# Proximity Pattern Definition

❑ **Characteristics**
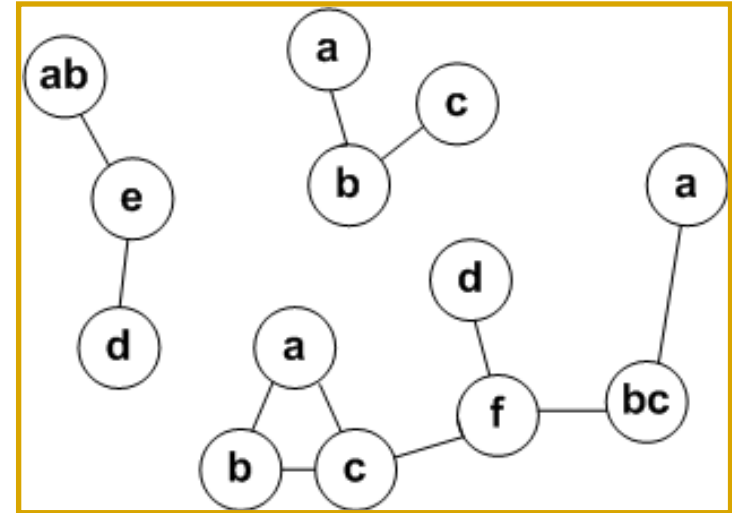
(1) Proximity

(2) Frequency

(3) Flexibility



**a, b – YES; a, b, c – YES**
**d, e, f - NO**

# Proximity Pattern Definition
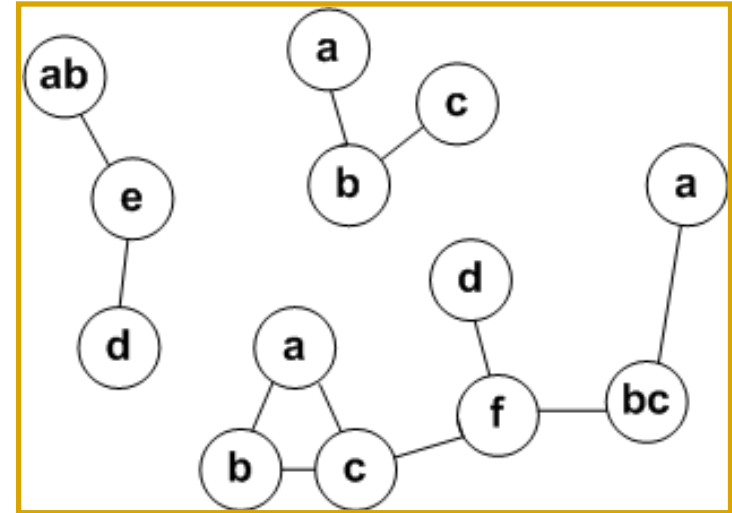
❑ Will Frequent Subgraph Mining Work?



Proximity pattern: {a, b, c}

# Proximity Pattern Definition

❑ Will Frequent Subgraph Mining Work?
- **NO !!!**
- **Lack of Flexibility**



Proximity pattern: {a, b, c}

# Proximity Pattern Definition

❑ Will Frequent Subgraph Mining Work?

- **NO !!!**
- **Lack of Flexibility**

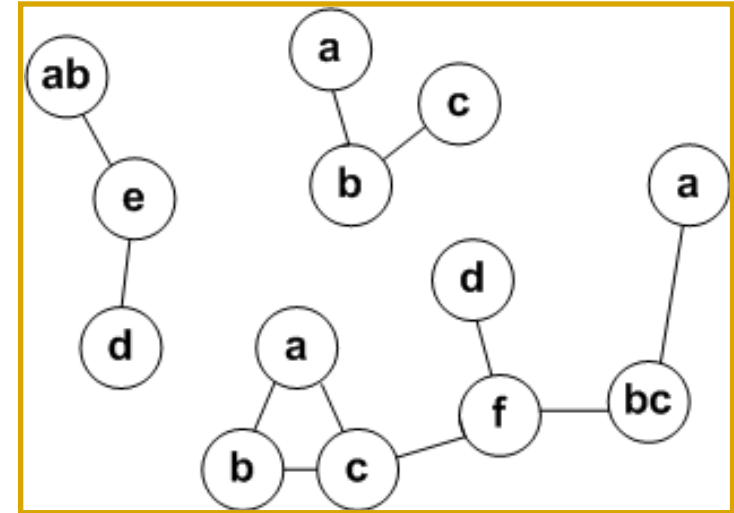❑ Will Frequent Itemset Mining Work?



Proximity pattern: {a, b, c}

# Proximity Pattern Definition

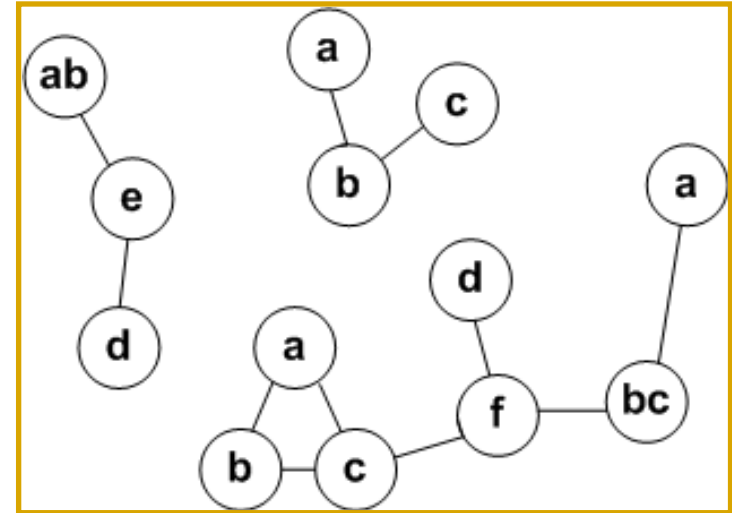❑ Will Frequent Subgraph Mining Work?
 - **NO !!!**
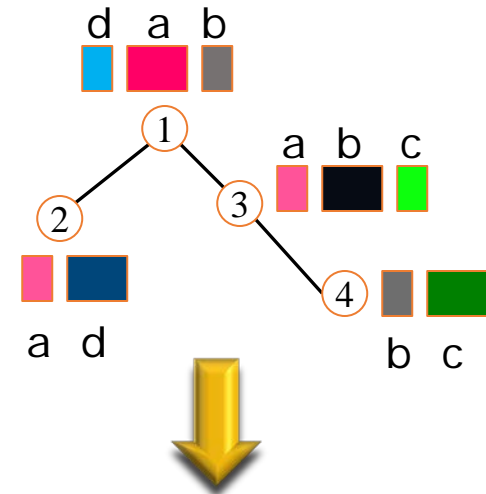 - **Lack of Flexibility**


❑ Will Frequent Itemset Mining Work?
 - **NO !!!**
 - **No Notion of Edge**



Proximity pattern: {a, b, c}

# Information Propagation Model

# Information Propagation Model



Information Propagation

- ☐ **Frequent-Pattern (FP) Tree** cannot handle fractional association values because of the new definition of Support.

- ☐ Modify FP Tree Structure and Algorithm.

- ☐ *C. C. Aggarwal et. al (KDD '09), Bernecker et. al (KDD '09).*

| | a | b | c | d |
|---|---|---|---|---|
| 1 | 1.00 | 0.12 | 0.00 | 0.12 |
| 2 | 0.19 | 0.00 | 0.00 | 1.00 |
| 3 | 0.12 | 1.00 | 0.12 | 0.00 |
| 4 | 0.00 | 0.19 | 1.00 | 0.00 |

**Frequent Itemset Mining (Probabilistic)**

# Top-k Interesting Patterns

❑ How to measure "Interesting-ness"?

 **----Randomization Test**

Generate graph *Q* from graph *G* by randomly swapping the labels among nodes. Let, *p* and *q* be the support values of itemset *I* in *G* and *Q* respectively. High difference indicates interestingness.

G-test Score:

$$p \cdot \ln \frac{p}{q} + (1 - p) \cdot \ln \frac{1 - p}{1 - q}$$

**----Proximity Patterns minus Frequent Itemset Patterns**

# Top-k Interesting Patterns

❑Last.FM

Proximity
Patterns

| # | Proximity Patterns | Score |
|---|---|---|
| 1 | Tiësto, Armin van Buuren , ATB | 0.62 |
| 2 | Katy Perry, Lady Gaga, Britney Spears | 0.58 |
| 3 | Ferry Corsten, Tiësto, Paul van Dyk | 0.55 |
| 4 | Neaera, Caliban, Cannibal Corpse | 0.52 |
| 5 | Lacuna Coil, Nightwish, Within Temptation | 0.47 |

❑ Britney Spears, Lady Gaga, Katy Perry – **American Female Pop Singers**
❑ ATB, Paul van Dyk – **German DJ**
❑ Tiesto, Ferry Corsten, Armin van Buuren – **Dutch DJ**
❑ Neaera, Caliban, Cannibal Corpse – **Death Metal Bands**

❑ Lucuna Coil, Nightwish, Within Temptation – **Gothic Metal Bands**

# Top-k Interesting Patterns

❏Intrusion

Proximity
Patterns

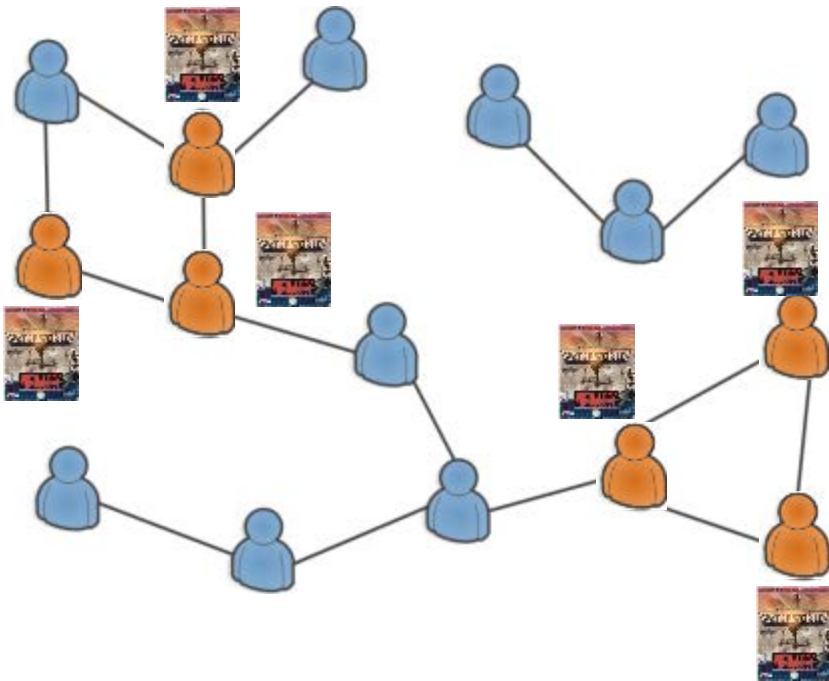| # | Interesting Patterns | Score |
|---|---|---|
| 1 | Ping_Sweep, Smurf_Attack | 2.42 |
| 2 | TFTP_Put, Audit_TFTP_Get_Filename, ICMP_Flood, Ping_Flood | 2.32 |
| 3 | TCP_Service_Sweep, Email_Error | 1.21 |
| 4 | HTML_Outlook_MailTo_Code_Execution, HTML_NullChar_Evasion | 1.15 |
| 5 | SQL_SSRP_Slammer_Worm, SQL_SSRP_StackBo | 0.88 |

Proximity Patterns
Minus
Frequent Itemsets

| # | Interesting Patterns | Score |
|---|---|---|
| 1 | ICMP_Flood, Ping_Flood | 0.94 |
| 2 | Email_Error, SMTP_Relay _Not_Allowed, HTML_Null Char_Evasion | 0.94 |
| 3 | Image_RIFF_Malformed, HTML_NullChar_Evasion | 0.90 |
| 4 | TFTP_Put, Ping_Flood, Audit_TFTP_Get_Filename | 0.80 |
| 5 | Email_Command_Overflow, Email_Virus_Double_Extension, Email_Error | 0.75 |

# Structural Correlational Pattern

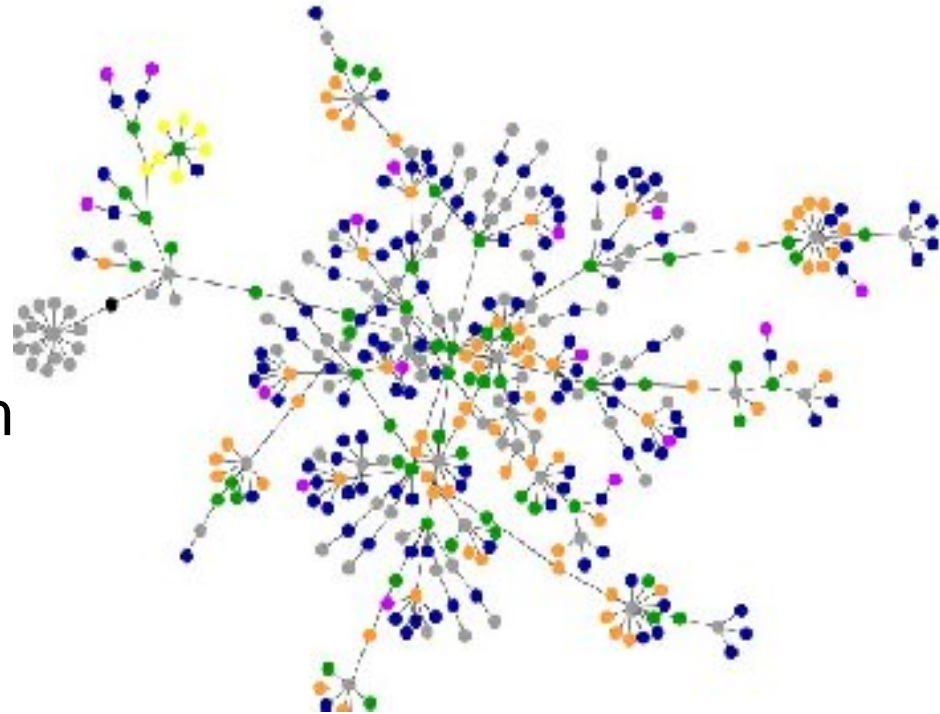Which product's sales is more correlated with the social network structure?



Waterworld game

Hp printer

# A General Situation

❑ Events taking place on nodes of a social graph
- Online shopping
- Blogging
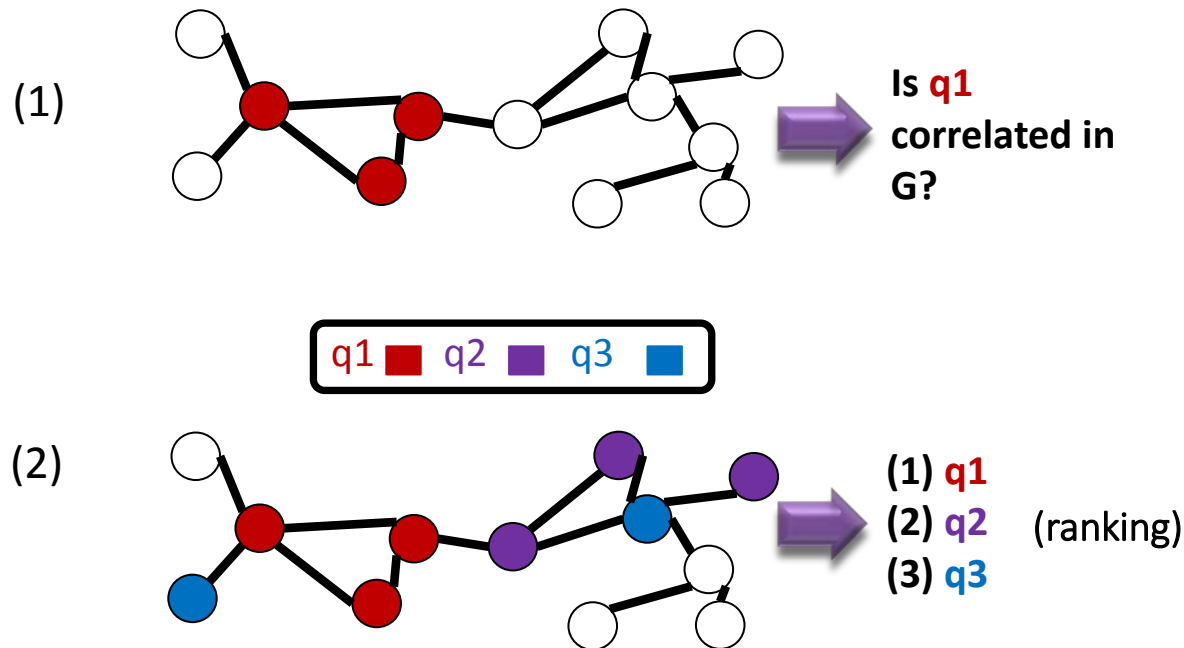- Virus infection

❑ Social influence vs. Random occurrence

# Why Measuring Such Correlations?

❑ Help understand the distribution of events in networks

❑ Help detect viral influence in the underlying network (e.g. product sales)

  o Correlation has to do with link type, event type and time

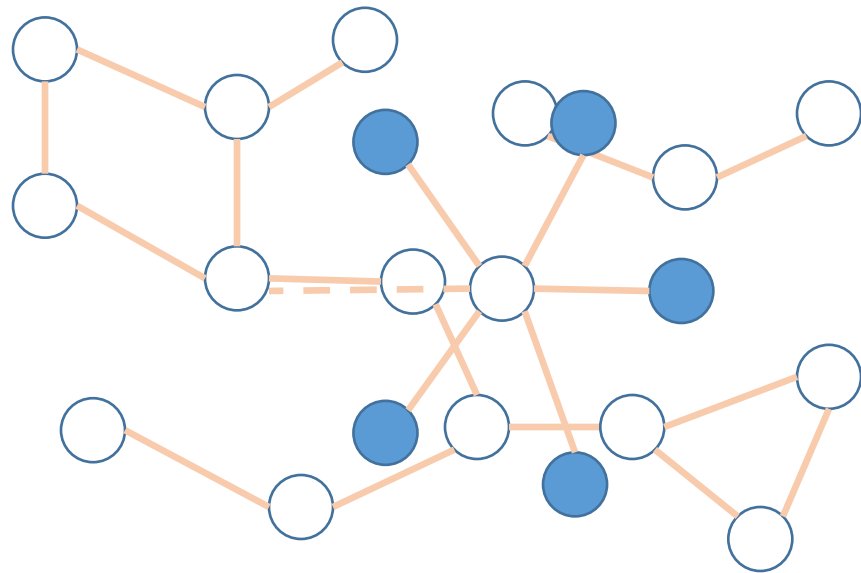  o Facilitate product promotion, online ads recommendation

# Problem Formulation

❑ A graph $G = (V, E)$ and an event set $Q = \{q_i\}$

❑ $V_q$--the set of nodes having event $q$. Let $|V_q| = m$, $|V| = n$

# How to Characterize Correlation?

❑ If correlated, blue nodes tend to stick together.

❑ A naïve approach: only look at neighborhood

❑ General idea: compute the **aggregated proximity** among blue nodes

# Measure Definition

❑ The measure is defined as

$$\rho(V_q) = \frac{\sum_{v \in V_q} s(v, V_q \setminus \{v\})}{\left| V_q \right|}$$

○ $V_q$: the set of nodes having event $q$; $s(\cdot)$ can be any graph proximity measure
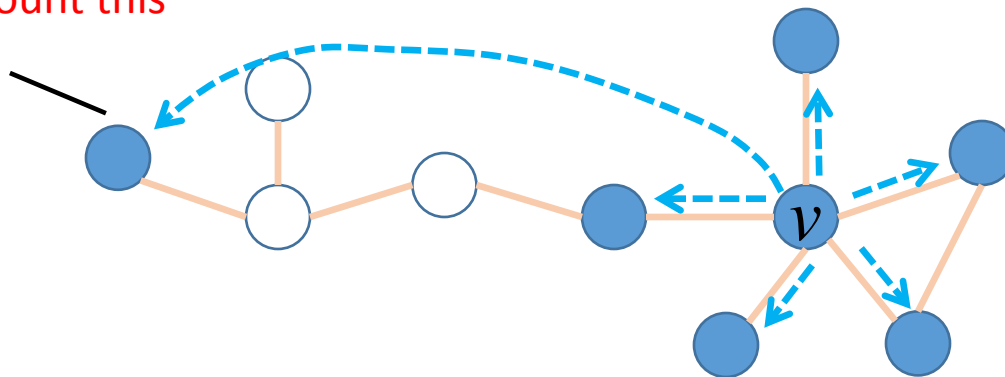
❑ We choose **Hitting Time.**

# Measure Definition

❑ **Hitting time**: expected number of steps to reach a target node via random walk:

$$h(v_i, B) = \sum_{t=1}^{\infty} t \Pr(T_B = t \mid x_0 = v_i)$$

o $B$: target node set; $\Pr(T_B=t|x_0=v_i)$: the probability that we start from $v_i$ and reach $B$ after $t$ steps

Hitting time will not count this node

# Hitting time & Decayed Hitting Time

❑ **Hitting time**: expected number of steps to reach a target node via random walk:

$$h(v_i, B) = \sum_{t=1}^{\infty} t \Pr(T_B = t \mid x_0 = v_i)$$

o $B$: target node set; $\Pr(T_B{=}t|x_0{=}v_i)$: the probability that we start from $v_i$ and reach $B$ after $t$ steps

# Hitting time & Decayed Hitting Time

❑ **Hitting time**: expected number of steps to reach a target node via random walk:

$$h(v_i, B) = \sum_{t=1}^{\infty} t \Pr(T_B = t \mid x_0 = v_i)$$

  o $B$: target node set; $\Pr(T_B{=}t|x_0{=}v_i)$: the probability that we start from $v_i$ and reach $B$ after $t$ steps

❑ **Decayed Hitting Time (DHT)**:

$$\tilde{h}(v_i, B) = \sum_{t=1}^{\infty} e^{-(t-1)} \Pr(T_B = t \mid x_0 = v_i)$$

  o Mapping $[1,\infty)$ to $[0,1]$, high value means high proximity
  o Emphasizing the importance of local neighborhood and reducing the impact of long paths

❑  $\rho(V_q) = \dfrac{\sum_{v \in V_q} s(v, V_q \setminus \{v\})}{|V_q|}$, with $s(v, V_q \setminus \{v\})$ instantiated as $\tilde{h}(v, V_q \setminus \{v\})$

# Measure Computation

❑ Directly compute $\rho(V_q)$ may be time consuming since $\left| V_q \right|$ may be large.

$$\rho(V_q) = \frac{\sum_{v \in V_q} s(v, V_q \setminus \{v\})}{\left| V_q \right|}$$

❑ Sampling: randomly select $c$ nodes from $V_q$ to estimate their DHTs to the remaining nodes.

# Top-5 Structural Correlated Products (TaoBao)

## Laptops and tablets

| # | Product | Bound for $\tilde{\rho}$ | m |
|---|---------|--------------------------|---|
| 1 | ThinkPad T400 | [554.43, 554,47] | 47 |
| 2 | Apple iPad | [227.56, 227.57] | 698 |
| 3 | ThinkPad X200 | [91.39, 91.42] | 60 |
| 4 | Toshiba L600 | [20.36, 20.41] | 31 |
| 5 | ThinkPad T410 | [−1.13, −1.09] | 72 |

## Mobile and handheld devices

| # | Product | Bound for $\tilde{\rho}$ | m |
|---|---------|--------------------------|---|
| 1 | iPod Touch 3 | [92.06, 92.09] | 484 |
| 2 | Nokia 6300 | [90.97, 90.99] | 188 |
| 3 | iPhone 4 | [69.07, 69.09] | 520 |
| 4 | Nokia N82 | [53.20, 53.24] | 84 |
| 5 | HTC G3 | [36.48, 36.49] | 732 |

## Other

| # | Product | Bound for $\tilde{\rho}$ | m |
|---|---------|--------------------------|---|
| 1 | Mamy Poko baby diapers | [238.50, 238.51] | 4892 |
| 2 | Beingmate Infant milk powder | [227.71, 227.72] | 163 |
| 3 | EVE game cards | [198.56, 198.58] | 374 |
| 4 | Mabinogi game cards | [189.56, 189.58] | 446 |
| 5 | Gerber cookies | [149.51, 149.52] | 1491 |

## Other (top 5 uncorrelated)

| # | Product | Bound for $\tilde{\rho}$ | m |
|---|---------|--------------------------|---|
| 1 | Tiffany rings | [2.71, 2.72] | 1092 |
| 2 | Jack&Jones suits | [−0.48, −0.46] | 311 |
| 3 | Ray-Ban sunglasses | [−0.78,−0.77] | 4958 |
| 4 | Swarovski anklets | [−0.88,−0.84] | 72 |
| 5 | Jack&Jones shirts | [−3.28,−3.27] | 1606 |

UCSB

SMU SINGAPORE MANAGEMENT UNIVERSITY

# Top-12 Structural Correlated keywords (DBLP co-author network)

| # | Keyword | Bound for $\tilde{\rho}$ | m |
|---|---------|--------------------------|---|
| 1 | Hadoop | [1225.22, 1225.46] | 57 |
| 2 | Microarray | [958.64, 958.67] | 4738 |
| 3 | OLTP | [912.52, 912.68] | 105 |
| 4 | AJAX | [857.61, 857.74] | 179 |
| 5 | Virus | [825.14, 825.21] | 905 |
| 6 | E-Learning | [811.85, 811.91] | 3552 |
| 7 | Database | [775.80, 775.83] | 19522 |
| 8 | Mining | [760.33, 760.36] | 15371 |
| 9 | Relational | [697.17, 697.22] | 6225 |
| 10 | Retrieval | [647.60, 647.64] | 13996 |
| 11 | Indexing | [624.94, 625.00] | 5069 |
| 12 | Computation | [586.58, 586.63] | 11187 |

# Attribute-Structure Correlated Patterns

❑ Examples
  o Densely connected webpages that share content
  o Groups of friends with common interests
  o Genes that interact and are expressed on the same issues.

# Attribute-Structure Correlated Patterns

❑ Definition
  o A dense subgraph induced by a particular attribute set.

❑Characteristics
  o High correlation between a given attribute set and the occurrence of dense subgraphs

$$\epsilon(S) = \frac{|\mathcal{K}_S|}{|\mathcal{V}(S)|}$$

  *S*: attribute set;
  *Ks*: the set of vertices having *S* in dense subgraphs
  *V*(*s*): the vectices with attribute set *S*

  o High-support attribute sets do not necessarily present high structural correlation.

# **Attribute-Structure Correlated Patterns**

❑ DBLP
   o 108,030 vertices
   o 276,658 edges
   o 23,285 attributes



Pattern induced by
{system , performance}

# Attribute-Structure Correlated Patterns

❑ Last.Fm
- o 272,412 vertices
- o 350,239 edges
- o 3,929,101 attributes



Pattern induced by
{Van Morrison}

# Attribute-Structure Correlated Patterns

❑ CiteSeer
  o 294,104 vertices
  o 782,147 edges
  o 206,430 attributes



Pattern induced by
{perform, system}

## Two-event Structural Correlations

# Two-event Structural Correlations

**How is the relationship between the sales of two products in a social network?**

**Attraction (positive correlation)**

**Repulsion (negative correlation)**

Video games

Computers

# A New Notion of Correlation

❑ Two-Event Structural Correlation (**TESC**)

- o Defined on graph structures

- o Capture relationships between distributions of two events on a graph

- o Events can be different things in different contexts:
  - ➤ Topics or products (social networks)
  - ➤ Virus (contact networks)
  - ➤ Intrusion alerts (computer networks)

# The Major Challenge

❑ **Simple idea**: compute the distance between the two events on the graph

- How near/far for **significant** positive/negative correlation?
- Hypothesis testing

# The Major Challenge

❑ **Simple idea**: compute the distance between the two events on the graph
- How near/far for **significant** positive/negative correlation?
- Hypothesis testing

# How To Measure?

**"reference nodes"**

■ Consistent   ■ Inconsistent



a:3->1
b:2->1

a:3->1
b:1->3

❑Concordance score

If the density changes are **consistent**

$$c(r_i, r_j) = \begin{cases} 1 & (s_a^h(r_i) - s_a^h(r_j))(s_b^h(r_i) - s_b^h(r_j)) > 0 \\ -1 & (s_a^h(r_i) - s_a^h(r_j))(s_b^h(r_i) - s_b^h(r_j)) < 0 \\ 0 & otherwise \end{cases}$$

If the density changes are **inconsistent**

❑ Density function   $s_a^h(r) = \dfrac{|V_a \cap V_r^h|}{|V_r^h|}$

# Kendall's Tau as The Measure

❑ **Kendall's Tau rank correlation** as a measure of TESC:

$$\tau(a,b) = \frac{\displaystyle\sum_{i=1}^{N-1}\sum_{j=i+1}^{N} c(r_i, r_j)}{\dfrac{1}{2}N(N-1)}$$

$N = |V_{a \cup b}^{h}|$ : the number of all reference nodes

❑ Costly computation!

**Sampling + significance testing**

44

# Real Events (DBLP)

**Highly positive pairs:**

| # | Pair | z-score | | | TC |
|---|------|---------|---------|---------|-----|
| | | h = 1 | h = 2 | h = 3 | |
| 1 | Texture **vs.** Image | 6.22 | 19.85 | 30.58 | 172.7 |
| 2 | Wireless **vs.** Sensor | 5.99 | 23.09 | 32.12 | 463.7 |
| 3 | Multicast **vs.** Network | 4.21 | 18.37 | 26.66 | 123.2 |
| 4 | Wireless **vs.** Network | 2.06 | 17.41 | 27.90 | 198.2 |
| 5 | Semantic **vs.** RDF | 1.72 | 16.02 | 24.94 | 120.3 |

**Highly negative pairs:**

| # | Pair | z-score | | | TC |
|---|------|---------|---------|---------|-----|
| | | h = 1 | h = 2 | h = 3 | |
| 1 | Texture **vs.** Java | -23.63 | -9.41 | -6.40 | 4.33 |
| 2 | GPU **vs.** RDF | -24.47 | -14.64 | -6.31 | 1.24 |
| 3 | SQL **vs.** Calibration | -21.29 | -12.70 | -5.45 | -0.62 |
| 4 | Hardware **vs.** Ontology | -22.31 | -8.85 | -5.01 | 3.38 |
| 5 | Transaction **vs.** Camera | -22.20 | -7.91 | -4.26 | 4.85 |

**z-score > 2.33: significant positive TESC**
**z-score < -2.33: significant negative TESC**

**TC: Transaction Correlation**

# Several Other Pattern Miners

❑ CoPaM: Cohesive Pattern Miner [Moser et al., SDM'09]

❑ Itemset-Sharing Patterns [Fukuzaki et al., 2013]

❑ Topological Pattern Miner [Prado et al., TKDE'13]

# CoPaM: Cohesive Pattern Miner

❑ Cohesive Pattern

A connected subgraph satisfying:

(1)  Feature subspace cohesion

(2)  Densely connected

Two cohesive patterns:

# Itemset-Sharing Patterns

❑ Goal
  o To find sets of subgraphs with common itemsets in a large graph.



**(a)**
**Itemset –attributed graph**

**(b)**
**Itemset-sharing subgraphs**
**with $\{i_1, i_2\}$.**

# Topological Pattern Miner

❑ A Topological Pattern

A set of vertex attributes and topological properties that strongly co-vary over the vertices of the graph

# Topological Pattern Miner

❑ A Topological Pattern
A set of vertex attributes and topological properties that strongly co-vary over the vertices of the graph



Prado et al., TKDE'13

$$P = \{h^+, i^-, \text{BETW}^+\}$$

The higher the value of attribute $h$, the lower the value of attribute $i$ and the higher the betweenness centrality of a vertex.

# Anomaly Detection in Graphs

❏ Various Interesting-ness/Anomaly Criteria

   e.g.,
- Bgp-lens: anomalies in internet routing updates.
  [Prakash et al., KDD'09]

- Oddball: anomalies in weighted graphs.
  [Akoglu et al., PAKDD'10]

- Heavy subgraphs in time-evolving networks.
  [Bogdanov et al., ICDM'11]

❏ Anomaly, Event, and Fraud Detection in Large Graph Datasets, Akoglu et al., http://www.cs.stonybrook.edu/~leman/wsdm13/

# Anomaly Vertices/ Regions



Comedy

Action

1. Target marketing

2. Recommendation systems

3. Social influence analysis

# Graph Iceberg (i.e., gIceberg)



(a) Original Graph

(b) Vertices Arranged by Aggregate Scores

# What aggregate functions?

❑ SUM and AVG [Yan et al., ICDE'10]



(a) x's 2-hop neighborhood    (b) y's 2-hop neighborhoo d

❑ Personalized PageRank Vector (PPV) [Page et al., Technical Report, 1999]
  o Reflect the proximity from one node to another w.r.t. the graph structure.

# Aggregate q-score:



$q$-score of $v$: $p_v(x) + p_v(y)$

$$\mathcal{P}_q(v) = \Sigma_{x|x \in V, q \in L(x)} p_v(x)$$

Aggregate over the Personalized PageRank Vector (PPV) to define the *closeness* of a vertex and an attribute *q.*

# Forward & Backward Aggregation



$q$-score of $v$: $p_v(x) + p_v(y)$

(a) FA

$q$-score of $v$: $f_x(p_x(v)) + f_y(p_y(v))$

(b) BA

# Results in Customer Network

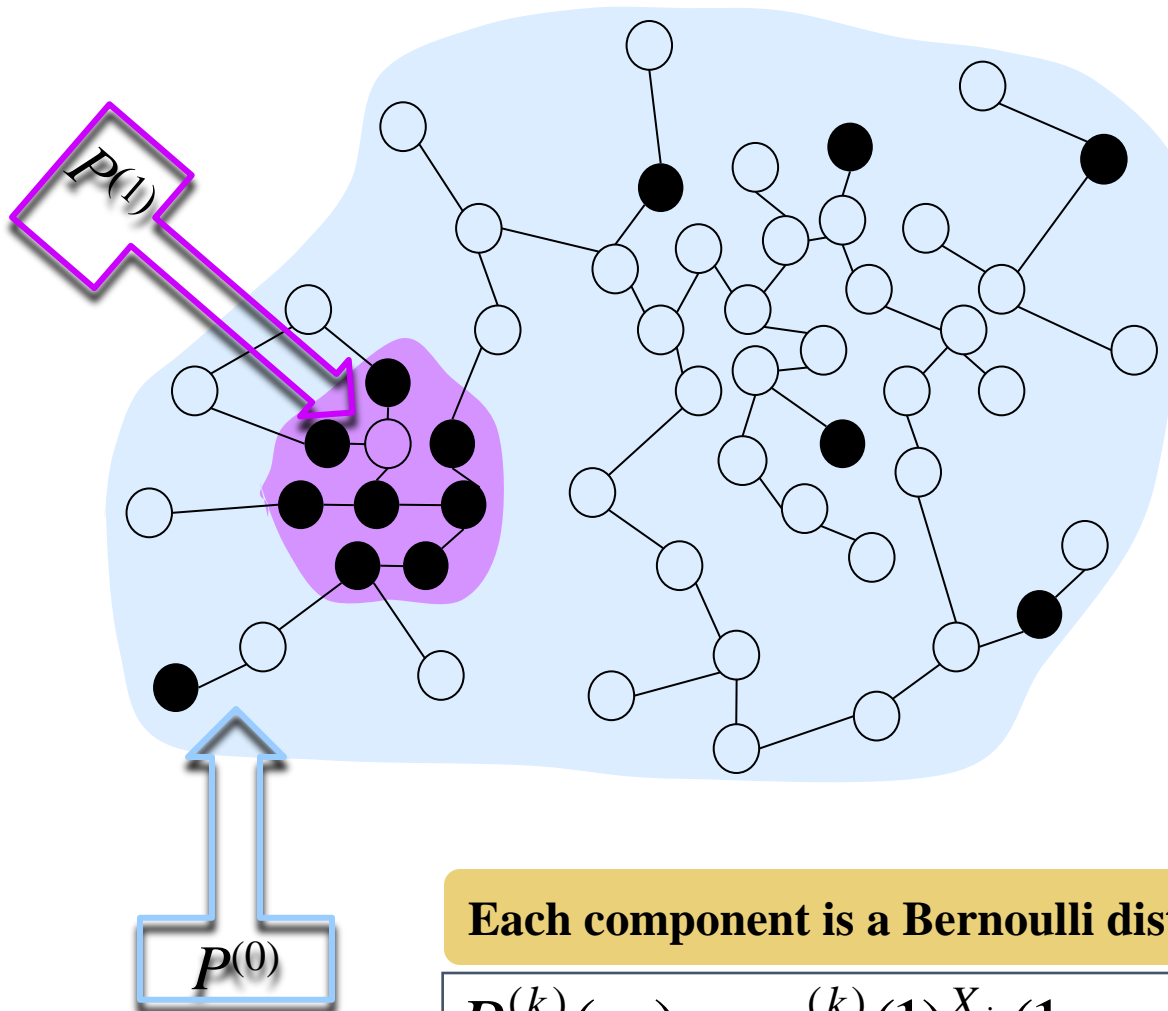❑ Recall indicates gIceberg's ability to retrieve real graph Iceberg vertices.

# Anomalous Regions (i.e., gAnomaly)



❑ Why does a disease occur more intensively in some portions of a network?

❑ Why do a subset of computers receive most of the attacks in the past day, and are they therefore targeted attacks?

# Data Models

Background: $V^{(0)}$  Anomaly: $V^{(1)}$



**Two generative processes: anomaly distribution & background distribution**

$$V = V^{(0)} \bigcup V^{(1)}, V^{(0)} \bigcap V^{(1)} = \varnothing$$

**One overall mixture**

$$P(v_i) = \sum_{k=0}^{1} \theta_i^{(k)} P^{(k)}(v_i)$$

With probability $\vartheta_i^{(0)}$, $v_i$ belongs to the background component $V^{(0)}$, and with $\vartheta_i^{(1)}$ the anomaly component $V^{(1)}$.

**Each component is a Bernoulli distribution.**

$$P^{(k)}(v_i) = \boldsymbol{p}^{(k)}(1)^{X_i} (1 - \boldsymbol{p}^{(k)}(1))^{1-X_i}$$

UCSB

SMU
SINGAPORE MANAGEMENT UNIVERSITY

# Regularized Data Likelihood

❑Un-regularized log data likelihood of vertex set *V*

$$L(V) = \sum_{v_i \in V} \log P(v_i) = \sum_{v_i \in V} \log \sum_k \theta_i^{(k)} P^{(k)}(v_i)$$

❑Regularized log data likelihood of vertex set *V*
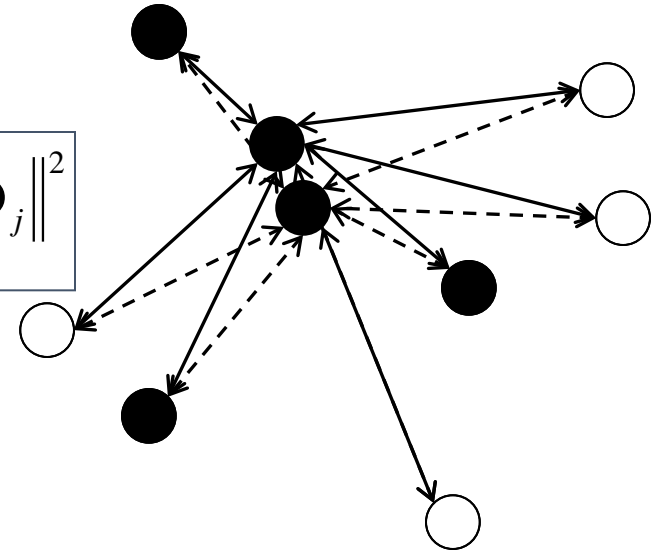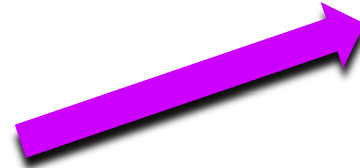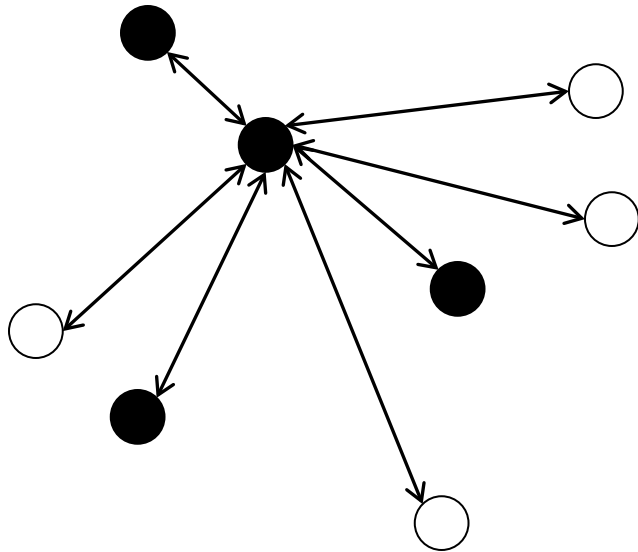
$$\widehat{L}(V) = L(V) - \lambda R_N(\Theta) + \gamma R_E(\Theta)$$

o Enhance connectivity within each component
  ➢ Network regularizer

o Enhance the interpretability of the mixture weights
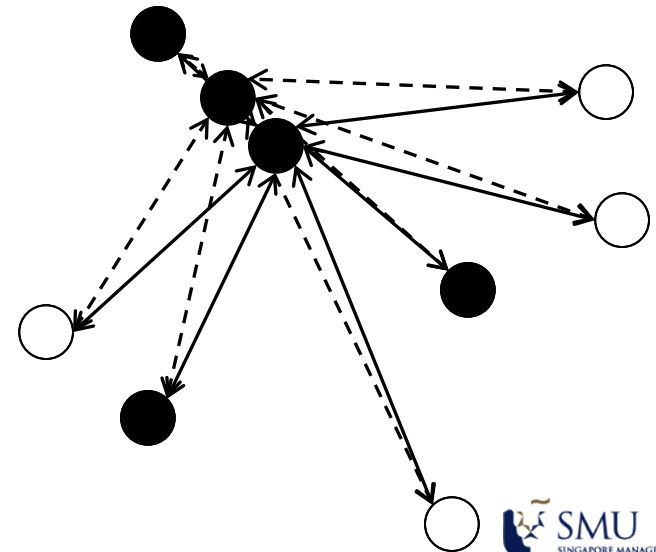  ➢ Entropy regularizer

# Network Regularizers

**Minimize the mean**

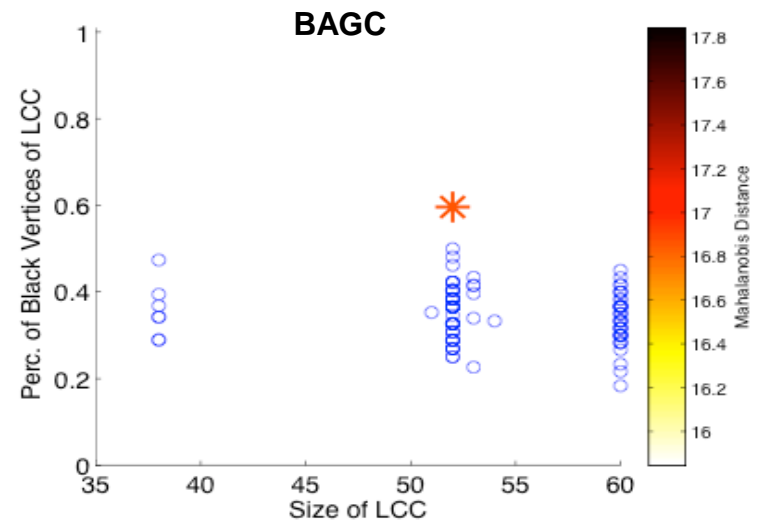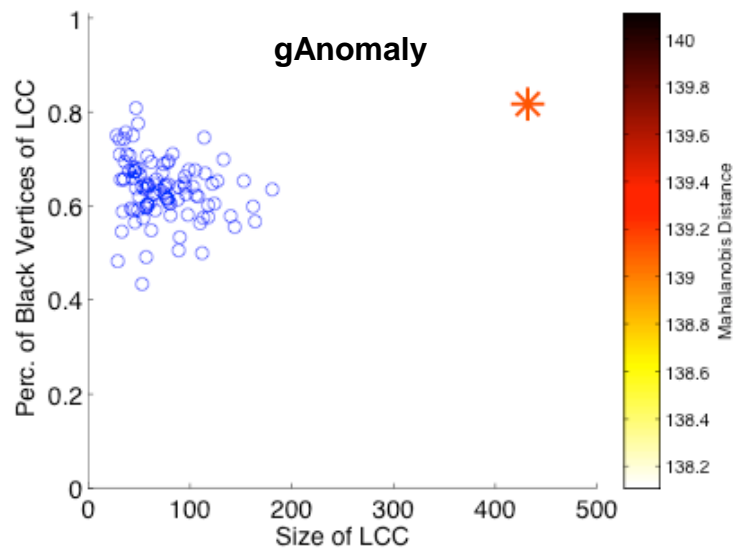$$R_N^{(1)}(\Theta) = \frac{1}{2} \sum_{v_i \in V} \frac{1}{|N(i)|} \sum_{v_j \in N(i)} \left\| \Theta_i - \Theta_j \right\|^2$$

**Minimize the minimum**

$$R_N^{(2)}(\Theta) = \frac{1}{2} \sum_{v_i \in V} \min_{v_j \in N(i)} \left\| \Theta_i - \Theta_j \right\|^2$$
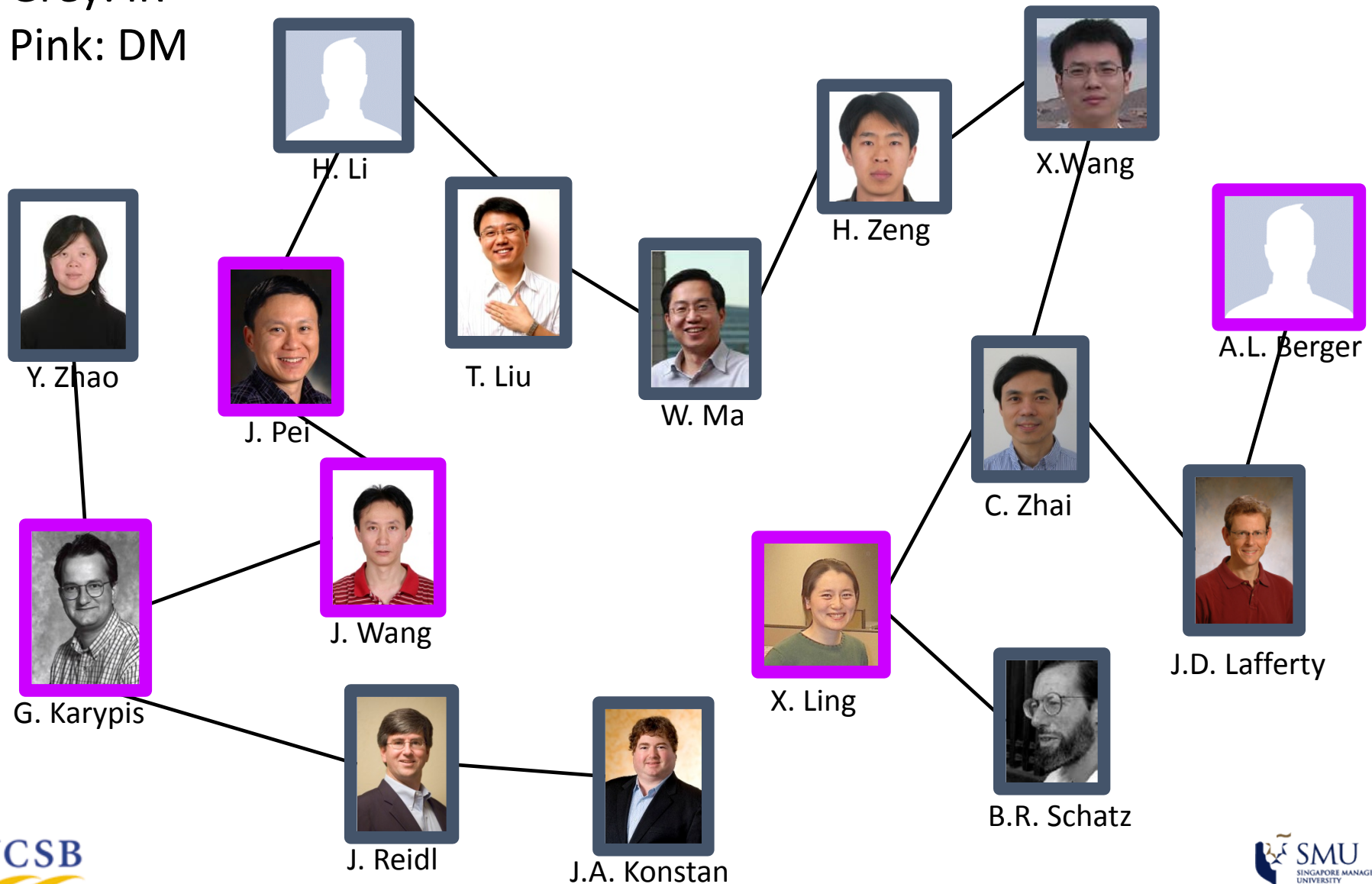
UCSB

SMU
SINGAPORE MANAGEMENT
UNIVERSITY

# Significance Evaluation

**M-distance Comparison**

## Case Study on DBLP-IR

Grey: IR
Pink: DM

# Part I: Brief Sum-Up

❑ Proximity Pattern ✔

❑ Attribute-Structure Correlations ✔

❑ Cohesive Pattern ✔

❑ Itemset-sharing Pattern ✔

❑ Graph Topological pattern ✔

❑ Graph Anomaly ✔

❑ Frequent Network Pattern (Part II)

# References

- Towards Proximity Pattern Mining in Large Graphs. [Khan et al., SIGMOD'10]

- Assessing and ranking structural correlations in graphs. [Guan et al., SIGMOD'11]

- Measuring Two-Event Structural Correlations on Graphs. [Guan et al., VLDB'11]

- Mining Attribute-structure Correlated Patterns in Large Attributed Graphs. [Silva et al., VLDB'12]

- Mining Cohesive Patterns from Graphs with Feature Vectors. [Moser et al., SDM'09]

- Finding Itemset-Sharing Patterns in a Large Itemset-Associated Graph. [Fukuzaki et al., PAKDD'10]

# References

- Mining graph topological patterns: Finding covariations among vertex descriptors. [Prado et al., TKDE'13]

- Bgp-lens: anomalies in internet routing updates. [Prakash et al., KDD'09]

- Oddball: Anomalies in Weighted Graphs. [Akoglu et al., PAKDD'10]

- Heavy Subgraphs in Time-Evolving Networks. [Bogdanov et al., ICDM'11]

- Giceberg: Towards Iceberg Analysis in Large Graphs. [Li et al., ICDE'13]

- A Probabilistic Approach to Uncovering Attributed Graph Anomalies. [Li et al., SDM'14]
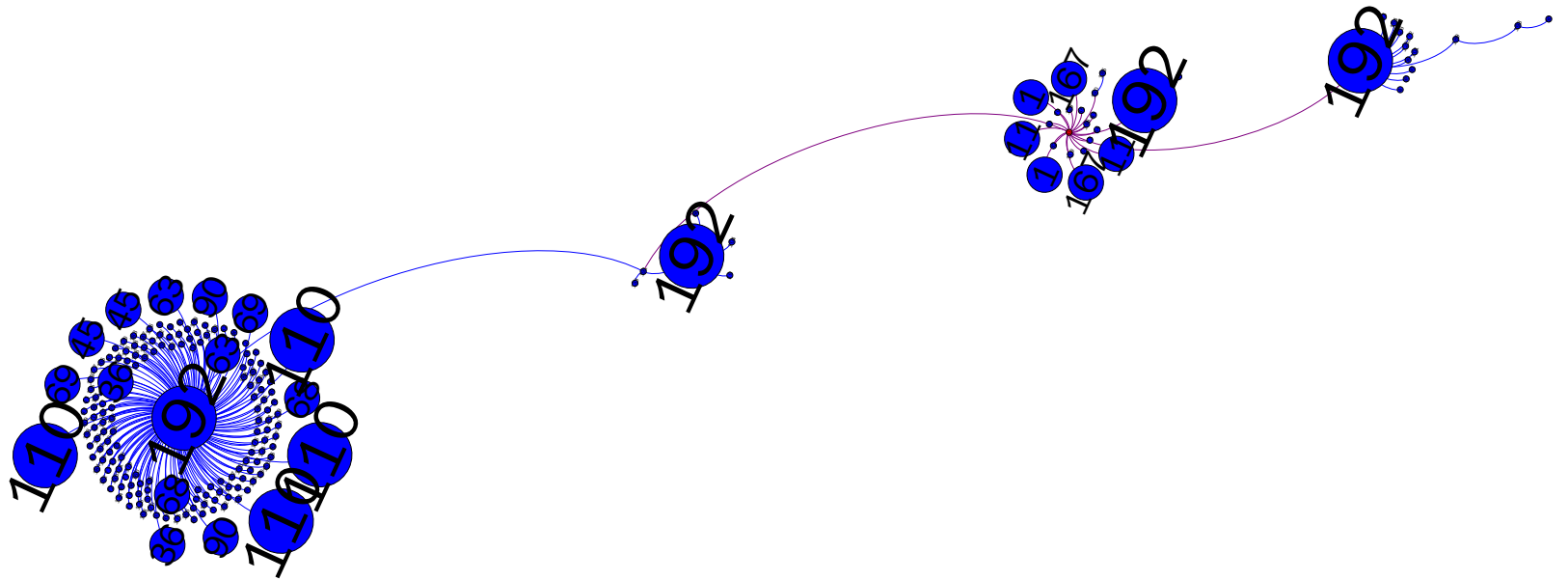
# Part II

## Frequent Network Pattern for Social Applications

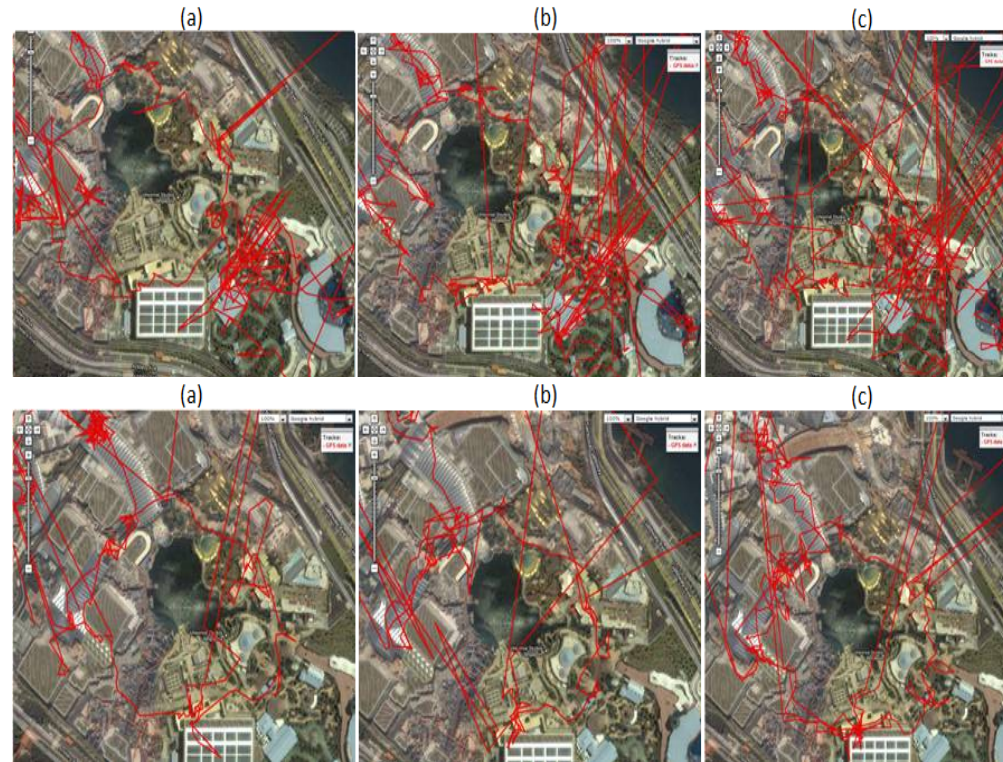# Frequent Network Patterns for Social Applications



**Community structure and dynamics Patterns**

# Frequent Network Patterns for Social Applications



**Information Diffusion and Viral Marketing Patterns**

# Frequent Network Patterns for Social Applications


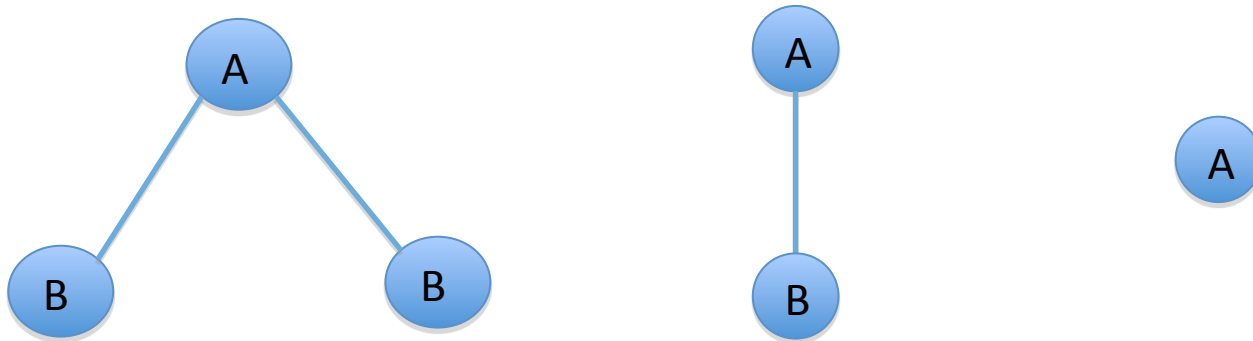
## Mobile User Trajectory and Location Patterns

# Some Important Characteristics

❑ Single-graph setting

❑ Huge network size

❑ User-specified constraints

# Single network mining is not easy

- Support counting is tricky when more than one pattern embedding would be considered in a single data graph.



- The anti-monotone property of support is violated.
- All the frequency-based pruning techniques must fail.
- No infrequent subgraph patterns can be pruned.

# How to fix it?

- Edge-disjoint

- Node-disjoint
    - Overlap graph
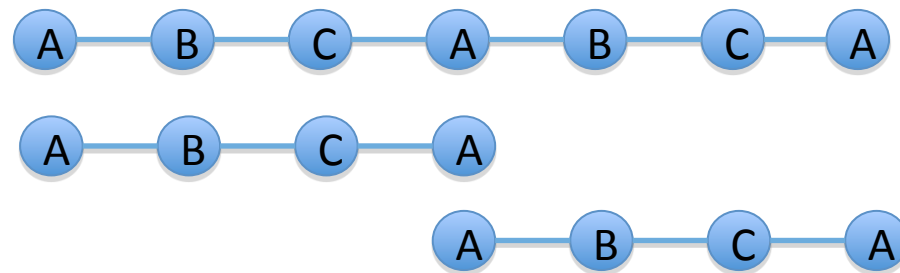    - Maximum Independent Set
    - MIS-support
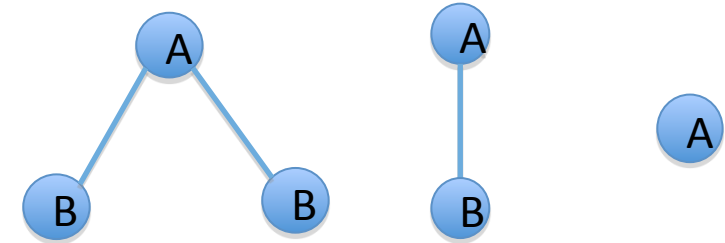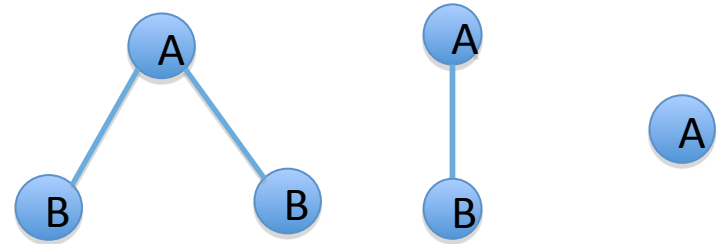        - Kuramochi and Karipus
    - MIS-support is antimonotone
    - But MIS-support is too restrictive!

- Harmful overlap support
    - Fiedler and Borgelt

# Some Important Characteristics

❑ Single-graph setting

❑ Huge network size

❑ User-specified constraints

# Not all patterns are wanted
## A Pattern Lattice Perspective

Pattern size

small

K

K+1

large

- ❑ Smaller patterns are affordable to be enumerated
- ❑ The number of patterns quickly explodes as pattern size grows larger
- ❑ The number of large patterns is small

**It's all about different ways to traverse this pattern lattice.**

# Some Important Characteristics
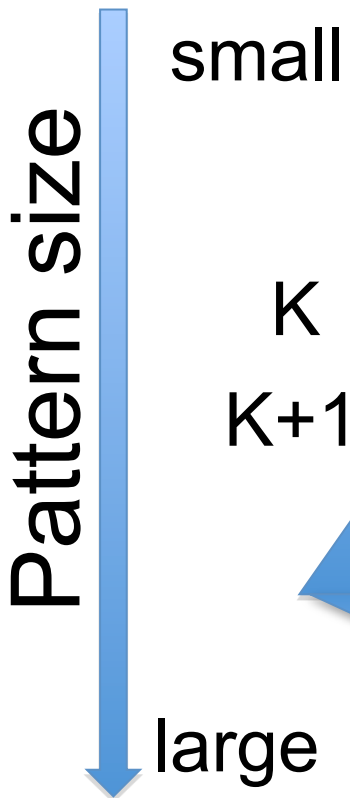
❑ Single-graph setting

❑ Huge network size

❑ User-specified constraints

# Traditional Solution for User-specified Constraints

☐ Mine all frequent patterns
  ☐ Pattern candidate generation
  ☐ Frequency checking in the data graph

☐ Check each frequent pattern for constraints

☐ Output those constraint-satisfying ones

*Increasing Pattern Size*

*Enumerate-and-check*

## Problem:  pattern explosion!

**What we really need ---- Blended frequent pattern mining with constraint-checking**

# Requirement Summary

❑ Single-graph setting
   ❑ Be able to handle both single-graph setting and graph-transaction setting.

❑ Huge network size
   ❑ Avoid enumerating all the pattern candidates before reaching those of interest.

❑ User-specified constraints
   ❑ The mining should be "direct".

# Representative Work

❑Graph-transaction setting
  ❑AGM, FSG, gSpan, FFSM, etc.
    ❑For mining *complete* pattern set.
    ❑Suffers from scalability issue for large patterns and input graphs due to exponential result size.
  ❑SPIN and MARGIN
    ❑Mining *maximal* patterns.
    ❑Still suffers from scalability issue as the number of maximal patterns could be formidable.
  ❑ORIGAMI
    ❑For mining a *representative* pattern set.
    ❑Returns a pattern set of mixed sizes.

# Representative Work

❑Single-graph setting
   ❑SUBDUE and SEuS
      ❑Pioneer work for mining *complete* pattern set.
      ❑Use different heuristics and work well for mining smaller patterns on certain classes of input graphs.
   ❑GREW
      ❑Mining *incomplete* pattern set.
      ❑Able to discover some large patterns, yet no guarantee on the pattern coverage of the answer set.
   ❑MoSS
      ❑State-of-the-art for mining *complete* pattern set.
      ❑Suffers from scalability issue for large patterns and input graphs due to exponential result size.

# SpiderMine: Frequent Large Patterns
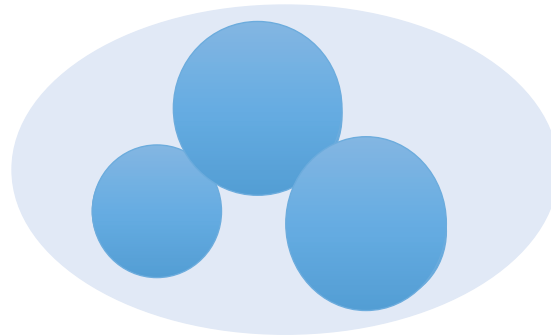## [Zhu et al. PVLDB'11]

- Graph data is getting ever bigger, and so are the patterns.
  - E.g., social networks like Facebook, Twitter, etc..

- Often, large patterns are more informative in characterizing large graph data.
  - E.g., in DBLP, small patterns are ubiquitous, larger patterns better characterize different research communities.
  - E.g., in software engineering, large patterns reveal software *backbones*

# Top-K largest frequent patterns

- To capture them exactly, no more and no less, we would have to generate all the smaller ones, which we cannot afford.

- A probabilistic solution, with user-defined error bound.

- SpiderMine --- designed to mine the top-K largest frequent patterns whose diameters are bounded by $D_{max}$ with a probability at least $1-\varepsilon$
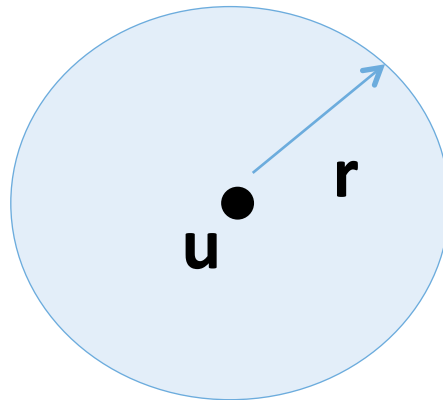
# Design Ideas

- How to capture large graph patterns?

- Observation:
  - Large patterns are composed of a large number of small components, called "spiders", which will eventually connect together after some rounds of pattern growth.

# r-Spider

- An *r-spider* is a frequent graph pattern P such that there exists a vertex u of P, and all other vertices of P are within distance r to u. u is called the head vertex.

# SpiderMine Overview

1) Mine the set S of all the r-spiders.

2) Randomly draw M r-spiders.

3) Grow these M r-spiders for t iterations. During the process, merge two patterns whenever possible.

4) Discard unmerged patterns.

5) Continue to grow the remaining ones to maximum size.

6) Return the top-K largest ones in the result.

◆ $t = D_{max}/2r$

# Large patterns vs small patterns

❑ Why can SpiderMine preserve large patterns and prune small ones with good chance?

1) Small patterns are much less likely to be hit in the random draw.
   ❖ First pruning at the initial random draw

2) Even if a small pattern is hit, it's even less likely to be hit multiple times.
   ❖ Second pruning after t iteration growth

3) The larger the pattern, the greater the chance it is hit and saved.

# How many r-spiders to draw?

LEMMA 2. *Given a network $G$ and a user-specified $K$, we have $P_{success} \geq \left(1 - (M+1)(1 - \frac{V_{min}}{|V(G)|})^M\right)^K$.*

With user-defined error threshold ε, we solve for M by setting

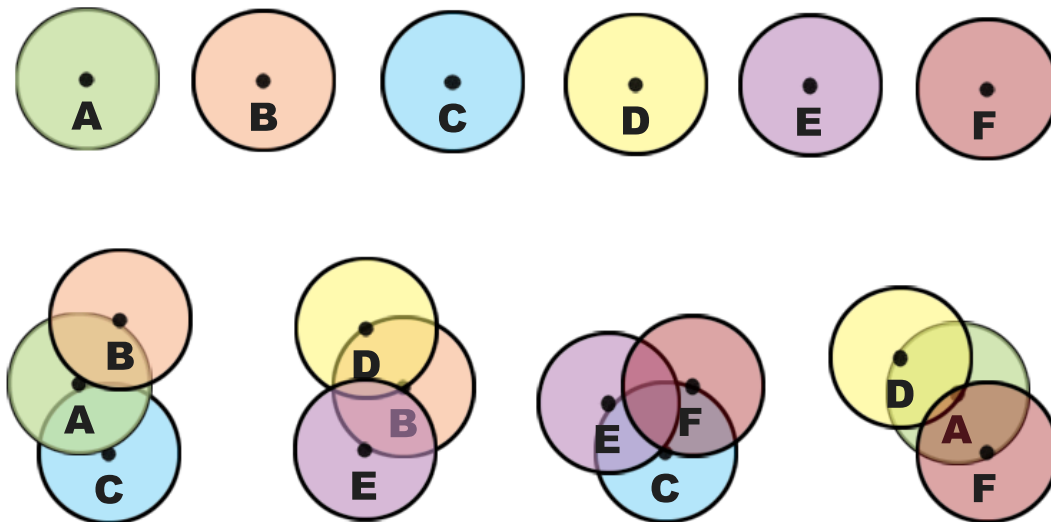$$\left(1 - (M+1)(1 - \frac{V_{min}}{|V(G)|})^M\right)^K = 1 - \epsilon$$

THEOREM 1. *Given a graph $G$, the error bound $\epsilon$, the diameter upper bound $D_{max}$, the support threshold $\sigma$ and $K$, with probability at least $1 - \epsilon$, SpiderMine returns a set $S$ of top-$K$ largest subgraphs of $G$ such that for each $P \in S$, $|P_{sup}| \geq \sigma$ and $diam(P) \leq D_{max}$.*

# Why Spiders?

(1) Reducing combinatorial complexity

Observation:

o Spiders are shared by many large patterns.

o Once obtained, they can be efficiently assembled to generate large patterns.
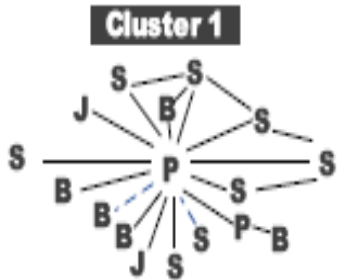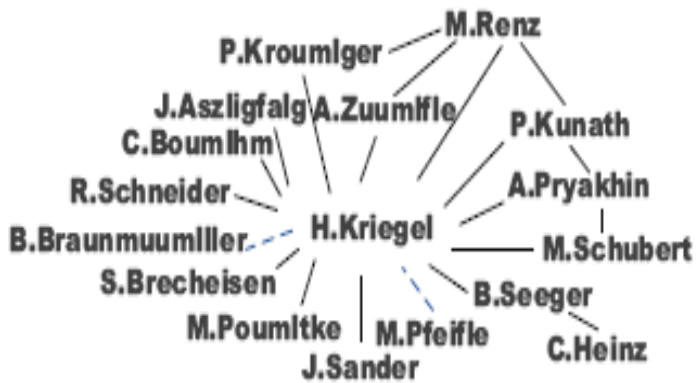
# Why Spiders?

(2) Reducing graph isomorphism checking

- ❑ We propose a novel graph pattern representation --- spider-set representation.

- ❑ A pattern is represented by the set of its constituent r-spiders.

- ❑ Two same patterns must have the same spider-set representation.

- ❑ Two same spider-set representations highly likely correspond to the same pattern.
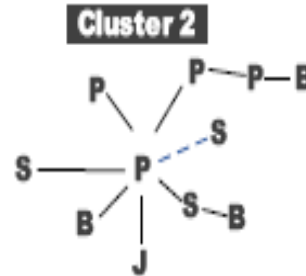
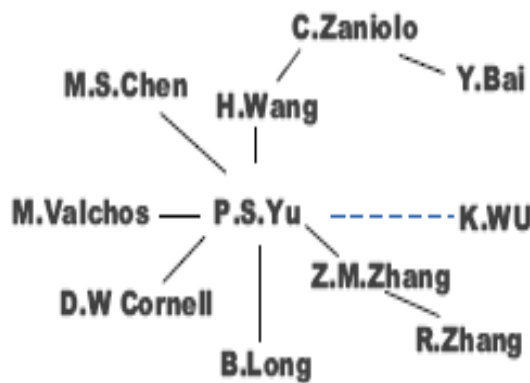- ❑ The larger the r, the more effective the pruning.
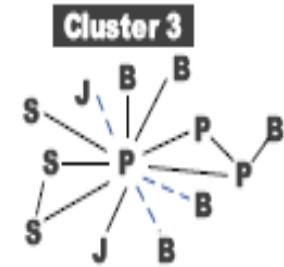
# DBLP Collaborative Patterns

# SpiderMine's response to the requirements

❑ Single-graph setting
- ❑ Handles overlapping embeddings

❑ Huge network size
- ❑ A probabilistic framework to capture only the top-k largest frequent patterns by leaping pattern growth.
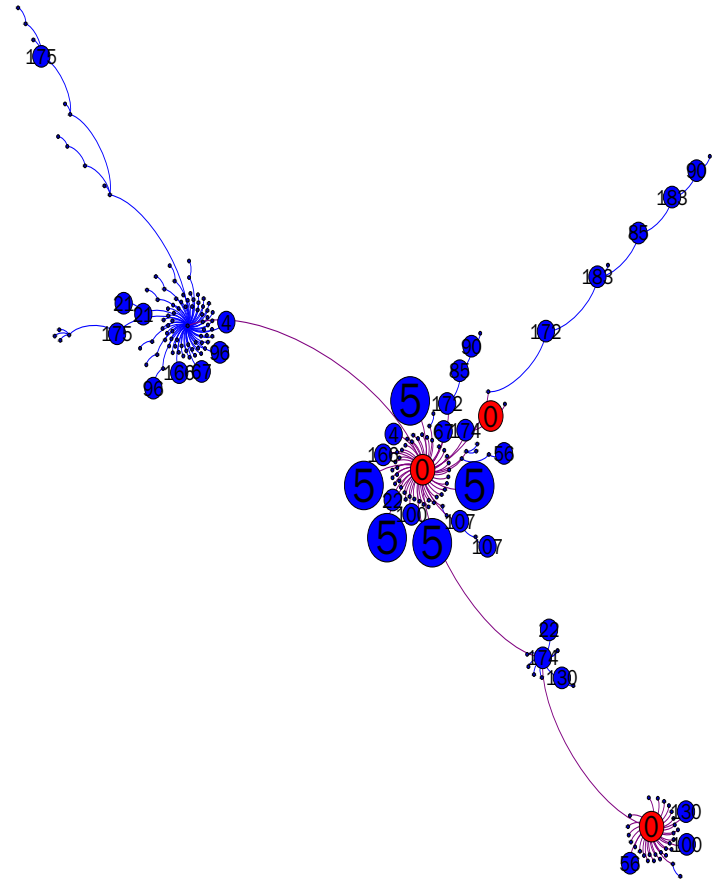
❑ User-specified constraints
- ❑ Pushed deep into the mining process by serving as the probabilistic bound

# SkinnyMine --- Catching the long and skinny ones
## [Zhu et al. SIGMOD'13]

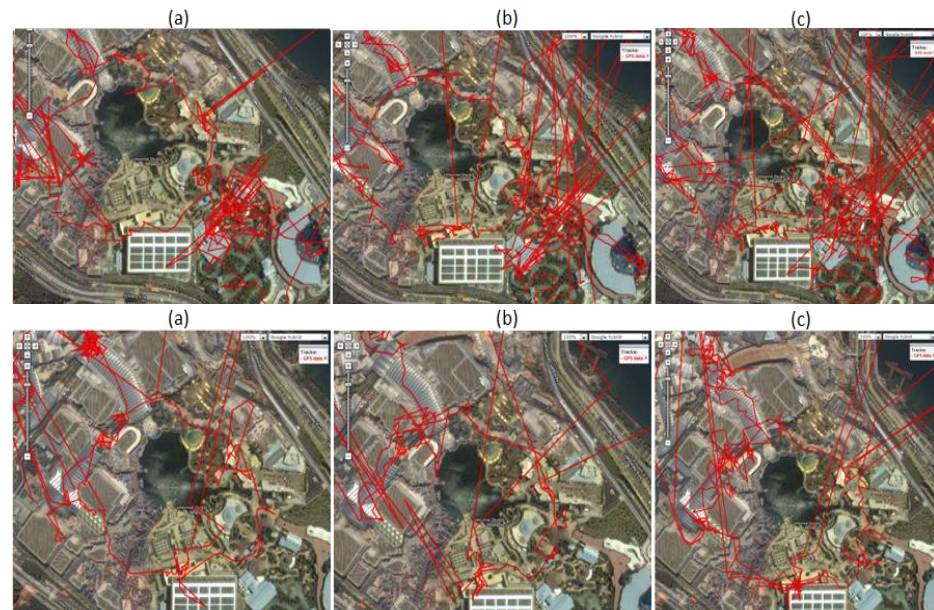❑Information diffusion patterns in social network

  ❑Viral topics, rumors, etc.

  ❑Diffusion paths are long backbones, associated ego-networks and local communities are of small radius

# A Motivating Example
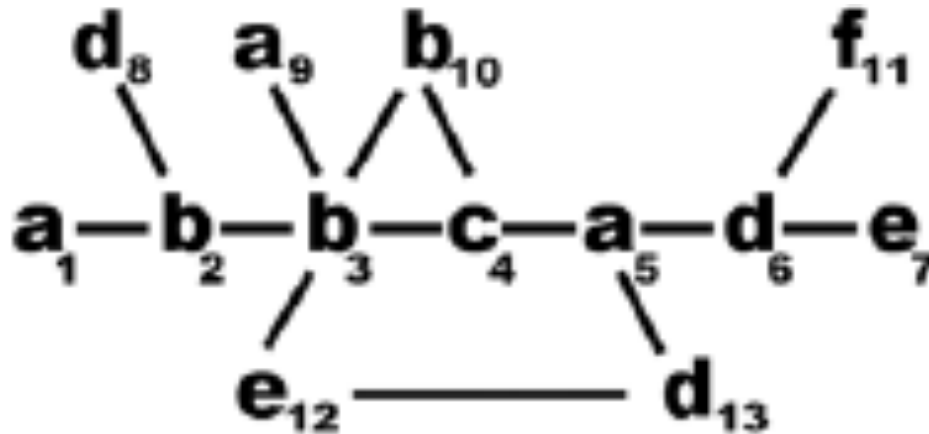
❑ Mobile trajectory patterns in LBS and Check-in data

  ❑ Trajectories are long backbones
  ❑ Associated point of interest

# Skinny Graph Patterns
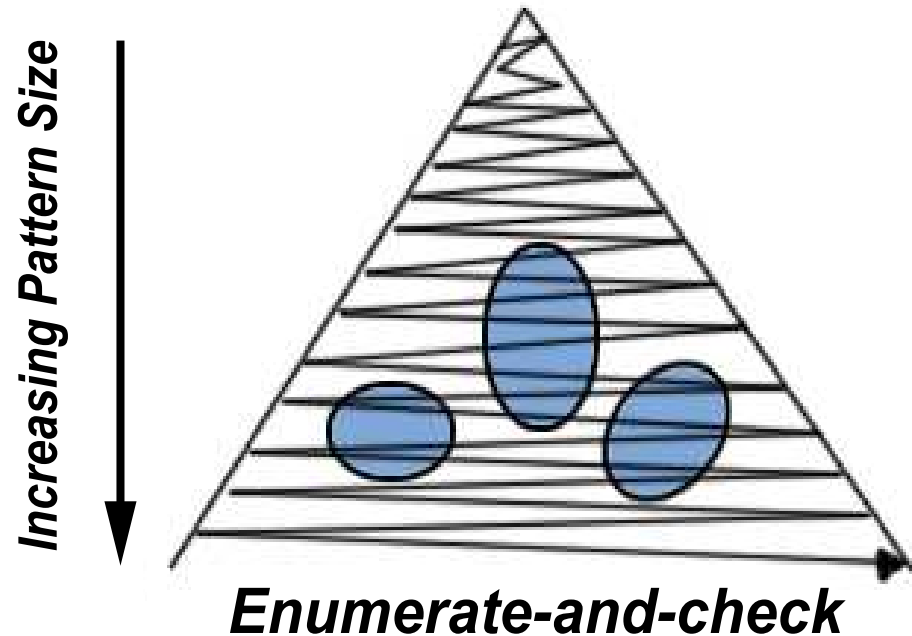
DEFINITION 8. **[$l$-Long $\delta$-Skinny Pattern Mining]** *Given a graph G, a frequency threshold $\sigma$, a length $l$ and a skinniness bound $\delta$, the problem of $l$-Long $\delta$-Skinny Pattern Mining $((l, \delta)$-SPM) is to find all $l$-long $\delta$-skinny subgraphs P of G such that $|E[P]| \geq \sigma$.*



A 6-long 1-skinny pattern
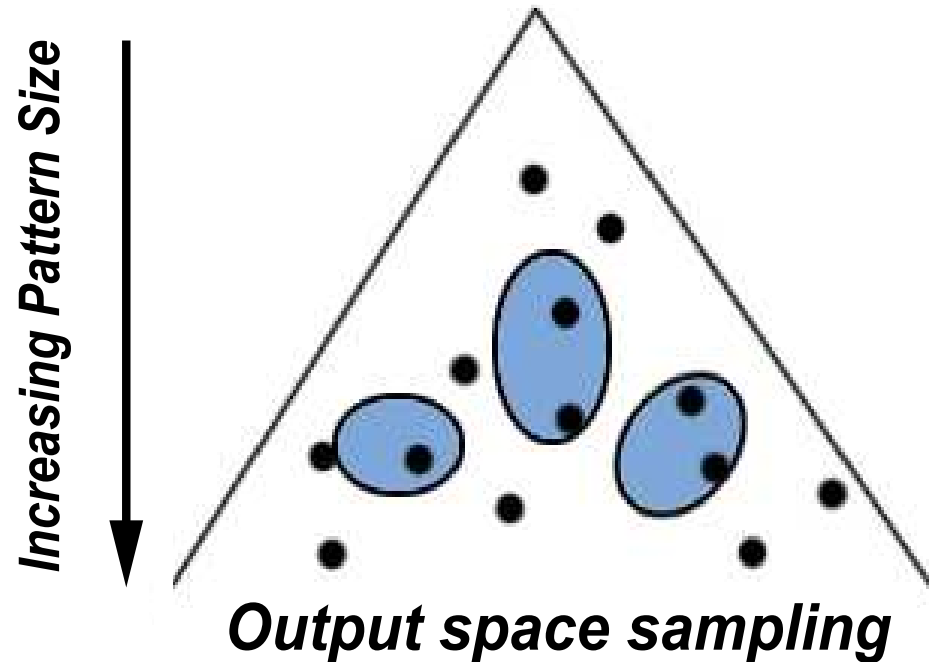
# Existing Mining Paradigm I

❑Enumerate-and-Check
  ❑Exhaustive enumeration of all frequent patterns and check if constraints are satisfied
  ❑Complete pattern result
  ❑Computationally unaffordable due to pattern explosion
  ❑Representative algorithms
    ❑gSpan, AGM, FSG,
    ❑Borgelt and Berthold, FFSM, etc.

*Increasing Pattern Size*



***Enumerate-and-check***

# Existing Mining Paradigm II

❑Output Space Sampling
- ❑A representative random sample of all frequent patterns and check if constraints are satisfied
- ❑Efficient
- ❑Returns only a sparse subset of the target result set
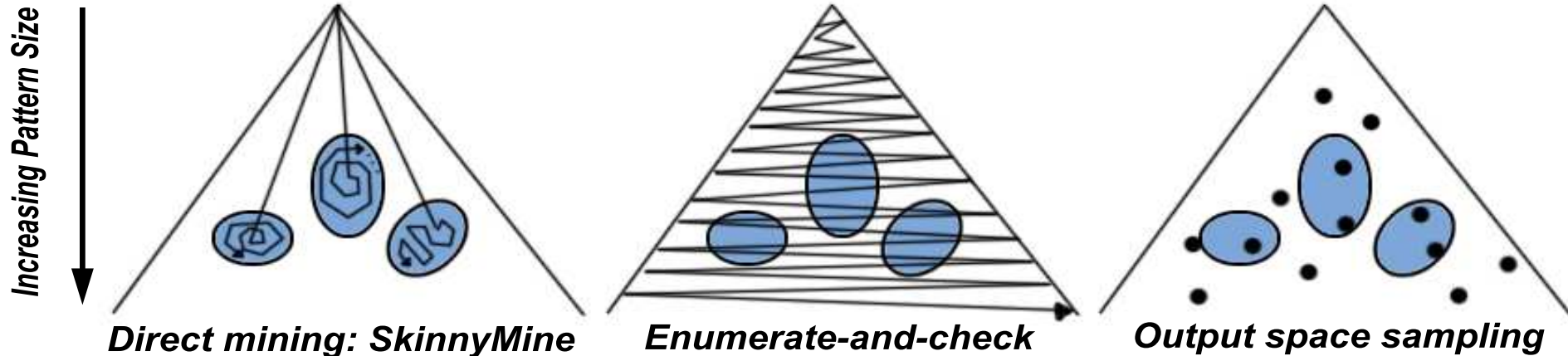
❑Representative algorithms
- ❑ORIGAMI, SpiderMine, etc.

*Increasing Pattern Size*

**Output space sampling**

# Paradigm Comparison

❏What we actually need
- ❏Complete result
- ❏Affordable computation time
- ❏Agile mining in terms of constraints

❏A mining framework that is "direct", "light" and "precise"

*A pattern lattice perspective for constrained frequent pattern mining*

*Increasing Pattern Size*



**Direct mining: SkinnyMine**     **Enumerate-and-check**     **Output space sampling**
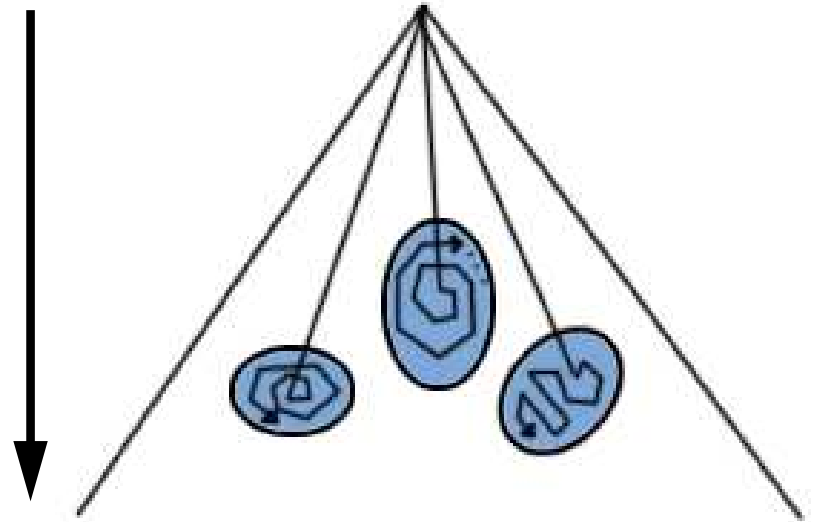
# Design Ideas

❑ Minimal Constraint-Satisfying Pattern
   ❑ Easy to compute with a low cost

❑ Constraint-Preserving Pattern Growth
   ❑ Non-redundant pattern generation

❑ How to generate each skinny pattern exactly once?
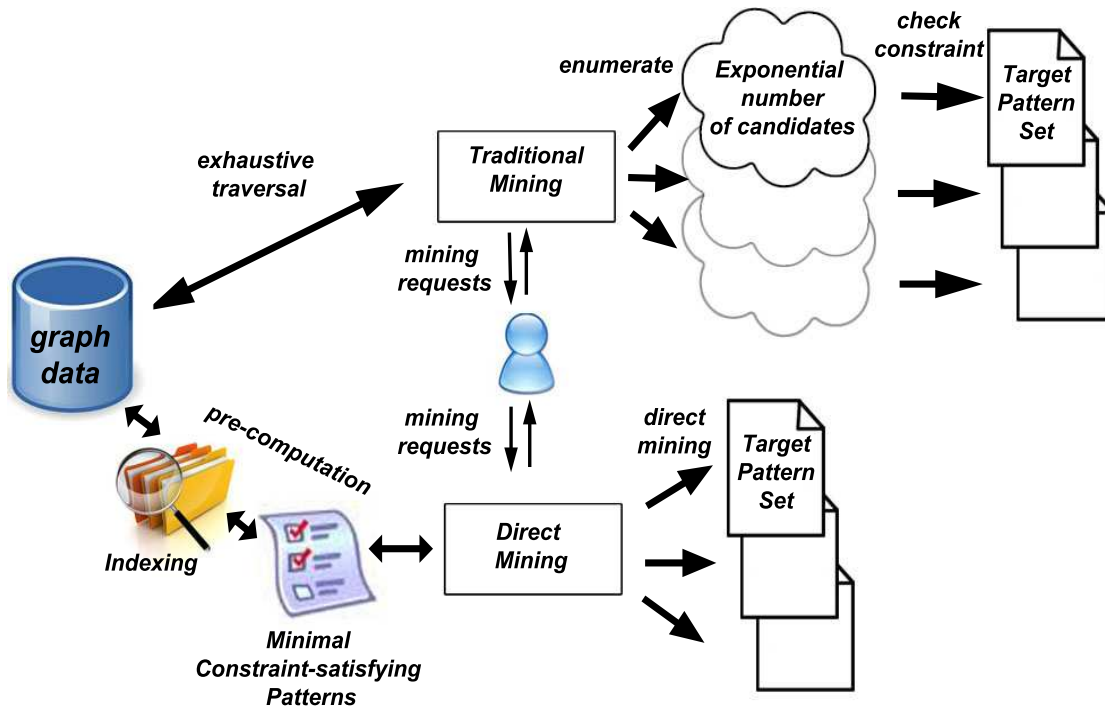
❑ Canonical diameters

*Increasing Pattern Size*

***Direct mining: SkinnyMine***

# Canonical Diameters

❑Each pattern has a unique canonical diameter.

❑The complete result set is partitioned by canonical diameters.

❑Patterns can be generated uniquely if they are generated only from the set corresponding to their canonical diameters.

# Algorithm Overview

❑Two stages
   ❑Stage I: Mining canonical diameters
   ❑Stage II:  Growing canonical diameters to patterns by level.

# Maintaining Canonical Diameters

1. **Constraint I: Diameter is not increased.**
   This means $\mathbb{D}(P') \leq \mathbb{D}(P)$, which guarantees that the edge extension does not create a longer diameter.

2. **Constraint II: $\mathbb{L}$ still realizes the shortest distance between $v_H$ and $v_T$.**
   This means in $P'$, $Dist(v_H, v_T) = |\mathbb{L}|$ and therefore $\mathbb{D}(P') \geq \mathbb{D}(P)$, which guarantees that the edge extension does not shorten the distance between $v_H$ and $v_T$.

3. **Constraint III: $\mathbb{L} \prec \mathbb{L}'$ for any newly generated diameter $\mathbb{L}'$ of the same length.**

# Efficient Constraint Maintenance

❑For each node v on the current growth frontier, we maintain two distances D^v_H and D^v_T. When adding new edge or new node u, we only need to check these two distances.

**Constraint I:  Diameter is not increased**

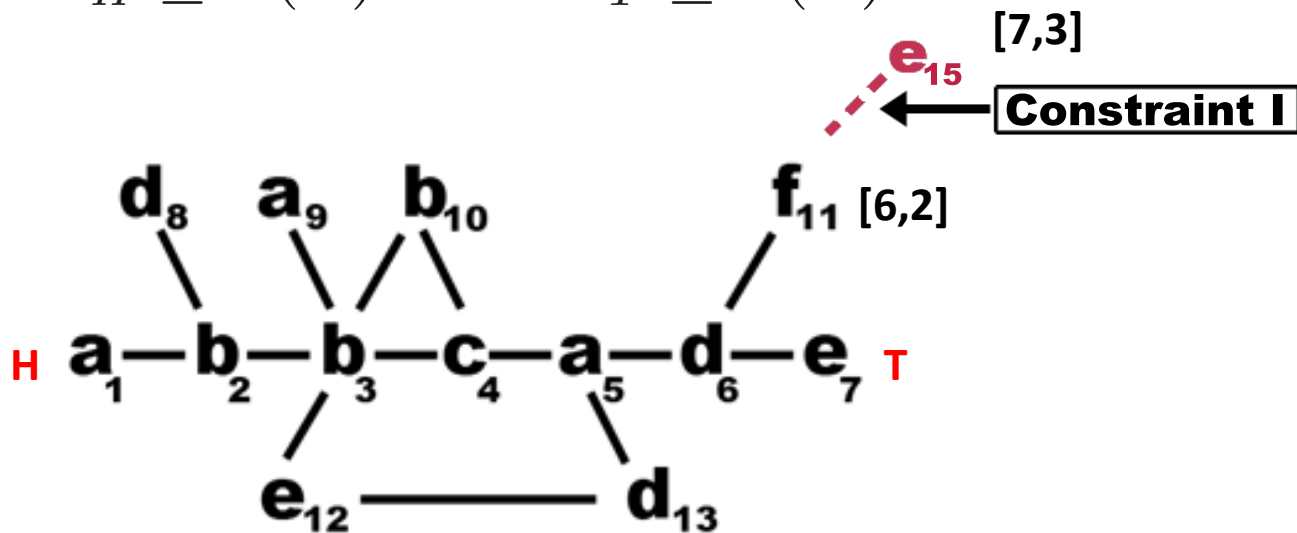$$D_H^u \leq \mathbb{D}(P) \ \ and \ \ D_T^u \leq \mathbb{D}(P)$$



D(P)=6

# Efficient Constraint Maintenance

❑For each node v on the current growth frontier, we maintain two distances D^v_H and D^v_T. When adding new edge or new node u, we only need to check these two distances.

**Constraint II: The shortest distance between H and T has not been shortened**

$$D_H^u + D_T^u \geq \mathbb{D}(P)$$

[3,2]

d$_8$    a$_9$    b$_{10}$ [3,4]    f$_{11}$

H    a$_1$ — b$_2$ — b$_3$ — c$_4$ — a$_5$ — d$_6$ — e$_7$    T

e$_{12}$ ——————— d$_{13}$

D(P)=6

# A General Direct Mining Framework

## ❑Two essential properties for constraints

PROPERTY 1. *[Reducibility]* *A graph constraint $C$ is called* reducible *if there exists a positive integer $k$ and a nonempty set $S$ of graph patterns such that for each $P \in S$, $|E(P)| \geq k$ and $f_C(P) = 1$, and for any pattern $P' \subset P$, $f_C(P') = 0$.*

PROPERTY 2. *[Continuity]* *A graph constraint $C$ is called* continuous *if for each pattern $P$ such that $f_C(P) = 1$, one of the following is true: (1)$P$ is a minimal constraint-satisfying pattern; (2)there exists at least one pattern $P'$ such that $P' \subset P$, $|E(P')| = |E(P)| - 1$ and $f_C(P') = 1$.*
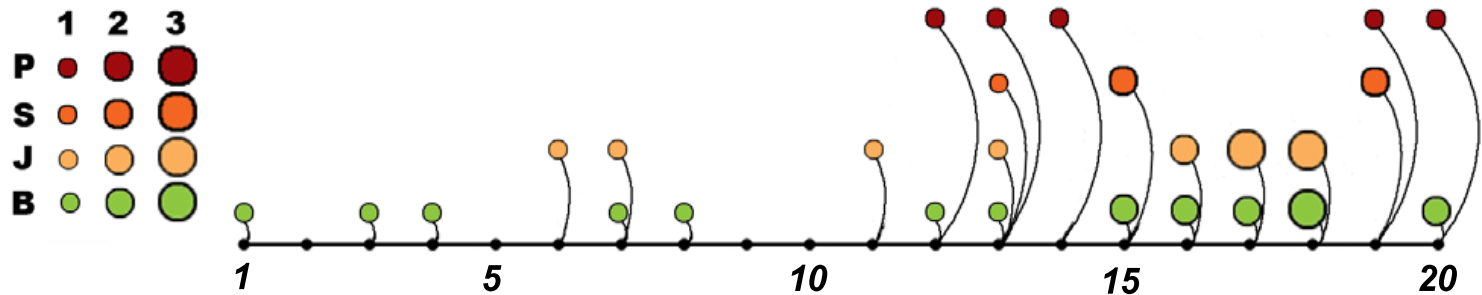
# DBLP Collaborative Patterns
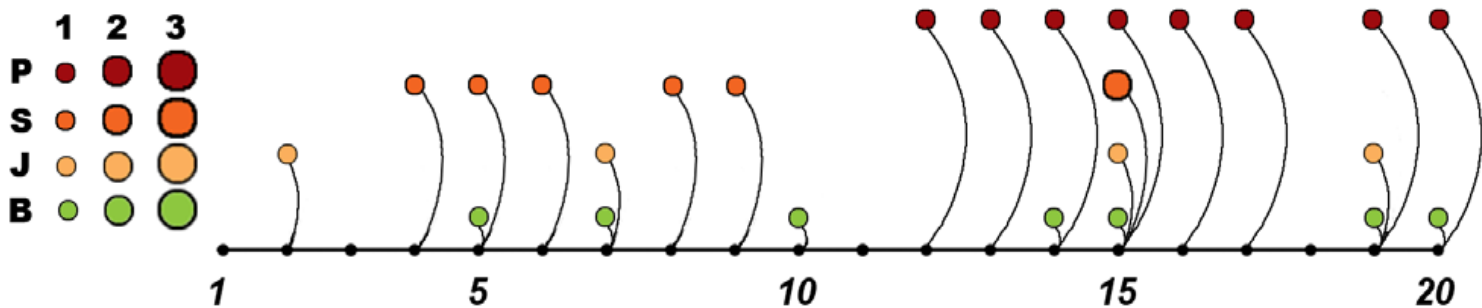
❑Real Data: DBLP



**Figure 21: DBLP skinny pattern example 1: A temporal collaborative pattern across 20 years.**

# Information Diffusion Patterns

❑Real Data:  Sina Weibo

# SkinnyMine's response to the requirements

❑ Single-graph setting
 ❑ Handles overlapping embeddings

❑ Huge network size
 ❑ A direct mining strategy to reach target patterns as quickly as possible, and locally identify only those of interest

❑ User-specified constraints
 ❑ Pushed deep into the mining process by guaranteeing the non-redundant pattern generation

# References

- F. Zhu, Q. Qu, D. Lo, X. Yan, J. Han, and P. S. Yu, *"Mining
- top-k large structural patterns in a massive network"*, PVLDB, 2011.
- F. Zhu, Z. Zhang and Q,Qu, *"A Direct Mining Approach To Efficient Constrained Graph Pattern Discovery"*. SIGMOD 2013.
- F. Zhu, X. Yan, J. Han, and P. Yu. *"gPrune: A constraint pushing framework for graph pattern mining"*, In PAKDD, 2007.
- M. A. Hasan, V. Chaoji, S. Salem, J. Besson, and M. J. Zaki. *"Origami: Mining representative orthogonal graph patterns"*. In ICDM 2007.

# Part III
# Collaboration Patterns

# Collaborations Are Everywhere

Question answering
in a social network

Customer service

Collaborator finding
in the academic world

Business referral

Software development

# Collaborations In This Talk

❑**Peer Interaction**
- o Patterns
- o Prediction
- o Optimization



❑**Team Joint Collaboration**
- o Patterns
- o Prediction
- o Optimization

# Collaborations In This Talk

**Peer Interaction**
- Patterns
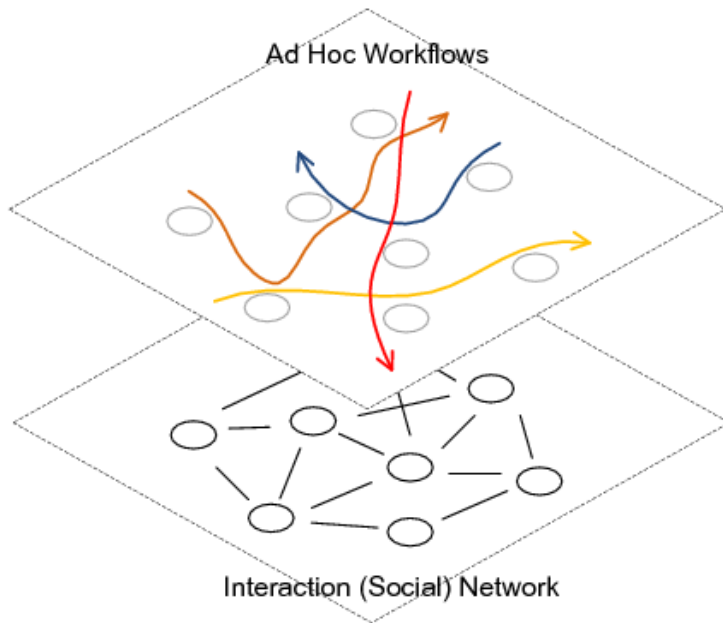- Prediction
- Optimization

**Team Joint Collaboration**
- Patterns
- Prediction
- Optimization

# Collaborative Networks

**Information Flow**



**Collaborative Network**

**Enterprise problem ticket**

| ID | Day | Entry |
|----|------|-------|
| 81 | 05-14 | New Ticket:   DB2 login fail |
| 81 | 05-14 | Transfer to Group SMRDX |
| 81 | 05-14 | Contacted Mary for recycling |
| 81 | 05-14 | Transfer to Group SSDSISAP |
| 81 | 05-14 | Status updated ... |
| 81 | 05-15 | Transfer to Group ASWWCUST |
| 81 | 05-15 | Web service checking |
| 81 | 05-18 | Could not solve the problem. |
| ... | ... | ... |
| 81 | 05-18 | Transfer to Group SSSAPHWOA |
| 81 | 05-22 | Resolved |

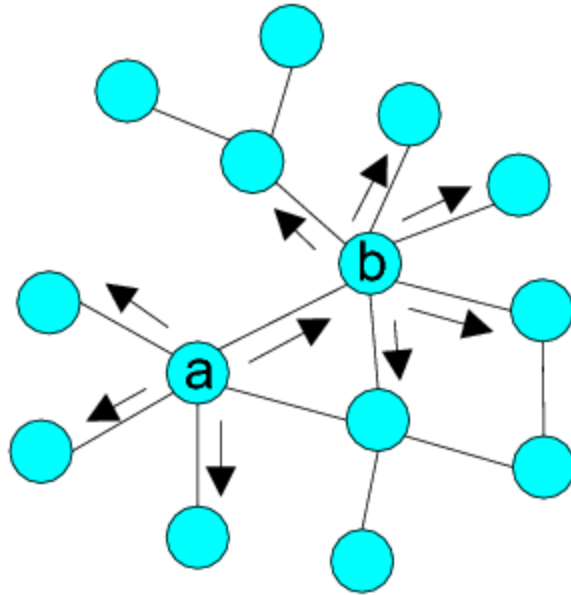**Eclipse bug record**

Table 1: Eclipse bug activity record.

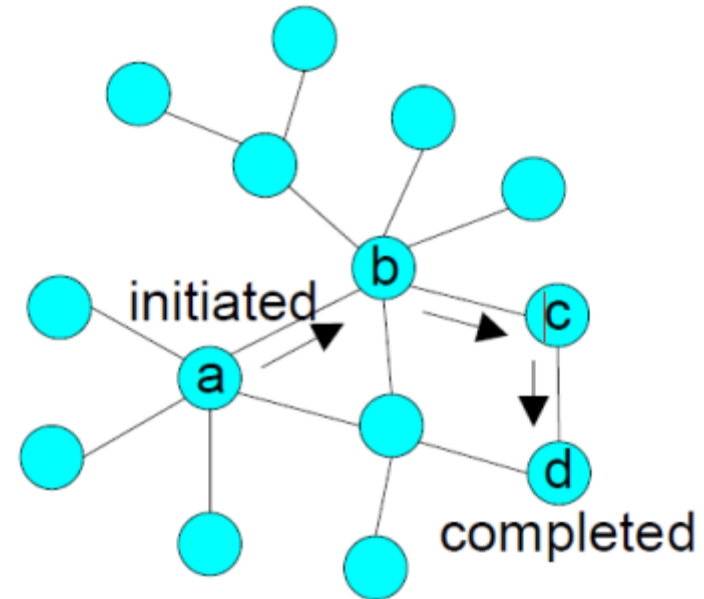| Bug description: NullPointerException referencing non-existing plugins. | | |
|------|------|------|
| **Who** | **When** | **Description** |
| dean | 2001-11-01 07:17:38 EST | Added component Core. Reassigned. |
| rodrigo | 2001-11-20 18:53:40 EST | Added component UI. Reassigned. |
| dejan | 2002-01-09 20:46:27 EST | Converted the unresolved plugin to a link. Fixed. |

https://bugs.eclipse.org/bugs/show_activity.cgi?id=325

# Social Networks VS. Collaborative Networks



**Traditional Social Networks**

**(Facebook, Twitter …)**

• Non-task driven information diffusion

**Collaborative Networks**

**(problem solving, team work ...)**

•Task-driven information flow

# We want to analyze…

- **How do experts make routing decisions?**

- **Who have made inefficient routing decisions?**

- **How to optimize the routing performance through targeted training?**

- **Can the completion time of a task be predicted so that one can act early for difficult tasks?**

   **In order to understand human factors in real task routing and the inefficiency bottleneck, take early actions, and further improve real collaborative networks.**

# We want to analyze…

- How do experts make routing decisions?
- Who have made inefficient routing decisions?
- How to optimize the routing performance through targeted training?
- Can the completion time of a task be predicted so that one can act early for difficult tasks?

**If we use intelligent algorithms to replace humans,**

- **Can we predict next best experts?**

- **How to model the collaborative networks?**

# We want to analyze…

- **How do experts make routing decisions?**

- **Who have made inefficient routing decisions?**

- **How to optimize the routing performance through targeted training?**

- **Can the completion time of a task be predicted so that one can act early for difficult tasks?**

- **Can we predict next best experts?**

- **How to model the collaborative networks?**

# Observations on Routing Decision Making

## Observation 1

Tasks with similar content, but different routing sequences

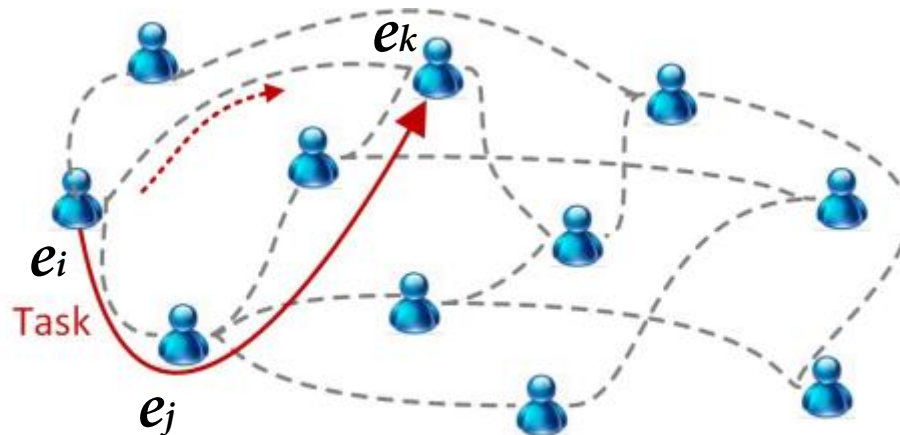e.g., two problem tickets in IBM IT service department

| Task ID | Task Content | Routing Sequence |
|---------|--------------|------------------|
| **492** | **Need password reset for kasperj on machine "pathfinder"". Route to NUS_N_DCRCHAIX** | **12 →505 → 1914→ 1915 →1916 →247** |
| **494** | **Need password reset for jhallacy on machine ""pathfinder"". Route to NUS_N_DCRCHAIX** | **12 → 13 → 86** |

Routing decision is not deterministic, given a certain task.

# Observations on Routing Decision Making

Observation 2
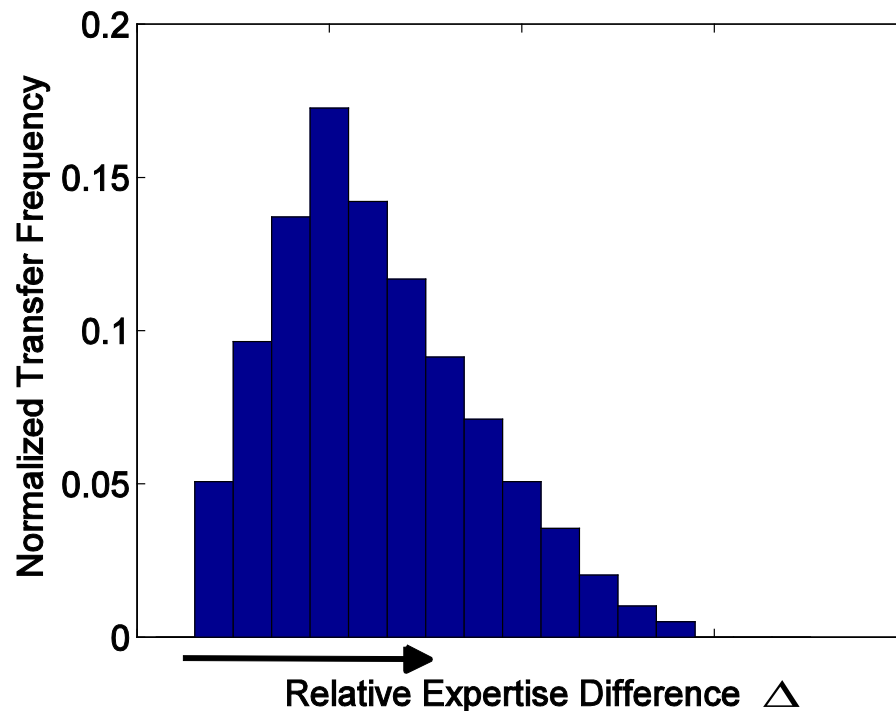An expert might not directly send a task to a resolver.



Is it because he does not understand the task very well, thus randomly routing it? Or

he believes the other expert has a better chance to solve it, or a better chance to find the right expert to solve it?

# Observations on Routing Decision Making

Observation 3

An expert tends to transfer a task to some expert whose expertise is *neither* too close *nor* too far from his own, i.e., not necessarily the final resolver!
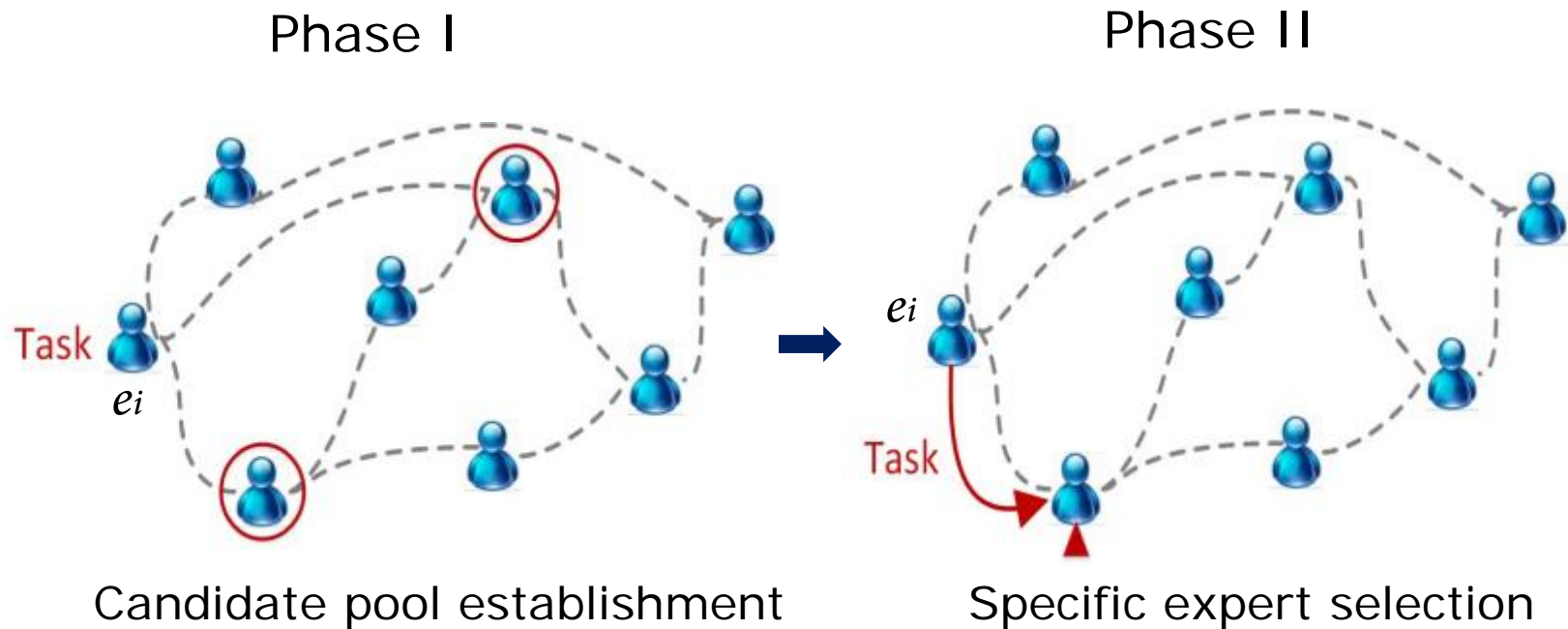
log-normal density:

$$f \propto \frac{1}{\Delta} e^{-\frac{[\ln \Delta - \mu]^2}{2\sigma^2}}$$

# Modeling Expert Decision Logic

❑ A Two-Phase Assumption

When an expert  transfers a task,

Phase **I**                                    Phase **II**



Candidate pool establishment            Specific expert selection

# Modeling Expert Decision Logic

❑6 particular routing patterns

**Phase I**

| | Task-Neutral Routing (TNR) | Task-Specific Routing (TSR) |
|---|---|---|
| Uniform Random (UR) | $TNR^{ur}$ | $TSR^{ur}$ |
| Volume-biased Random (VR) | $TNR^{vr}$ | $TSR^{vr}$ |
| Expertise Difference (EX) | $TNR^{ex}$ | $TSR^{ex}$ |

**Phase II**

# Modeling Expert Decision Logic

❏Generative model

Assume: the routing decision of an expert is generated through a mixture of 6 routing patterns.

Decision generation process:

For each expert $e_i$ to transfer tasks,
-Draw the mixture weights of 6 routing patterns:

$$\theta_i \sim \text{Dirichlet}(\alpha)$$

(reflecting $e_i$'s preferences over adopting different routing patterns)

-For each task $t$ to be transferred by expert $e_i$,
* Draw a pattern label: $Z_{i,t} \sim \text{Mult}(\theta_i)$
* Draw an expert from the candidate pool to receive $t$.

# Evaluation

❑ Evaluation Measure
- Completion Time (CT): Number of experts contacted before a task is resolved

- Mean Absolute Error (MAE)

$$\mathbf{MAE} = \frac{1}{|\mathbf{Test\ Set}|} \sum_{t \in \mathrm{Test\ Set}} |\widetilde{\mathbf{CT}_t} - \mathbf{CT}_t|$$

❑ Results

| DB2 | |
|---|---|
| **Models** | **MAE** |
| **TNR+TSR** | **0.08** |
| Miao *et al.* | 0.68 |
| Support Vector Regression | 0.80 |
| Bayesian regression | 0.84 |

# Experiments

❑ Resolution Efficiency of Routing Patterns TNR vs. TSR

| Experts favoring TNR | Weight on TNR patterns $>$ Weight on TSR patterns. |
|---|---|
| Experts favoring TSR | Weight on TNR patterns $<$ Weight on TSR patterns |

**DB2**

# Application: Optimizing Collaborations

❑ Which expert should be trained first to adopt the most efficient routing patterns? How much efficiency improvement can we expect?

- o Random: random selection

- o Frequent transferor: select the expert who transfers the most tasks

- o Least efficient: select the least efficient expert

| Methods | Efficiency Improvement (%) |
|---|---|
| Random | 0.27 |
| Frequent Transferor | 0.91 |
| Least Efficient | 1.21 |
| Recommendation using our model | 2.75 |

# We want to analyze…

- **How do experts make routing decisions?** ✔

- **Who have made inefficient routing decisions?** ✔

- **How to optimize the routing performance through targeted training?** ✔

- **Can the completion time of a task be predicted so that one can act early for difficult tasks?** ✔

- **Can we predict next best experts?** ⬅

- **How to model the collaborative networks?**

# Ticket Resolution

| ID | Time | Entry |
|---|---|---|
| 28120 | 2007-05-14 | New Ticket:  DB2 login failure |
| 28120 | 2007-05-14 | Transferred to Group SMRDX |
| 28120 | 2007-05-14 | Contacted Mary for recycling WAS |
| 28120 | 2007-05-14 | Transferred to Group SSDSISAP |
| 28120 | 2007-05-14 | Status updated ... |
| 28120 | 2007-05-15 | Transferred to Group ASWWCUST |
| 28120 | 2007-05-15 | Web service checking |
| 28120 | 2007-05-18 | Could not solve the problem. |
| ... | ... | ... |
| 28120 | 2007-05-18 | Transferred to Group SSSAPHWOA |
| 28120 | 2007-05-22 | Resolved |

❑ Managing problem tickets is a key issue in IT service industry.

❑ The efficiency of solving tickets highly depends on how they are routed.

❑ How to develop efficient automated algorithms for ticket routing?

# Automate Ticket Resolution:
## Sequence Mining Approach

A set of tickets reported to the expert network $\mathcal{T} = \{t_1, t_2, ..., t_m\}$



An interconnected network of experts

$$\mathcal{G} = \{g_1, g_2, ..., g_L\}$$

Goal: Minimize the average length of routing sequences

$$S = \frac{\sum_{i=1}^{m} |R(t_i)|}{m}$$

Routing sequence of tickets

$$R(t) = g_{init}(t) -> ... -> g_{res}(t)$$

# **Markov Modeling**

❑Capture ticket transfer patterns embedded in historical resolution sequences

- o State: an expert group holding the ticket
- o Transition probability: given the previous groups, the probability of the ticket being transferred to group $g_i$, in the next step

$$
P(g_i \mid S_{(k)}) = \left\{ \begin{array}{ll} N(g_i, S_{(k)}) / N(S_{(k)}) & if \quad N(S_{(k)}) \\ 0 & else \end{array} \right.
$$

# of instances that a ticket is transferred to group *gi , after being processed by S(k)*

# of instances with a set of group transfers

# Routing Algorithms

❑First-order Memoryless (FM)
  o Relies on the <u>current state</u> to make transfer decisions

# Routing Algorithms

❑First-order Multiple active State (FMS)
 o Ticket transfer decision is based on <u>any one of the past states</u> (instead of the current state).

# Routing Algorithms

❑ Variable order Multiple active State (VMS)
  o With multiple active states, use <u>higher-order</u> whenever it can make more confident prediction
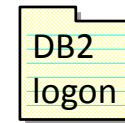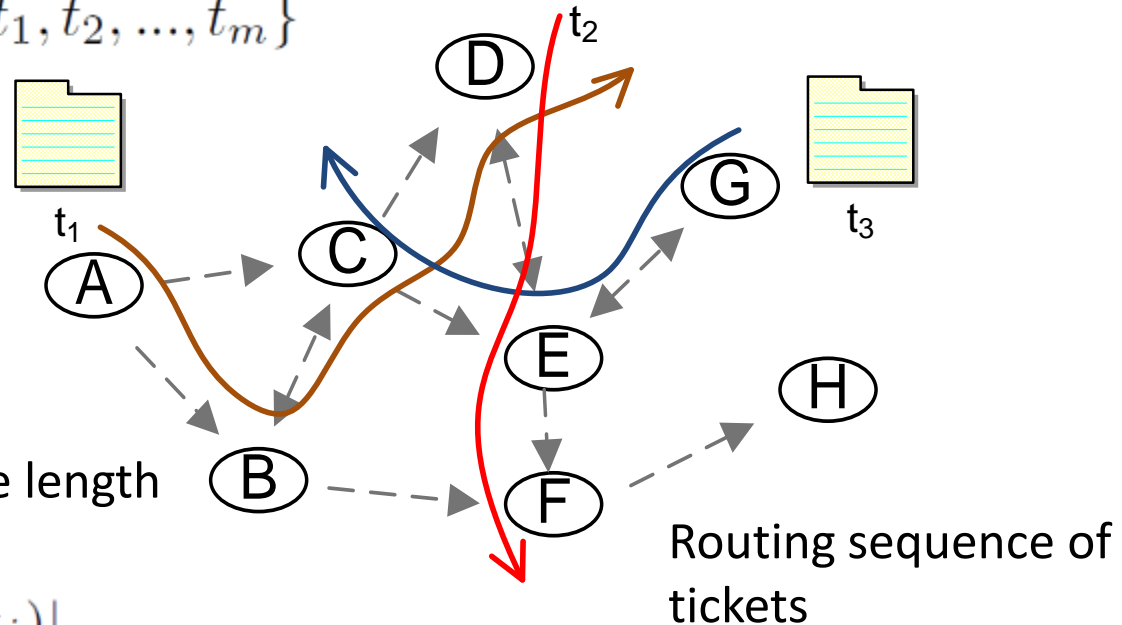


(a) Routing Steps in the model

| $S_{(2)}$ | Next group $g_i$ | $P(g_i \mid S_{(2)})$ |
|---|---|---|
| {A,C} | E | 0.6 |
| {A,C} | D | 0.2 |
| {A,C} | B | 0.1 |
| {A,C} | G | 0.1 |
| {E,C} | F | 0.6 |
| {E,C} | G | 0.4 |
| ... | ... | ... |

(b) $2^{nd}$ order probability table

# Incorporate Ticket Content:
## Generative Models

A set of tickets reported to the expert network

$$\mathcal{T} = \{t_1, t_2, ..., t_m\}$$

DB2 logon

**Word description of tickets**

$$\mathcal{W} = \{w_1, w_2, ..., w_n\}$$

An interconnected network of experts

$$\mathcal{G} = \{g_1, g_2, ..., g_L\}$$

Goal: Minimize the average length of routing sequences

$$S = \frac{\sum_{i=1}^{m} |R(t_i)|}{m}$$

Routing sequence of tickets

$$R(t) = g_{init}(t) - > ... - > g_{res}(t)$$

# Resolution Model (RM)

❑Each expert has an expertise profile
  o An expert is likely to be able to resolve tickets similar to what he/she has resolved previously



Tickets resolved by expert E

$$P_{g_i} = [P(w_1|g_i), P(w_2|g_i), ..., P(w_n|g_i)]^T$$

# Transfer Model (TM)

❑Expertise awareness between experts

  o An expert transfers similar tickets to another expert



Tickets transferred from expert B to expert F

$$P_{e_{ij}} = [P(w_1|e_{ij}), P(w_2|e_{ij}), ..., P(w_n|e_{ij})]^T$$

# Optimized Network Model (ONM)

❑Transfer profiles optimized for the entire expert network

$$\mathcal{L} = \prod_{t \in \mathcal{T}} P(R(t)|t)$$

$$P(R(t)|t) \quad = \quad P(g_1|t)P(g_2|t,g_1)P(g_3|t,g_2)P(g_3|t,g_3)$$

$$P(g_j|t,g_i) \quad = \quad \frac{P(t|e_{ij})P(g_j|g_i)}{Z(t,g_i)}$$

$$= \quad \frac{(\prod_{w_k \in t} P(w_k|e_{ij})^{f(w_k,t)})P(g_j|g_i)}{Z(t,g_i)}$$

$$Z(t,g_i) = \sum_{g_j \in \mathcal{G}} P(t|e_{ij})P(g_j|g_i)$$

UCSB

SMU
SINGAPORE MANAGEMENT
UNIVERSITY

# Routing Algorithms: Ranked Resolver

❑Match the ticket content with the expertise profiles



P(D) = 20%
DB2 1%
Logon 2%
failure 20%

A: 2.2 e-4

P(A) = 30%
DB2 10%
Logon 5%
failure 15%

DB2
Logon
failure

D: 8.0 e-6

P(F) = 10%
DB2 5%
Logon 1%
failure 1%

F: 5.0 e-7

$$P(g_i|t) = \frac{P(g_i)P(t|g_i)}{P(t)} \propto P(g_i) \prod_{w_k \in t} P(w_k|g_i)^{f(w_k,t)}$$

# Routing Algorithms: Greedy Transfer

❑Match the ticket with the transfer profiles



$$Rank(g_j) \propto \max_{g_i \in R(t)} P(g_j|t, g_i)$$

$$P(g_j|t, g_i) = \frac{P(g_j|g_i) \prod_{w_k \in t} P(w_k|e_{ij})^{f(w_k, t)}}{\sum_{g_l \in \mathcal{G}} P(g_l|g_i) \prod_{w_k \in t} P(w_k|e_{il})^{f(w_k, t)}}$$

# Routing Algorithms:
## Holistic Routing

❑All possibilities are explored

# Evaluation

# We want to answer…

- **How do experts make routing decisions?** ✔

- **Who have made inefficient routing decisions?** ✔

- **How to optimize the routing performance through targeted training?** ✔

- **Can the completion time of a task be predicted so that one can act early for difficult tasks?** ✔

- **Can we predict next best experts?** ✔

- **How to model the collaborative networks?** ⬅
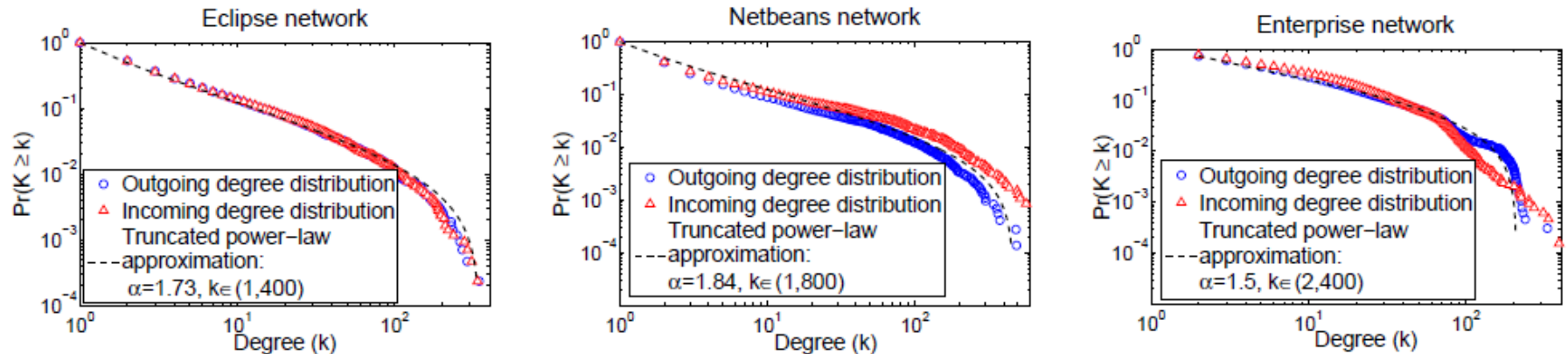
# Collaborative Network Observations


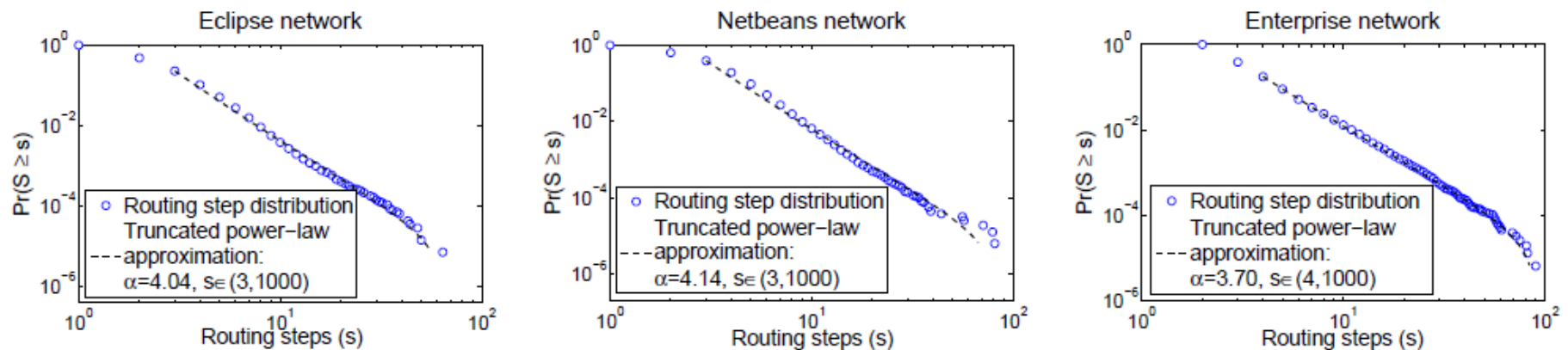
Figure 2: Degree distribution of collaborative networks.

Figure 3: Routing steps distribution of problem solving in collaborative networks.

# Collaborative Network Observations

❑ Clustering effect

    ○ Eclipse or Netbeans network

- Developers are loosely organized
- <u>Lower clustering coefficients</u>

    ○ Enterprise network

- Agents are more organized in groups by specialty
- Agents inside a group is more likely to collaborate with each other
- <u>Higher clustering coefficient</u>

❑ Clustering coefficient of studied networks

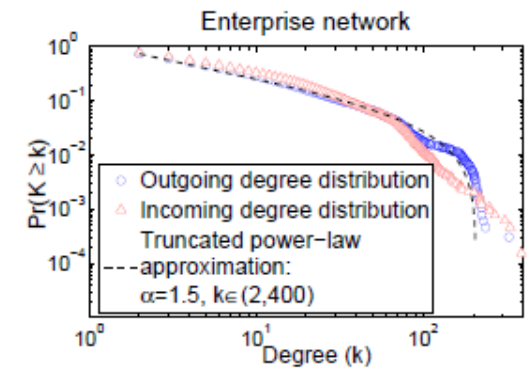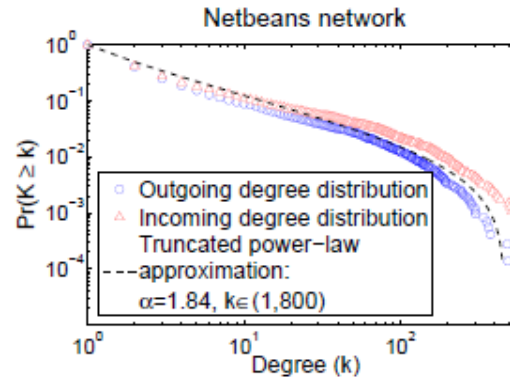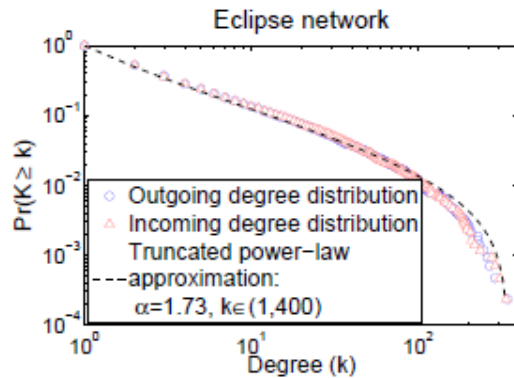| Eclipse Network | NetBeans Network | Enterprise Network |
|:---:|:---:|:---:|
| 0.198 | 0.21 | 0.35 |

# Modeling Approach

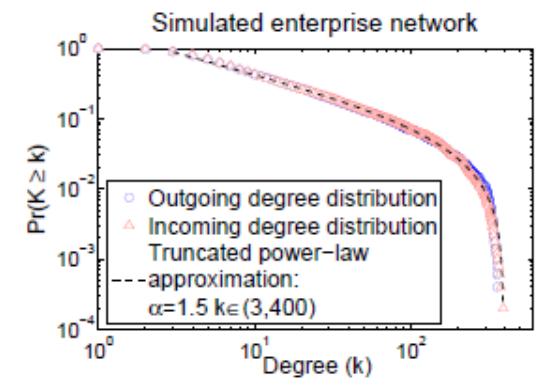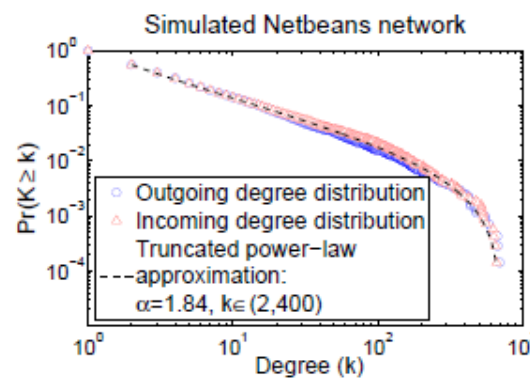❑ Network Model
 o Simulates static network connectivity

❑ Routing Model
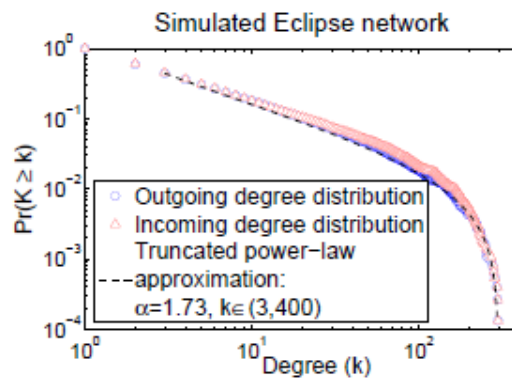 o Simulates dynamic human behavior in information routing

# Network Model: Simulation results
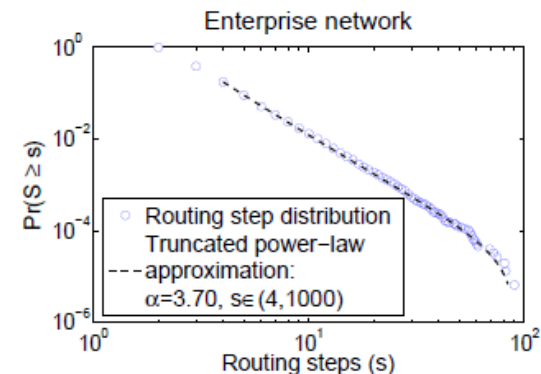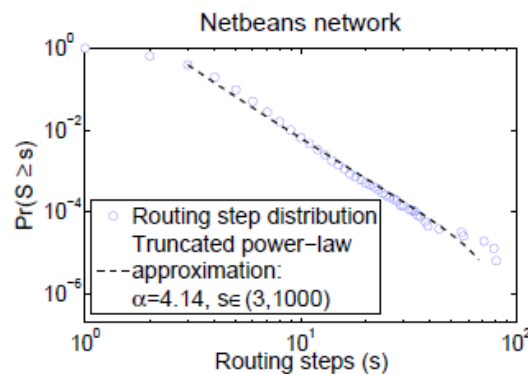
# Routing Model: Simulating Task Execution

❑Intuitions

- o If an expert cannot resolve a problem, she tends to forward it to a neighbor "closer to resolver"

- o If an expert doesn't know which neighbor is closer to resolver, she tends to choose a well-connected neighbor who is more likely to be connected to the resolver

- o Among all neighbors, an expert randomly chooses one to forward the task based on some probability distribution

# Routing Model: Simulation Results



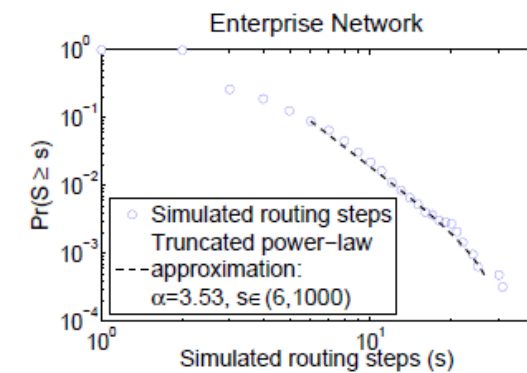Note: Simulation results were obtained by running proposed routing model on a 2-D network topology, which is a neighborhood-preserving representation of real collaborative networks, using the method of spectral embedding.

# Agent pool optimization: Putting the models in use

❑ **Environment**

   o A problem management function of an IT service provider, consisting of 5000 agents with varying expertise

❑ **Business Problem**

   o Constantly restructure agent pools to accommodate the evolving workload and human resource availability, while maintaining efficiency in problem resolution

❑ **Trade-offs**

   o Dividing agents into smaller pools

   o Dividing agents into larger pools

   o Correct initial dispatching (with probability $p$) is also critical to improving routing efficiency

# How to effectively organize workforces?

❑Problem Definition
- o Each task is associated with one difficulty value $x$
- o Each worker $i$ is characterized with their ability $w\_i$
- o Each worker can solve the task if and only if $x<=w\_i$
- o If a worker cannot solve a given task, the task must be forwarded to a worker of greater ability.

**Goal:**

Given a set of worker abilities $W$ and a distribution $P$ over task difficulty,  design efficient forwarding structure,  that minimizes both the maximum workload of any worker, and the number of workers that need to attempt a task.

# How to effectively organize workforces?

❑Results

- o Omniscient workload
  - ➢ Assign a task to any worker that is capable of solving it.
  - ➢ For any **W** and task difficulty distribution **P**, precisely characterize the optimal omniscient maximum workload **M**

- o Near-Optimal depth-workload tradeoffs
  - ➢ By hierarchically arranging the workers in balanced **b**-ary trees, we can simultaneously obtain a multiplicative factor of **b²** of **M** in terms of maximum workload, with a resulting depth of *log(n)/log(b)*.

# Collaborations In This Talk

❑ **Peer Interaction**
- o Patterns
- o Prediction
- o Optimization



❑ **Team Joint Collaboration**
- o Patterns
- o Prediction
- o Optimization

# Collaborative Teams

❑ Work together to produce music, movies, games, and other cultural products.

   o Wikipedia
   o Open Source Software (OSS)

**Questions:**

❑ How do teams form?

❑ What factors will influence the success of teams?

# The Formation and Success of Online Creative Collaborations

❑ Existing theories on common bond, social identity, and social exchange [Kraut et al., 2007, Taifel et al, 1982, Emerson et al., 1976]

    o Communication richness

    o Shared interests

    o Status within the community

    o Balance of efforts

# The Formation and Success of Online Creative Collaborations

❏ Data
- o February Album Writing Month (FAWM)
  - 6116 users; 39,103 FAWM songs

❏ Goal: to predict whether a pair of users posted a collaborative song to the website or not.

❏ Method
- o Logistic regression

# The Formation and Success of Online Creative Collaborations

❑ Features:

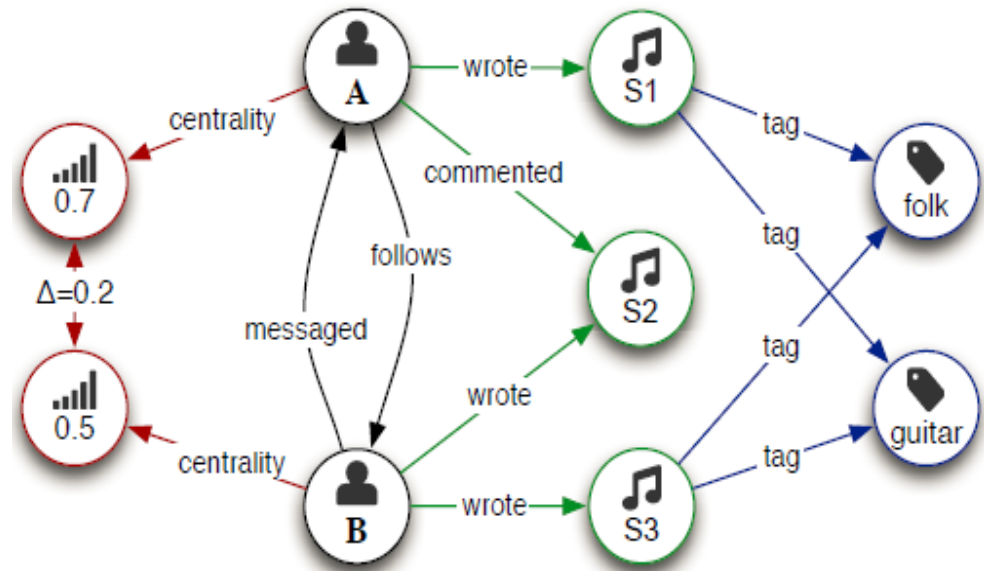Various Paths that can connect two users in FAWM network e.g., A-->follows--B, A<--messaged--B



Figure 1. A small example subset of the FAWM network.

# The Formation and Success of Online Creative Collaborations

❑Results

| Method | AUC | PP@1 |
|---|---|---|
| Path-based regression [31] | **0.990** | **0.344** |
| Baseline random walk [30] | 0.982 | 0.087 *** |
| Adamic & Adar [4] | 0.953 * | 0.005 *** |
| Matrix factorization (SVD) [28] | 0.865 ** | 0.133 ** |

*** $p < 0.001$    ** $p < 0.01$    * $p < 0.05$

Table 1. Evaluation of link-prediction methods, in terms of area under the ROC curve (AUC) and personalized precision at rank one (PP@1).

# The Formation and Success of Online Creative Collaborations

❑ Results

Key findings on team success:

1. Balance of effort improves satisfaction.

2. Higher-status partners may enjoy Collabs less.

3. Frequent communication helps (usually).

# The Formation and Success of Online Creative Collaborations

❑Results

Key findings on team formation:

1. Communication Exchanges predict collaboration, such as following a partner's song feed, direct messaging, and commenting on the partner's songs.

2. Collabs form out of shared interests, but different skills

3. Small status differences are positively associated with collabs. (Neither too different, nor exactly the same. similar to our finding in KDD'14)

# Success Factors in Online Creative Collaboration

❑ Collaborative animated movies (Collab)
   o Newgrounds: 2,200,000 registered users; 180,000 animated movies and games


❑ Success definition:
   o The completion of a Collab

# Success Factors in Online Creative Collaboration

❑Three factor types:
- o Planning & Structure
  - ➢ Collabs with initial planning and structure, especially technical specifications, are more likely to be successful.

- o Reputation & Experience
  - ➢ Animators who are well-known in the community are more likely to lead successful collabs.
  - ➢ Animators who have experience with Flash and past collabs are more likely to lead successful collabs.

- o Communication & Dedication
  - ➢ Collabs whose members frequently communicate are more likely to be successful
  - ➢ Collabs whose leaders frequently communicate are more likely to be successful.

# Success Factors in Online Creative Collaboration

❑Experiments

| Variable | Completed Collabs | | | Failed Collabs | | |
|---|---|---|---|---|---|---|
| | M | SD | N | M | SD | N |
| Theme** | 0.94 | 0.24 | 112 | 0.73 | 0.45 | 780 |
| Specs** | 3.74 | 1.70 | 112 | 2.03 | 2.02 | 780 |
| Authorship** | 0.29 | 0.53 | 112 | 0.11 | 0.41 | 780 |
| Restrictions* | 0.63 | 0.99 | 112 | 0.40 | 0.81 | 780 |
| Gatekeeping* | 0.39 | 0.61 | 112 | 0.25 | 0.54 | 780 |
| Communication** | 0.83 | 0.90 | 112 | 0.52 | 0.79 | 780 |

Table 1. Comparison of first post attributes for completed vs. failed collabs, * $p < 0.05$, ** $p < 0.01$.

Each row stands for one category of planning/structural elements in the first posts of Collab threads. Completed Collabs tend to have more planning and structure than failed ones.

# Part III: Sum-Up

**Collaborations:**

❑Peer Interaction

❑Team Joint Collaboration

**Questions:**

- How do experts make routing decisions?
- Who have made inefficient routing decisions?
- How to optimize the routing performance through targeted training?
- Can the completion time of a task be predicted so that one can act early for difficult tasks?
- Can we predict next best experts?
- How to model the collaborative networks?

- How do teams form?
- What factors will influence the success of teams?

# References

- Analyzing Expert Behaviors in Collaborative Networks. [Sun et al., KDD'14]

- Efficient Ticket Routing by Resolution Sequence Mining. [Shao et al., KDD'08]

- Generative Models for Ticket Resolution in Expert Networks. [Miao et al., KDD'10]

- Understanding Task-Driven Information Flow in Collaborative Networks. [Miao et al., WWW'12]

- Depth-Workload Tradeoffs for Workforce Organization. [Heidari et al., AAMAS'12]

- Why It Works (When It Works): Success Factors in Online Creative Collaboration. [Luther et al., GROUP'10]

- The Formation and Success of Online Creative Collaborations [Settles et al., CHI'13]

# Part IV
# Relationship Mining

# Relationships embedded in network structure

- Network links represent general relationships of longer term
  - Friendship of mutual consent, e.g., Facebook.
  - Follow link, e.g., Twitter.
  - Collaboration, e.g., DBLP.
  - Negative relationships, e.g., distrust (Epinion), foe (slashdot).

- Can we discover more specific and implicit relationships from these links?

# Example: Role Discovery



Information network without role/relationship info, e.g. a company's email network

**Latent relationship graph**

How to infer

**CEO**

**Manage**r

**Employee**

Courtesy of Chi Wang for this slide

# Twitter is a unique social platform

1. The "follow" links are established without mutual consent.
   - An explosion of social links
     - Everyone has a large number of followers and followees.
     - A huge number of tweets are generated everyday. (~300 million tweets/daily)
   - A further shrinkage of the network diameter
     - Information diffusion is much faster
2. It is a mixture of social network and news media
   - H.Kwak et al  WWW 2010

# Follow network = real-life social network?

## How much of this follow network reflects a user's real-life offline social network?

- Mutual follow links do not necessarily indicate real-life interaction.

- The number of followees and followers varies significantly.

**Problem:** Given a Twitter follow network of a target user, identify the user's offline community by examining the follow linkage alone.

# Motivation

1. More accurate and robust user interest modeling

2. Online vs Offline relationship understanding

3. User identity alignment across different platforms

4. Spam, Zombie account detection

5. Business competitive analysis

And more ...

# Principle I: Mutual Reachability

**Information should be able to flow in both directions within a small distance between real-life friends.**

# Principle II: Friendship Retainability

**The size of a user's offline community has an upper-bound threshold σ related to Dunbar's number**

# Principle III: Community Affinity

**A user's off-line friends usually group into clusters such that within each cluster members know each other**

# Approach

- Random Walk
- Closeness Score
- Iterative Off-line Community Discovery

# A Case Study

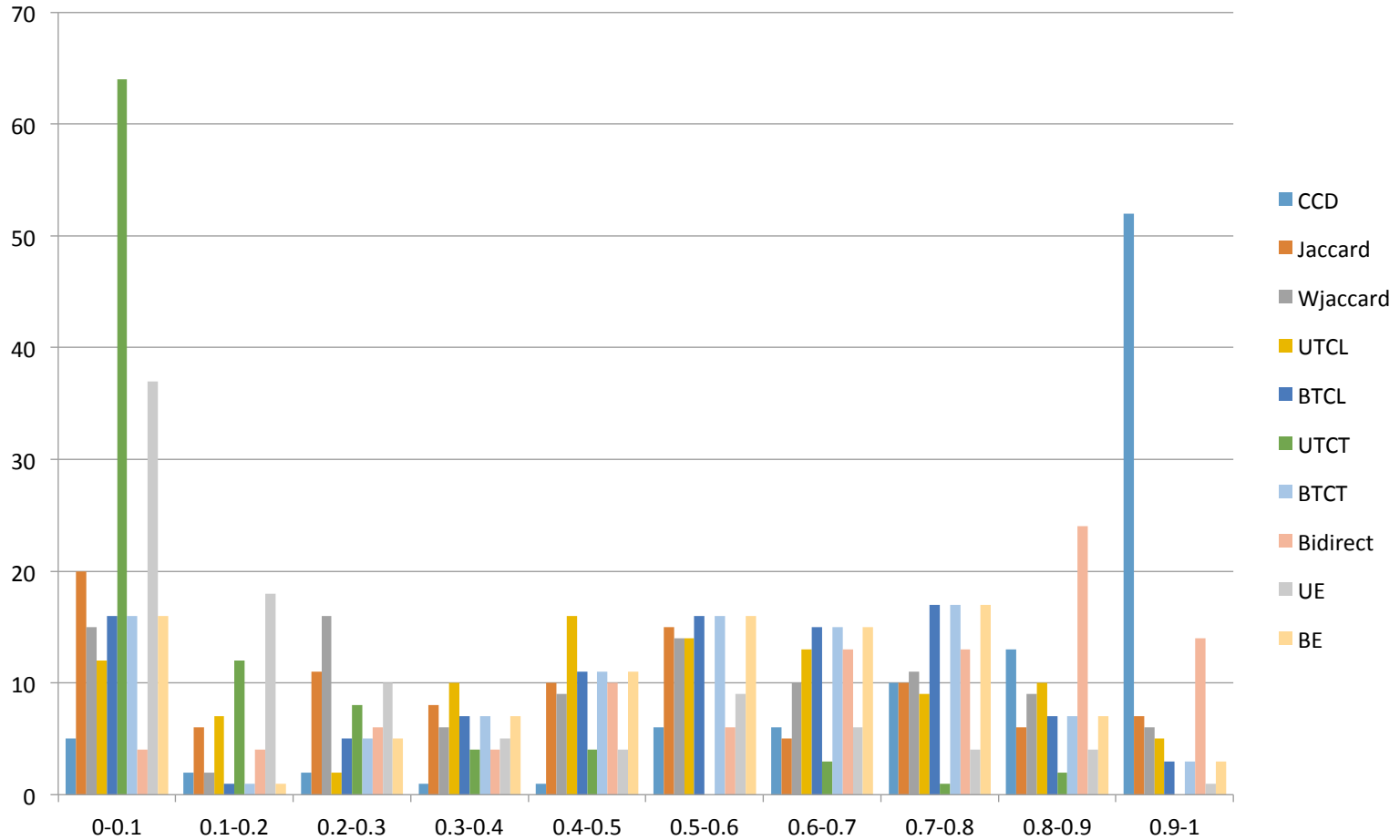**A real Twitter user: following 385 users; followed by 107 users**

# Precision

**Our method, Core Community Discovery (CCD), performs the best.**

# Recall

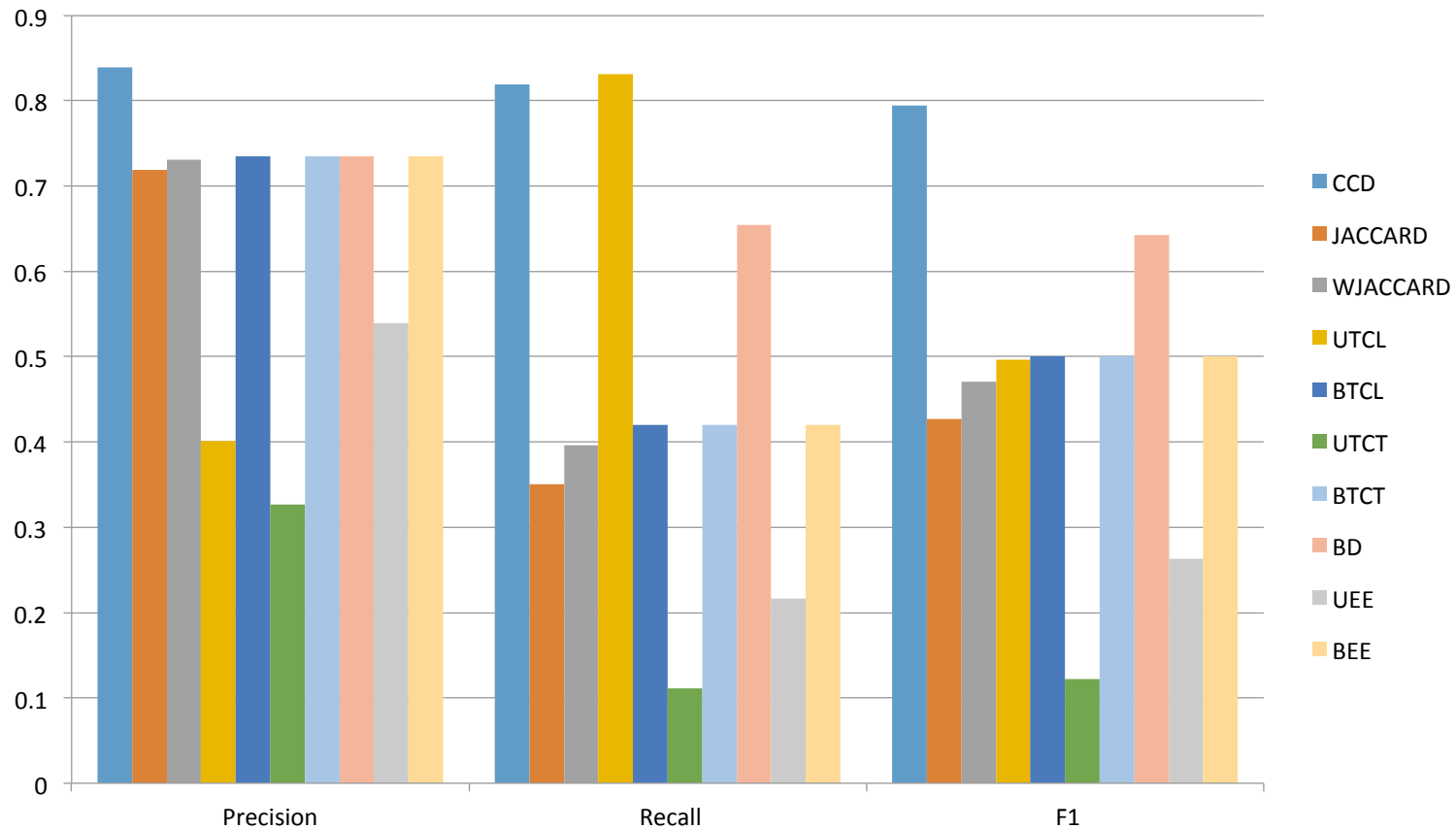**Our method, Core Community Discovery (CCD), performs the best.**

# F-score

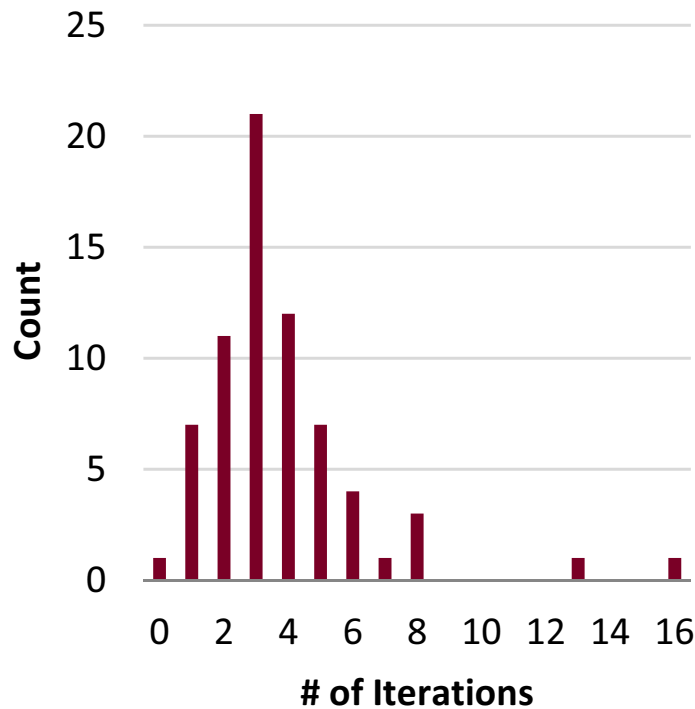**Our method, Core Community Discovery (CCD), performs the best.**

# ALL

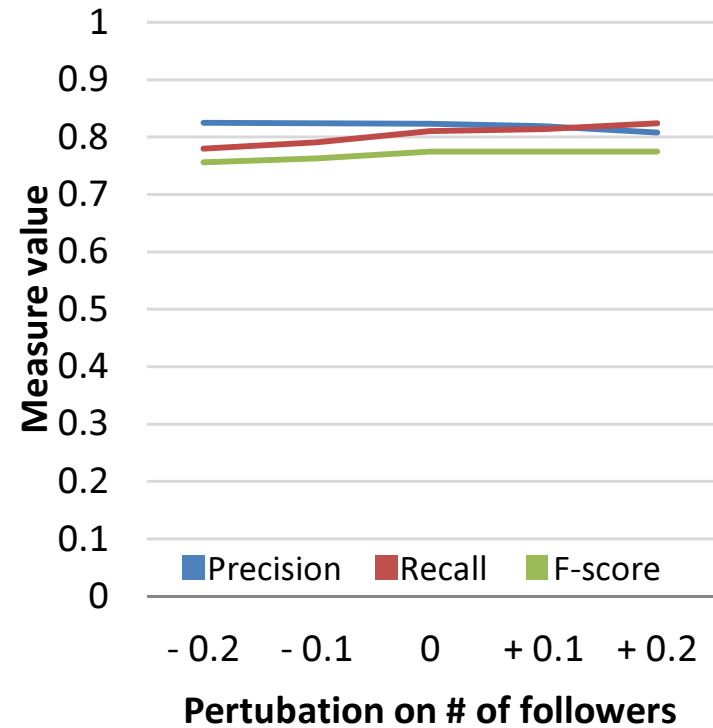**Our method, Core Community Discovery (CCD), performs the best.**

# Parameters

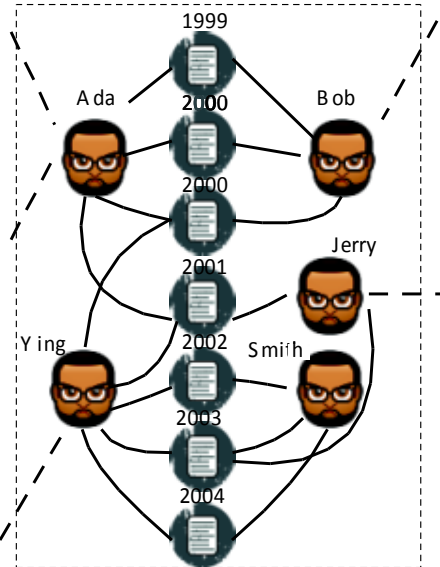■ **On # of Iterations**

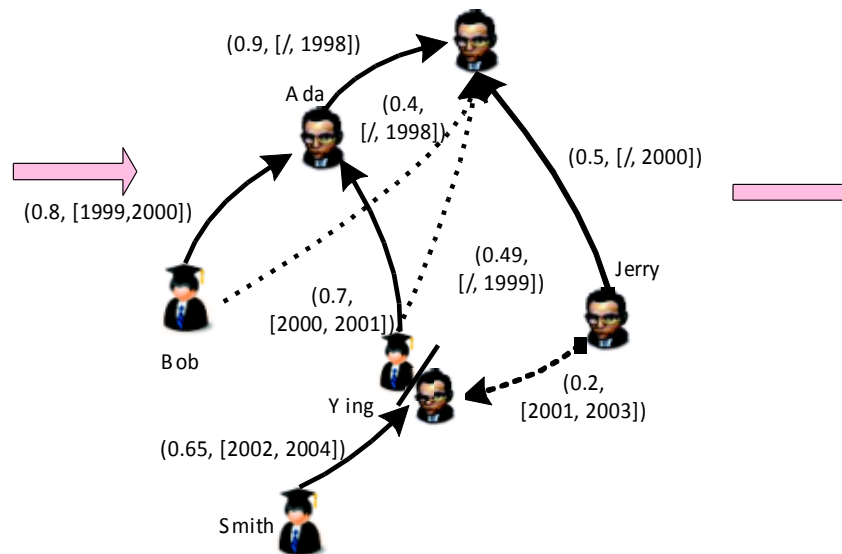■ **On Robustness**

# Mining Advisor-advisee Relationship
## [Wang et al, KDD'10]

- Input: research publication network.
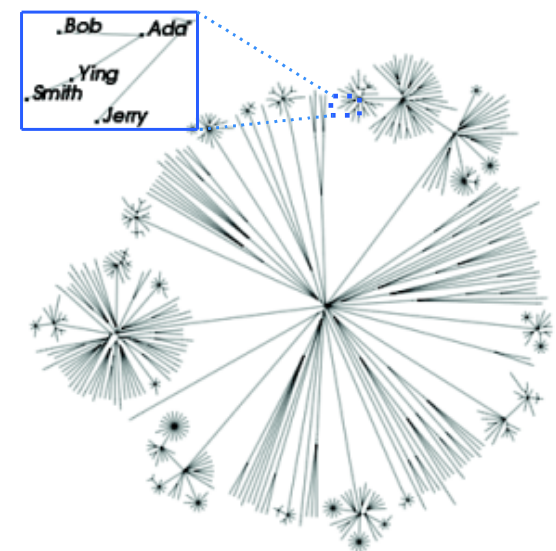- Output: potential advising relationship and their ranking – (r, [st, ed])
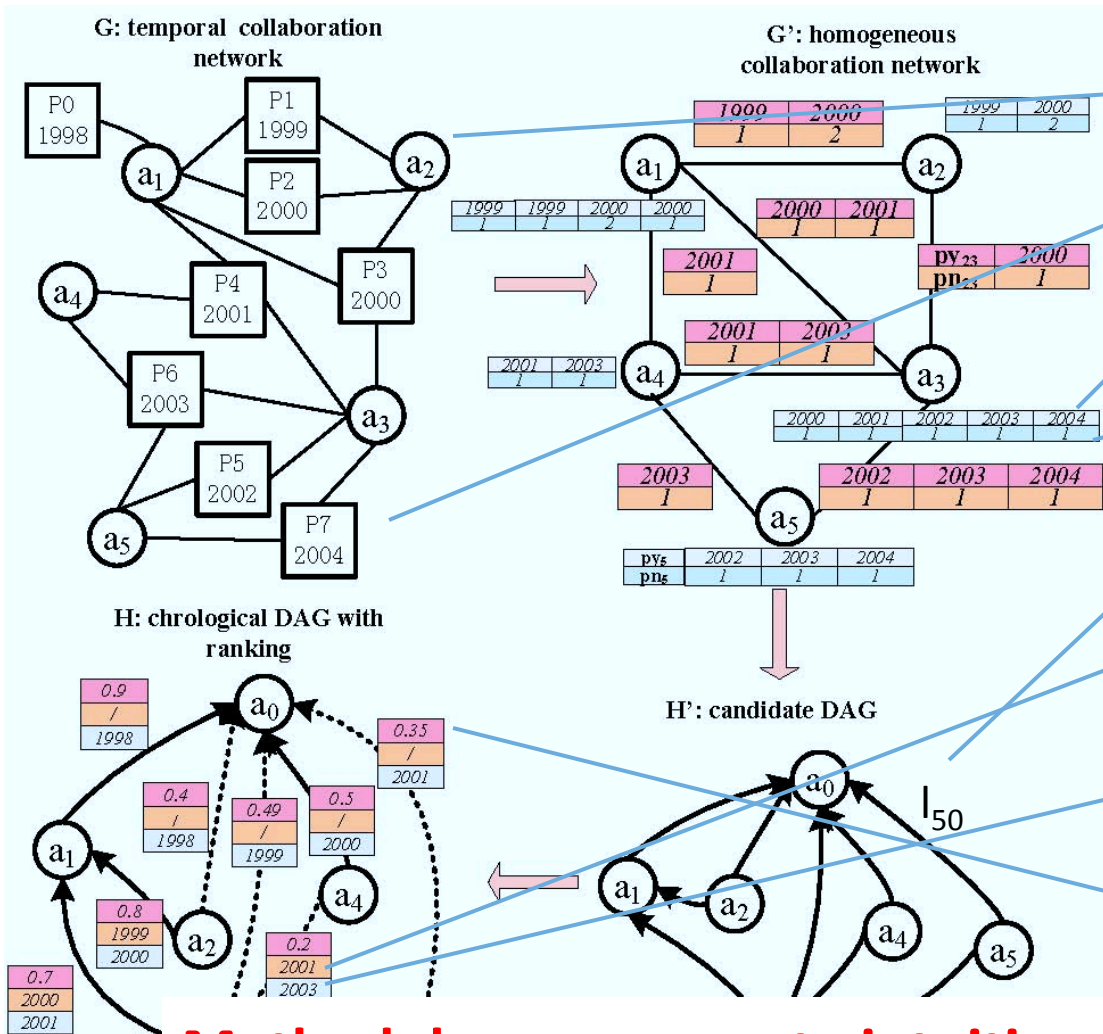


Input: Temporal collaboration network

Output: Relationship analysis

Visualized chorological hierarchies

Courtesy of Chi Wang for this slide

# Overall Framework



- $a_i$: author i
- $p_j$: paper j
- py: paper year
- pn: paper#
- $l_{ij}$: local feature
- $st_{i,yi}$: start time
- $ed_{i,yi}$: end time
- $r_{i,yi}$: ranking score

**Methodology: propagate intuitive rules and constraints over the whole network.**

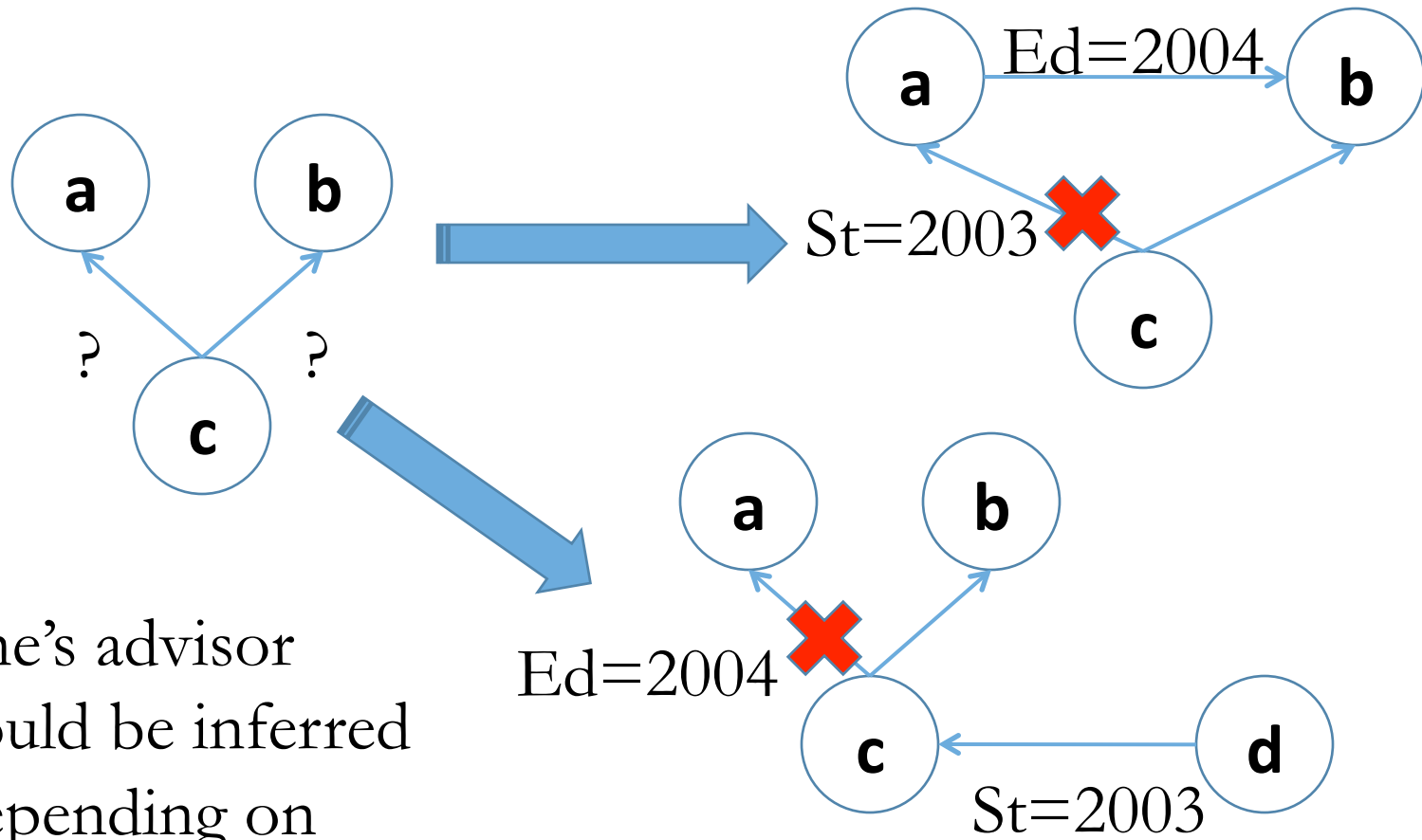18  Courtesy of Chi Wang for this slide

# Basic Constraints

- If $a_y$ advises $a_x$ since the year $st_x$
  - $a_y$ can only advise $a_x$ after it graduated
    - ➔ $ed_y < st_x < ed_x$
  - $a_y$ must have a longer history of publication than $a_x$ before $st_x$.
    - ➔ The candidate graph H' is a DAG.

The model can incorporate other intuitions as factor functions into a Time-constrained Probabilistic Factor Graph (TPFG)

Courtesy of Chi Wang for this slide

# Why is network structure helpful?

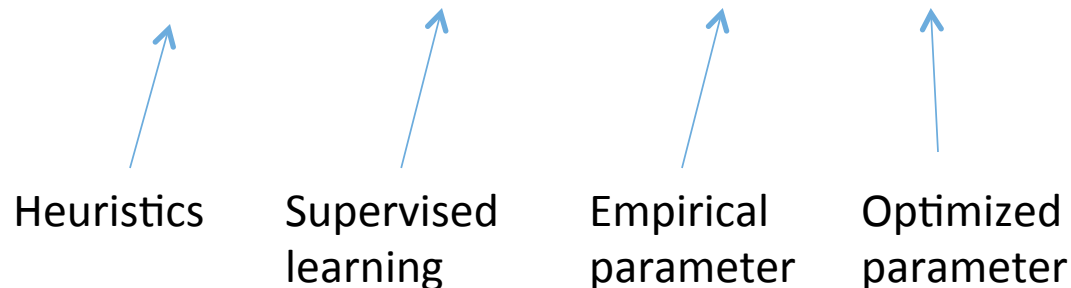- More than pairwise features: interdependency



one's advisor could be inferred depending on others' advisor!
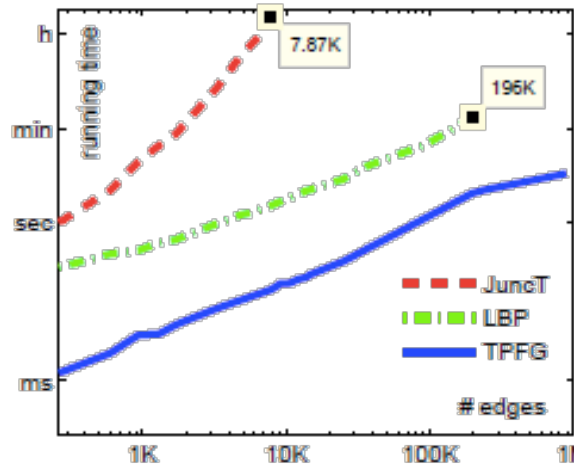
# Experiment on DBLP

- DBLP data: 654, 628 authors, 1076,946 publications, publishing time provided.

- Labeled data: MathGenealogy Project; AI Gealogy Project; Faculty Homepage

| Datasets | RULE | SVM | TPFG | |
|----------|------|------|------|------|
| TEST1 | 69.9% | 73.4% | 80.2% | **84.4%** |
| TEST2 | 69.8% | 74.6% | 81.5% | **84.3%** |
| TEST3 | 80.6% | 86.7% | 88.8% | **91.3%** |

Heuristics    Supervised learning    Empirical parameter    Optimized parameter

Courtesy of Chi Wang for this slide

# Case Study & Scalability

| Advisee | Top Ranked Advisor | Time | Note |
|---------|-------------------|------|------|
| David M. Blei | 1. Michael I. Jordan | 01-03 | PhD advisor, 2004 grad |
| | 2. John D. Lafferty | 05-06 | Postdoc, 2006 |
| Hong Cheng | 1. Qiang Yang | 02-03 | MS advisor, 2003 |
| | 2. Jiawei Han | 04-08 | PhD advisor, 2008 |
| Sergey Brin | 1. Rajeev Motawani | 97-98 | "Unofficial advisor" |



(a) Time        (b) Space

Courtesy of Chi Wang for this slide

# Application: Expert Finding

**An example on a real
system: Arnetminer**

**Performance
improvement**

Courtesy of Chi Wang for this slide

# Representative work

- "Relation Mining" [Kadri 03, Rinaldi 06, Coppola 08]
  - Mainly text mining and language processing on text data and structured data.

- "Relational Learning" [Getoor 07, Tang 09]
  - The classification when objects and entities are presented in multiple relations

- Relationship with semantic meaning
  - [Diehl 07]: a supervised approach

- Positive-Negative relationship
  - Leskovec et al. [ WWW'10]

- Social circle in ego network
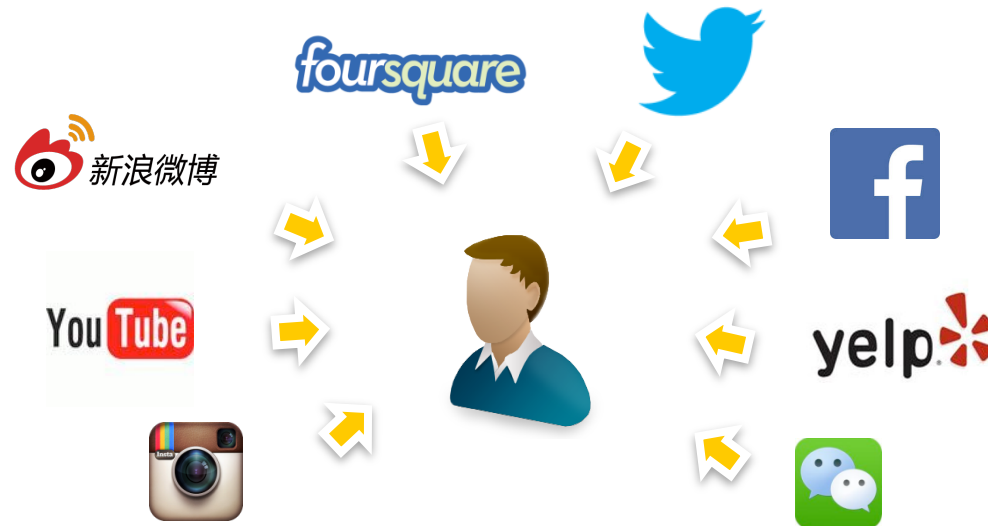  - McAULEY et al. [TKDD'14]

# References

- W. Xie, C. Li, F. Zhu, E. Lim and X. Gong, ***"When a Friend In Twitter is a Friend In Life"***, ACM Int. Conf. on Web Science 2012.
- C. Wang, J. Han, Y. Jia, J. Tang, D. Zhang, Y. Yu and J. Guo , ***"Mining advisor-advisee relationships from research publication networks"***. KDD 2010.
- J. Leskovec, D. Huttenlocher and J. Kleinberg, ***"Predicting Positive and Negative Links in Online Social Networks"***, WWW 2010.
- J. Mcauley and J. Leskovec, ***"Discovering social circles in ego networks"***, TKDD, 2014.

# Part V
## User Identity Linkage

# User Identity Linkage

Link up all the accounts of the **same** user across **different** social platforms

# Why do we care?

- **Completeness**
  - Cross-platform user linkage would enrich an otherwise fragmented user profile to enable an all-around understanding of a user's interests and behavior patterns.

- **Consistency**
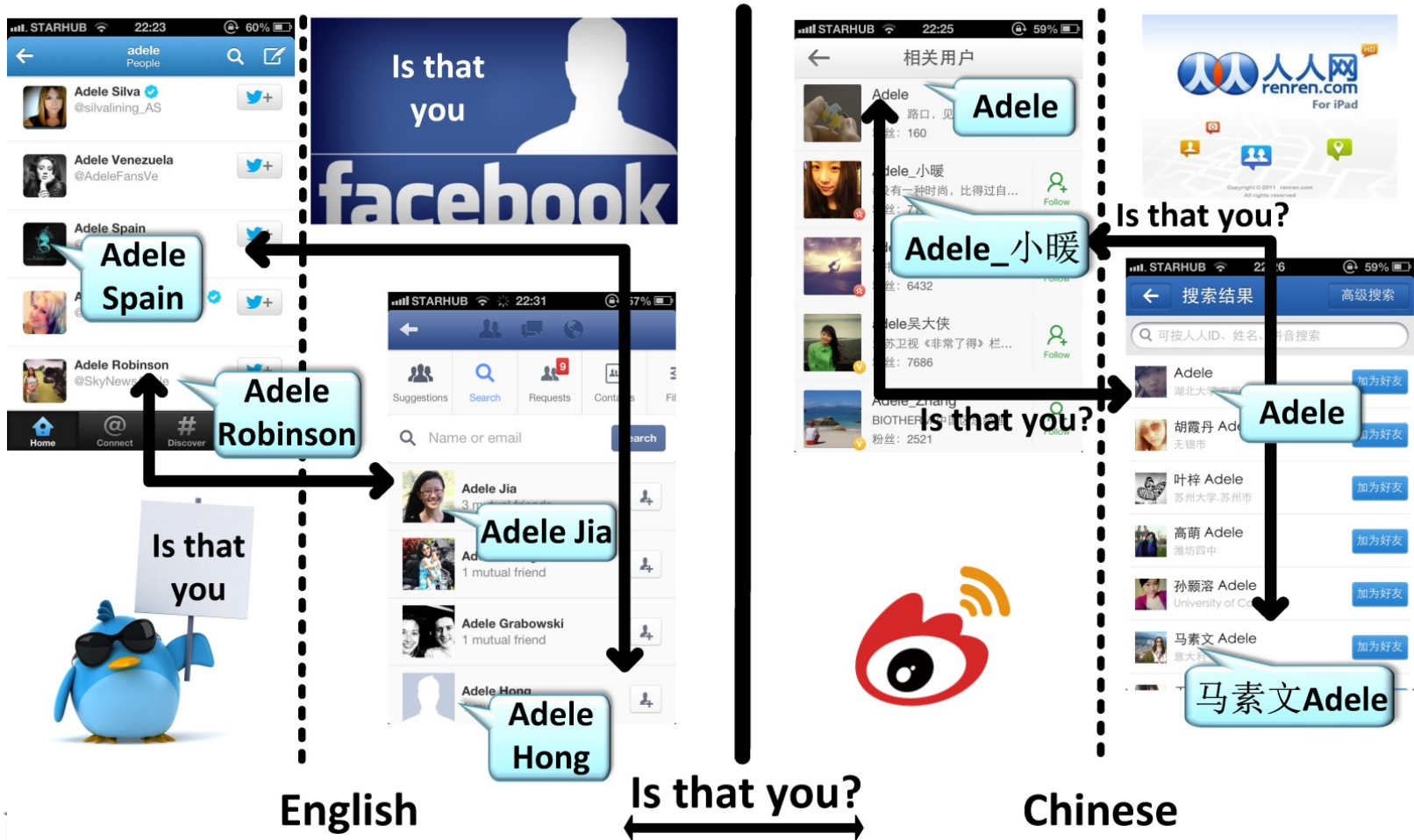  - Cross-checking among multiple platforms helps improve the consistency of user information.

- **Continuity**
  - User identity linkage makes it possible to integrate useful user information from those platforms that have over time become less popular or even abandoned.

# Example: A Not So Easy Case



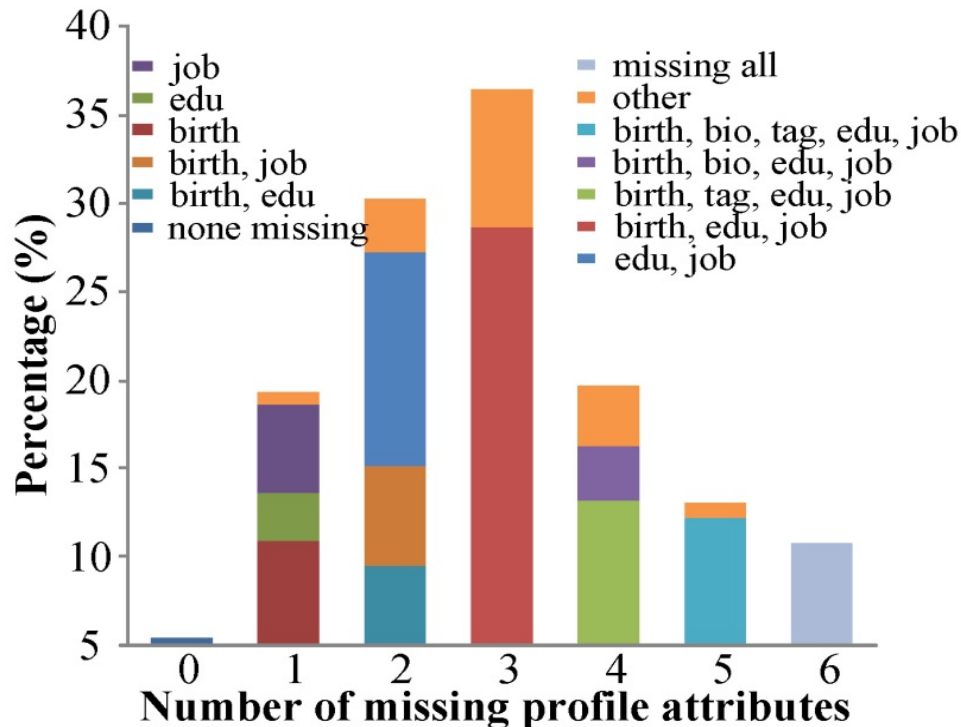Search name *Adele* in different social media platforms

English

Is that you?

Chinese

# Research Challenges

- **Usernames are not always reliable**
    - Traditional approaches that heavily rely on username parsing to link users may fail on more diversified communities.
    - Statistical models (e.g. SVM) or rule based models constructed with mere username and attribute analysis are far from being robust to accurately identify user linkage across online social communities.
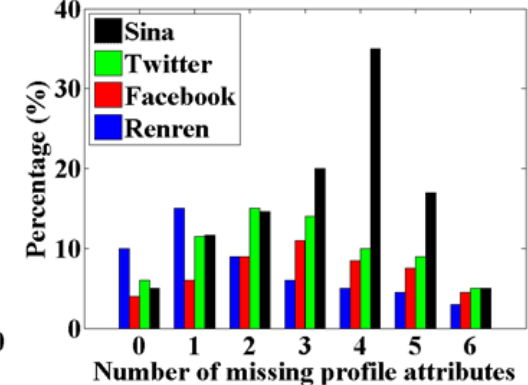
# Research Challenges
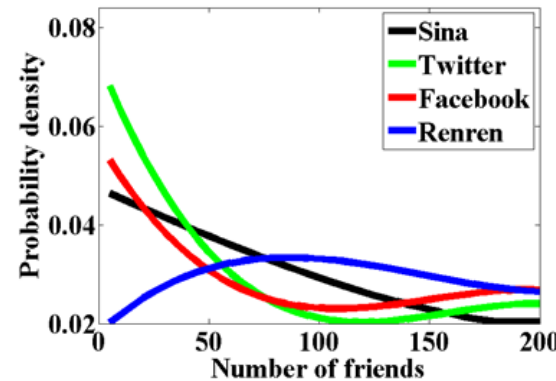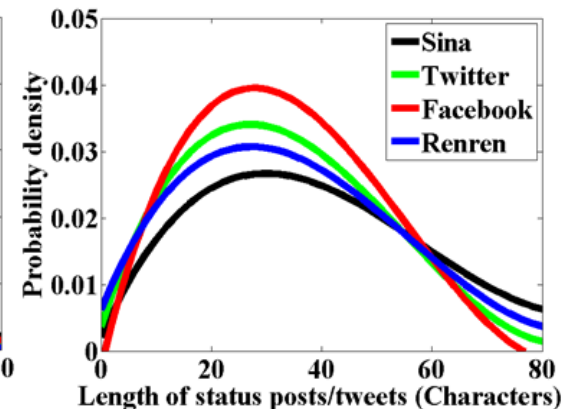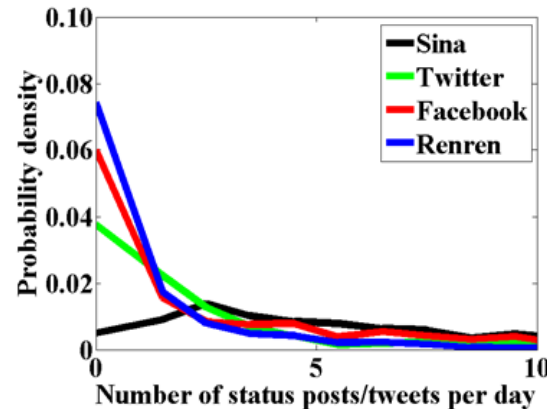
- **Missing Information**
  - At least 80% of users are missing at least 2 profile attributes out of the 6 most popular ones, and merely 5% of users have all attributes filled up.

# Research Challenges

- **Data Misalignment**
  - Platform Difference.
  - Heterogeneous behavior:
    - The user behavior can be represented by various types of media, e.g., locations, blogs, tweets, videos and images
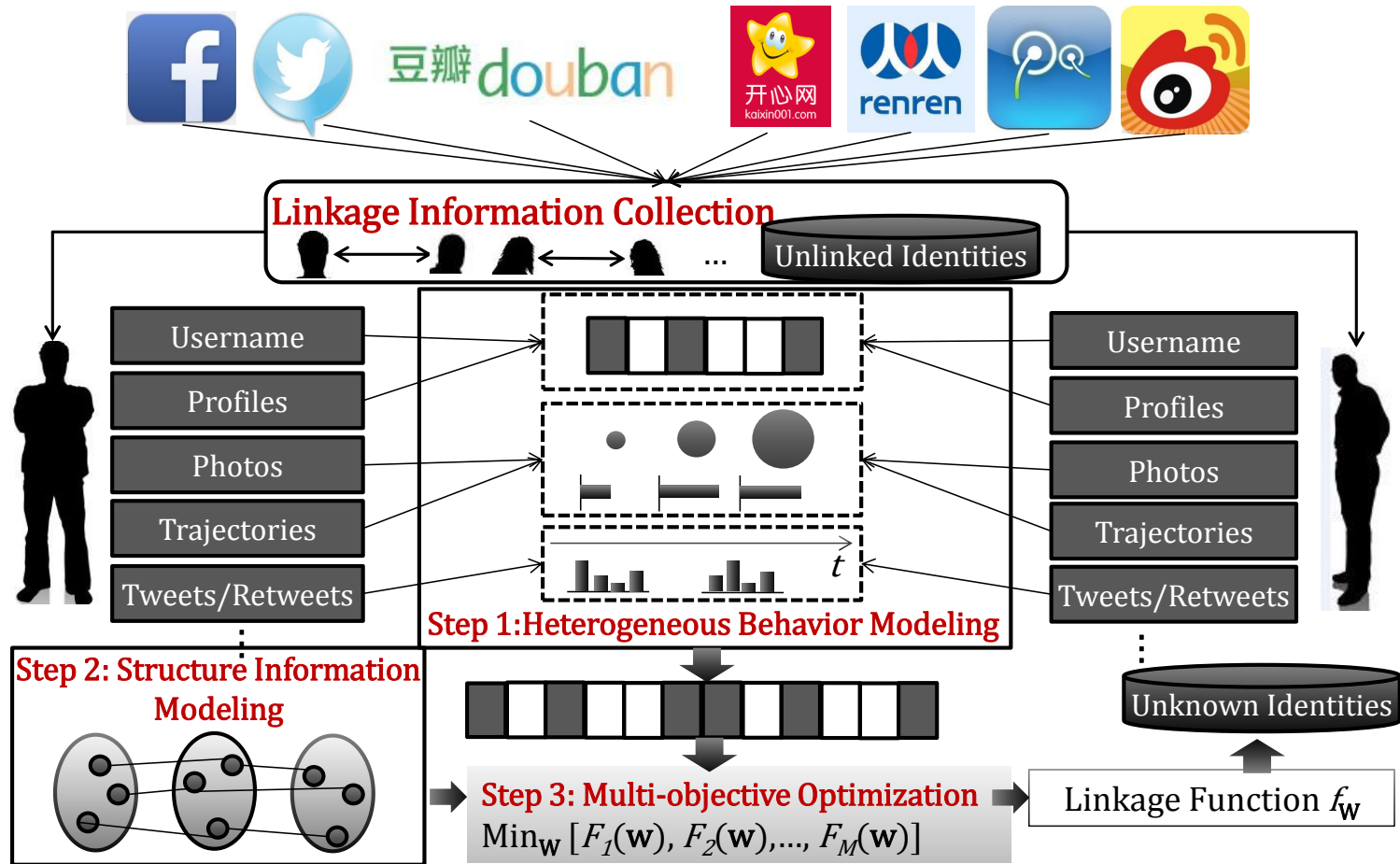  - Behavior Asynchrony.
  - Data Imbalance.

# Where are the hopes in social data?

**Two important features unique to social data:**

- User behavior trajectory along temporal dimension
  - Over a sufficiently long period of time, a user's social behavior exhibits a high level of consistency across different platforms.

- User's core social network structure
  - A user's core social network structures across different platforms share great similarity and offer a highly discriminative characterization of the user.

# HYDRA: a user identity linkage framework
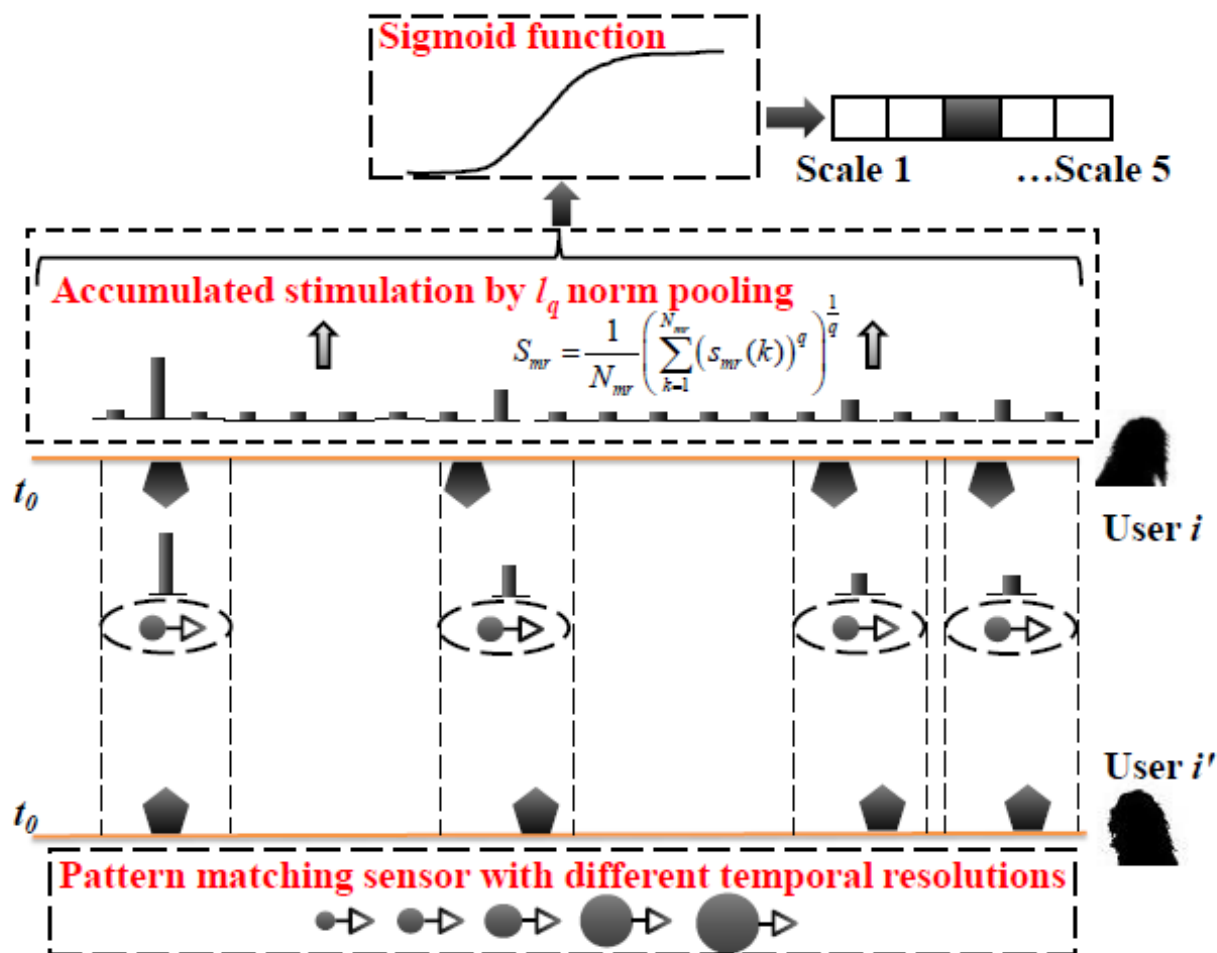## [Liu et al. SIGMOD'14]

# User Online Identity Data

- User attributes (numeric, categorical)
  - Demographics, location, personal interest, etc.

- User Generated Content (topics, sentiments)
  - Reviews, tweets, ratings, multimedia, etc.

- User Core Social network (snapshot/static view)
  - Friend network, followers/followees network, communities/interest groups, etc.

- User Behavior trajectory (dynamic, evolutionary)
  - content sharing history, social interaction pattern, network formation, etc.

# Main Stages

- Behavior Similarity Modeling
  - Calculate the multi-dimensional similarity vector between two users of a pair for all user pairs via heterogeneous behavior modeling.

- Structure Information Modeling
  - Construct the structure consistency graph on user pairs by considering both the core network structure of the users and their behavior similarities.

- Multi-objective Optimization with Missing Information
  - A two-class classification model via optimizing two kinds of objective functions simultaneously.
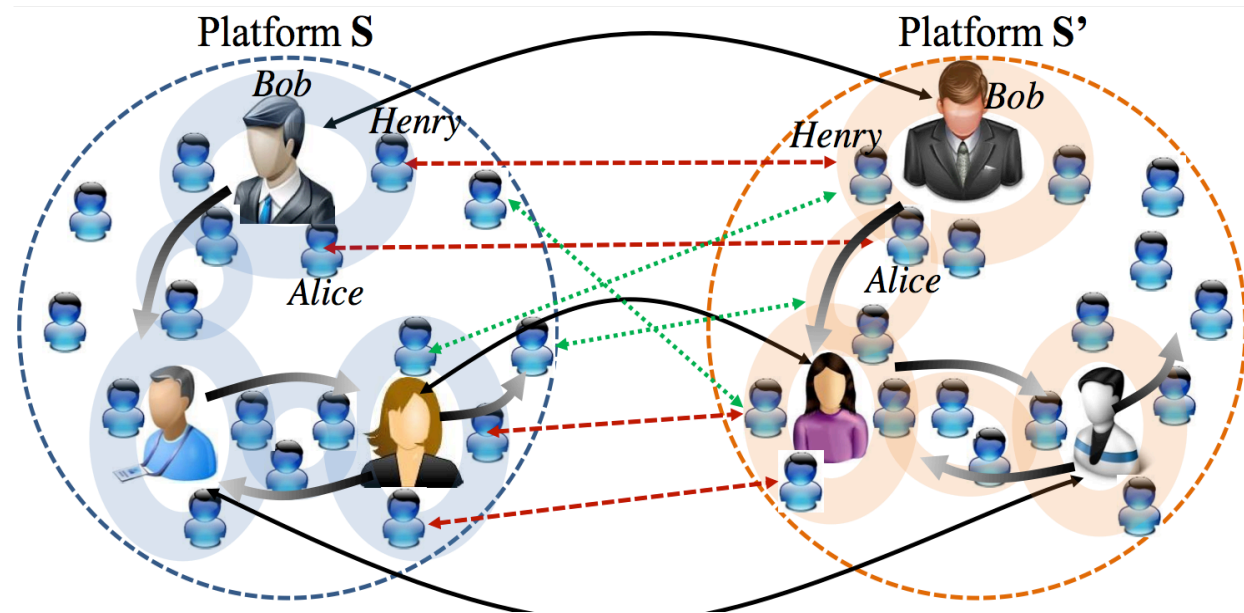
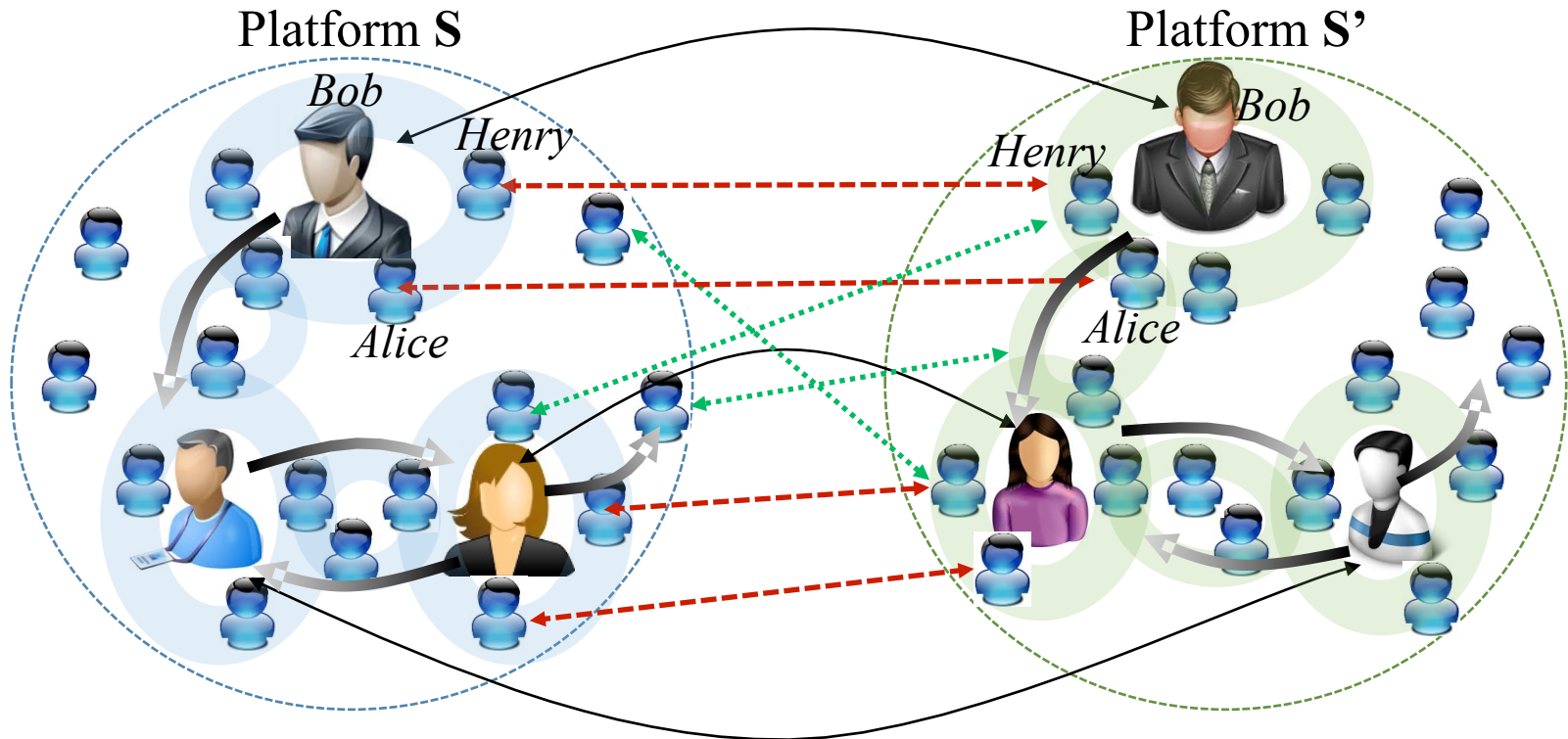# Multi-resolution Behavior Modeling

# Core Social Network Modeling

- People's closest friends are similar across different social platforms.

- Behavior similarity aggregation of the most frequently interacting friends of users provides insights into user identity linkage.

# A Structure Consistency modeling framework



Platform **S**

Platform **S'**

Bob

Henry

Alice

Henry

Bob

Alice

Black arrows: the ground-truth linkage information.
Red arrows: the correct linkage.
Green arrows: the falsely linked persons.

14

# Multi-objective Optimization Framework

- Supervised Learning

- Structure Consistency Modeling

- Multi-objective Optimization

A two-class classification problem --- construct multi-objective optimization which jointly optimizes the **prediction accuracy** on the labeled user pairs and **multiple structure consistency measurements** across different platforms.

# Multi-objective Optimization Framework

- ## Decision Model on Pairwise Similarity
  Support vector machine:

  $$f(x) = \mathbf{w}^T x + b$$

  $$F_D(\mathbf{w}) = \frac{\gamma_L}{2} \|\mathbf{w}\|^2 + \sum \xi_{ii'}$$
  $$s.t. \, y_{ii'}(\mathbf{w}^T x_{ii'} + b) \geq 1 - \xi_{ii'}$$

- ## High Order Structure Consistency

  An eigen-decomposition on the structure consistency graph

  $$\max_{\mathbf{w}} \mathbf{w}^T X^T \mathbf{M} X \mathbf{w}$$
  $$\Leftrightarrow \min_{\mathbf{w}} \mathbf{w}^T X^T (\mathbf{D} - \mathbf{M}) X \mathbf{w}$$
  $$\Leftrightarrow \min_{\mathbf{w}} \mathbf{w}^T X^T \Theta X \mathbf{w}$$
  $$s.t. \|\mathbf{w}\|^2 \leq s$$

- ## Multi-objective Optimization

  A generalized semi-supervised learning framework by optimizing the these two objective functions.

  $$\min_{\alpha, \beta} \left\{ \frac{1}{2}\alpha^T \left( 2\gamma_L \mathbf{K} + \frac{2\gamma_M}{|\mathbb{P}_l \cup \mathbb{P}_u|^2} \mathbf{K}(\mathbf{D} - \mathbf{M})\mathbf{K} \right) \alpha - \alpha^T \mathbf{K} \mathbf{J} \mathbf{Y} \beta + \beta^T \mathbf{1} \right\}$$
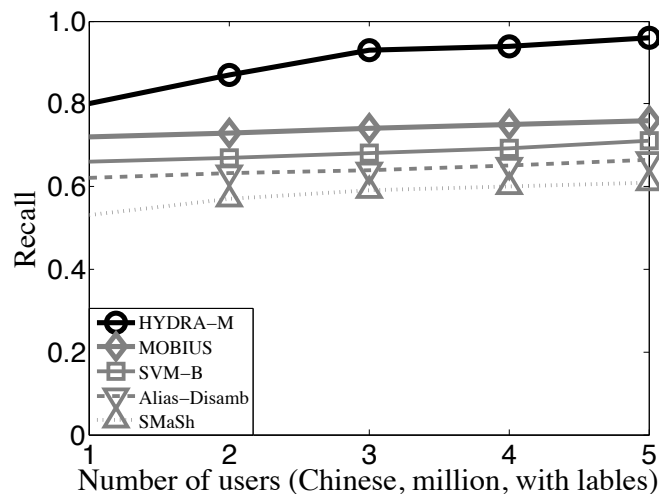
# Experiments

1. **Sina Weibo:** (www.weibo.com) A hybrid of Twitter and Facebook with a user base of 500 million users and 47 million daily active users by December 2012.
2. **Tecent Weibo:** (t.qq.com) Another twitter-like micro-blogging service with 500 million users and over 100 million daily active users.
3. **Renren:** (www.renren.com) A social network service dubbed as the Facebook of China with 162 million registered users.
4. **Douban:** (www.douban.com) A social network service for people to share content on topics of movies, books, music, and other off-line events in Chinese cities, with over 100 million monthly unique visitors.
5. **Kaixin:** (www.kaixin001.com) A social network service with 160 million registered users.
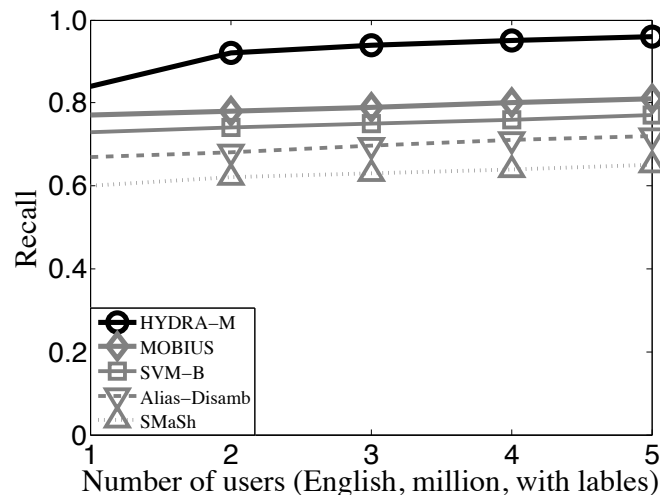
# Performance Evaluation

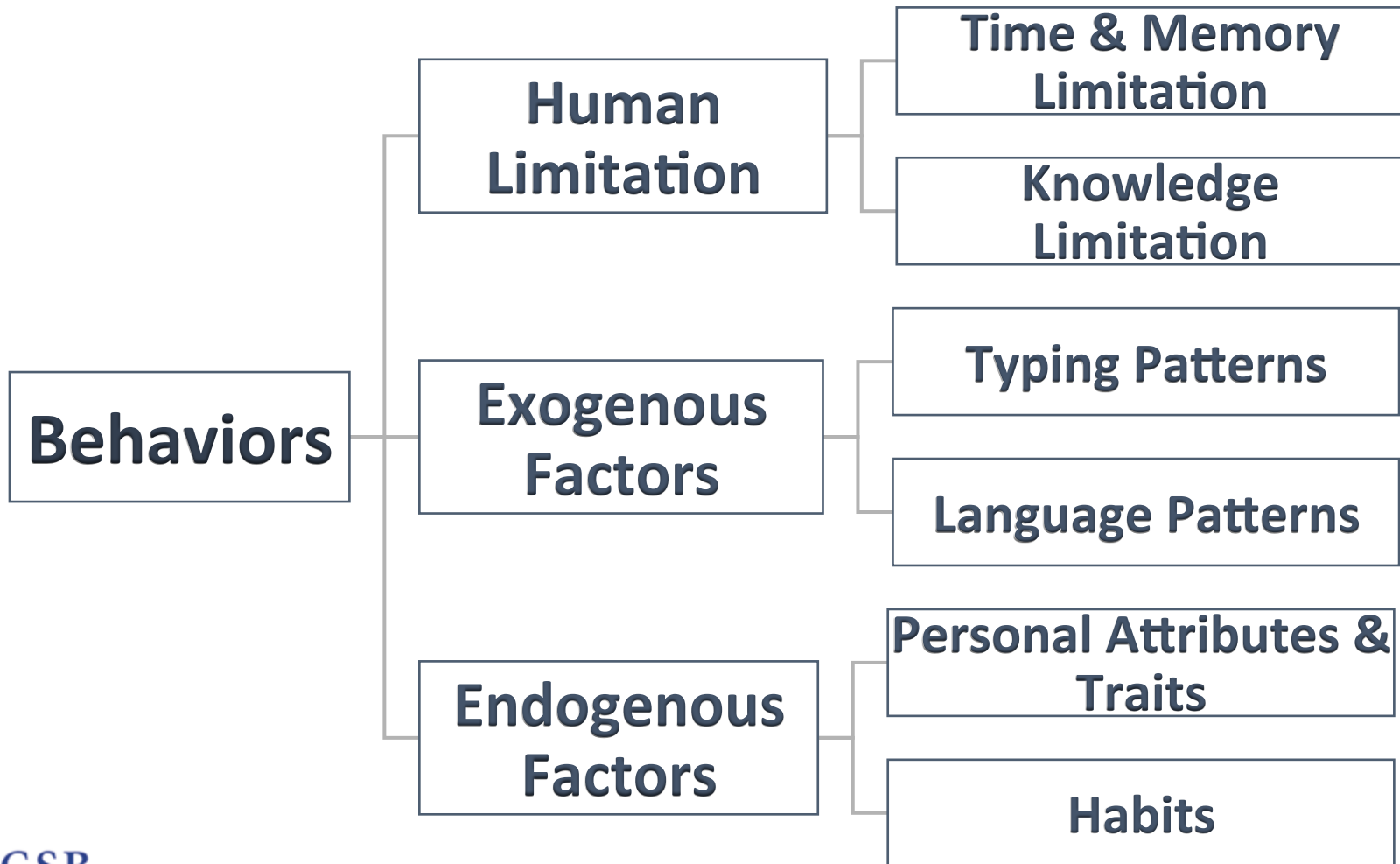

(a) Precision in *Chinese*

(b) Recall in *Chinese*

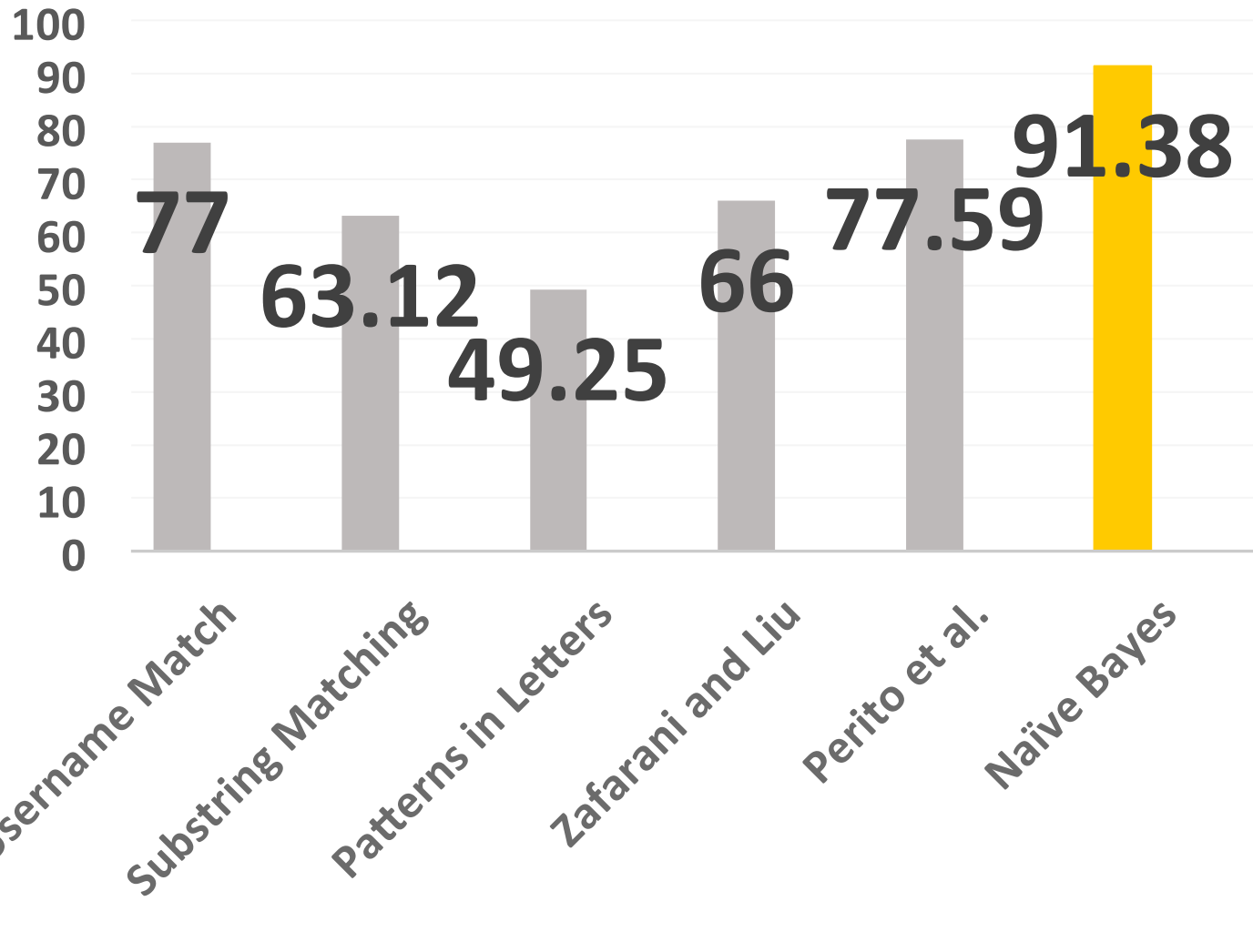(c) Precision in *English*

(d) Recall in *English*

# Connecting Users across Social Media Sites: A Behavioral-Modeling Approach
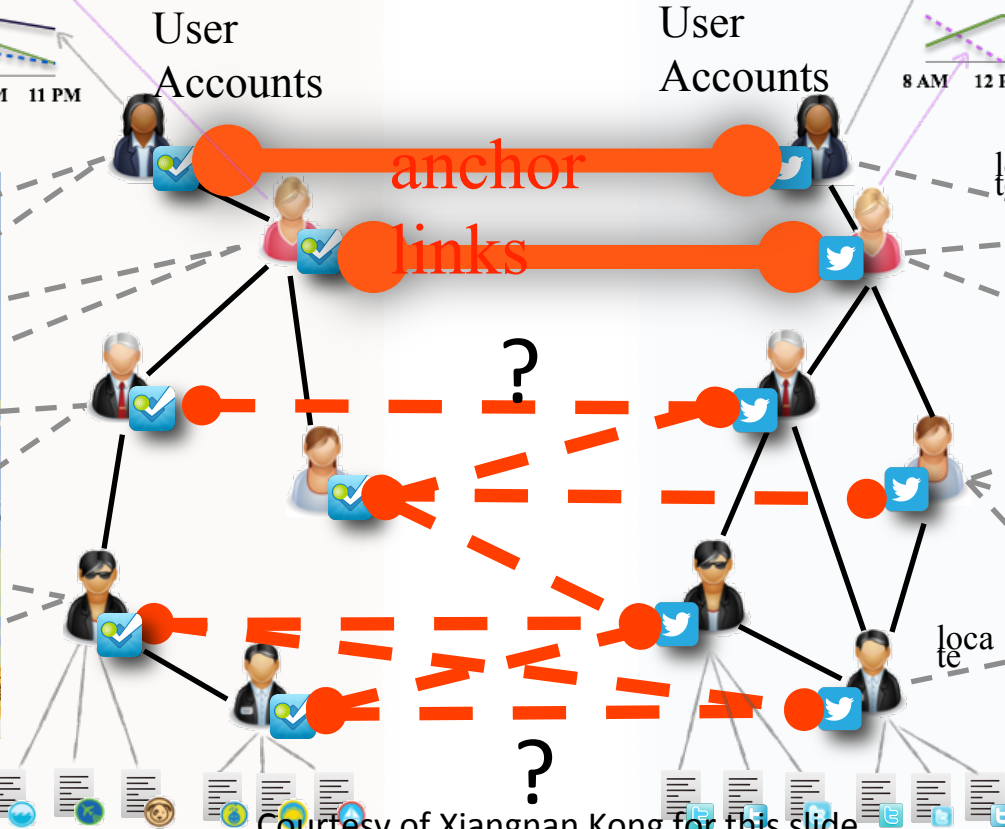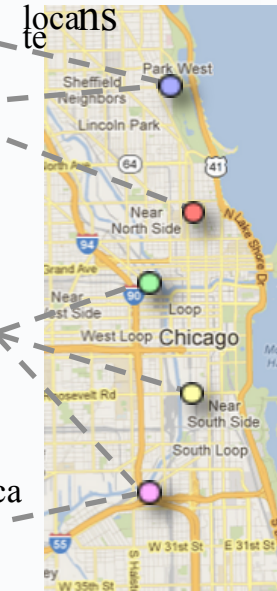[Zafarani et al. KDD'13]

# MOBIUS Performance

[Zafarani et al. KDD'13]

Courtesy of Reza Zafarani for this slide
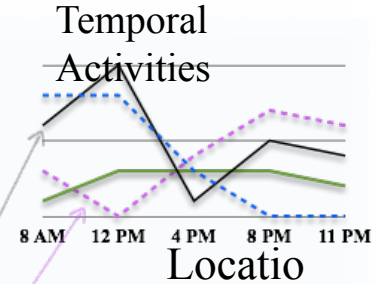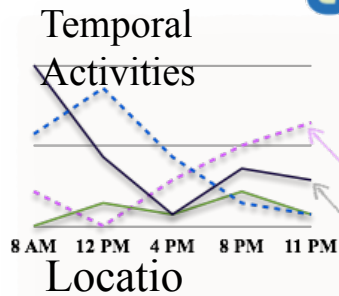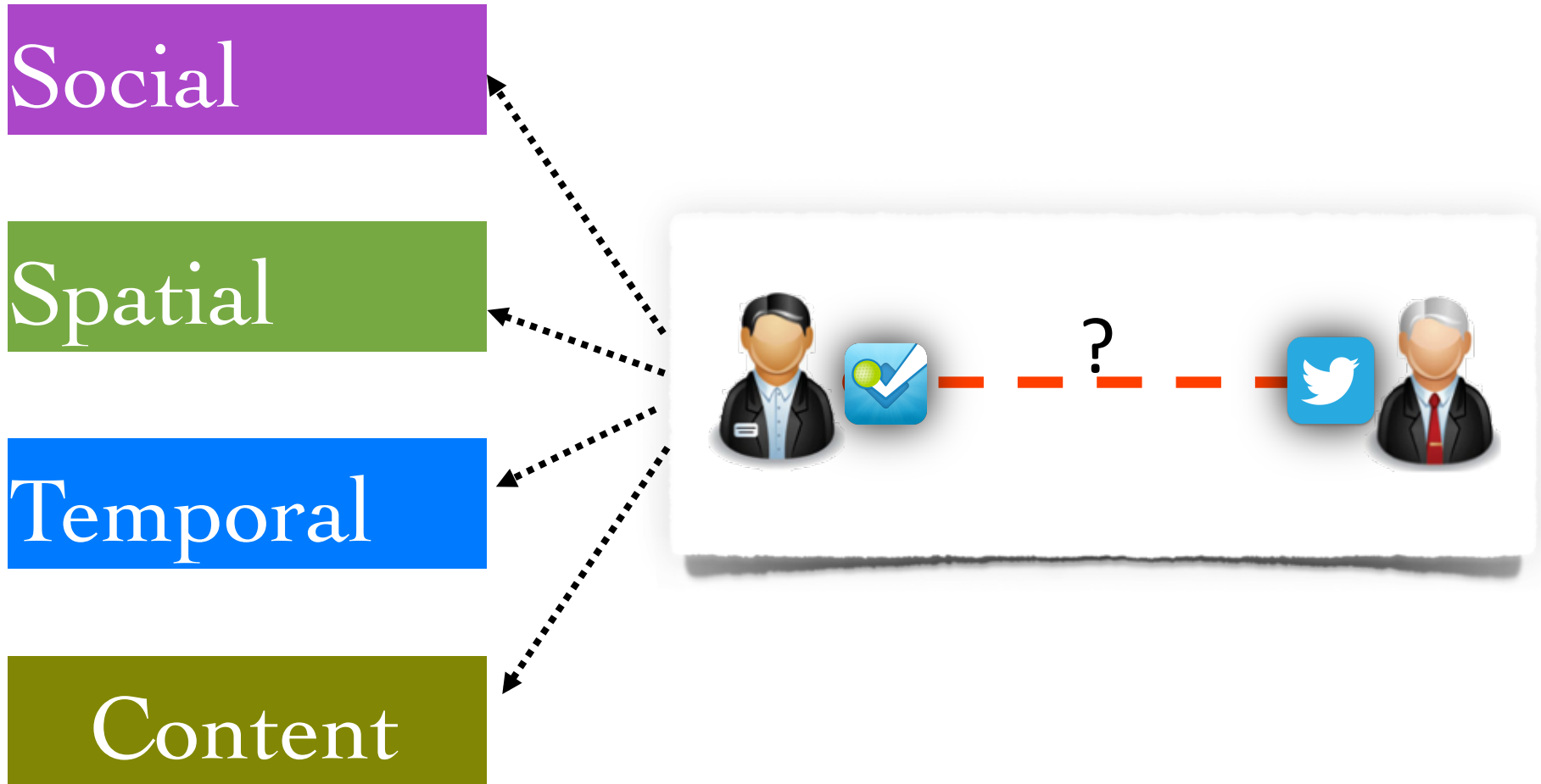
# Inferring Anchor Links across Multiple Heterogeneous Social Networks
## [Kong et al. CIKM'13]


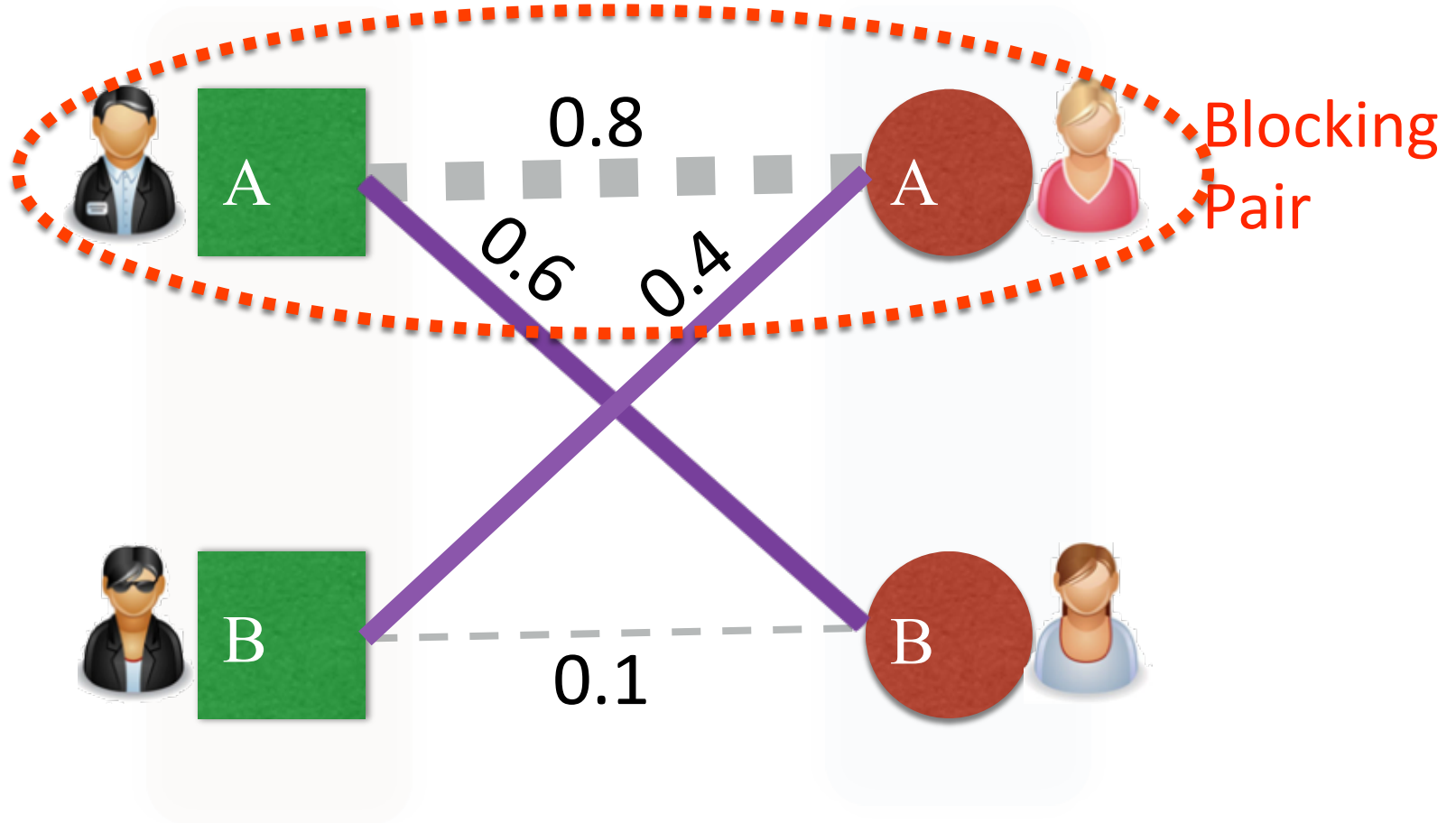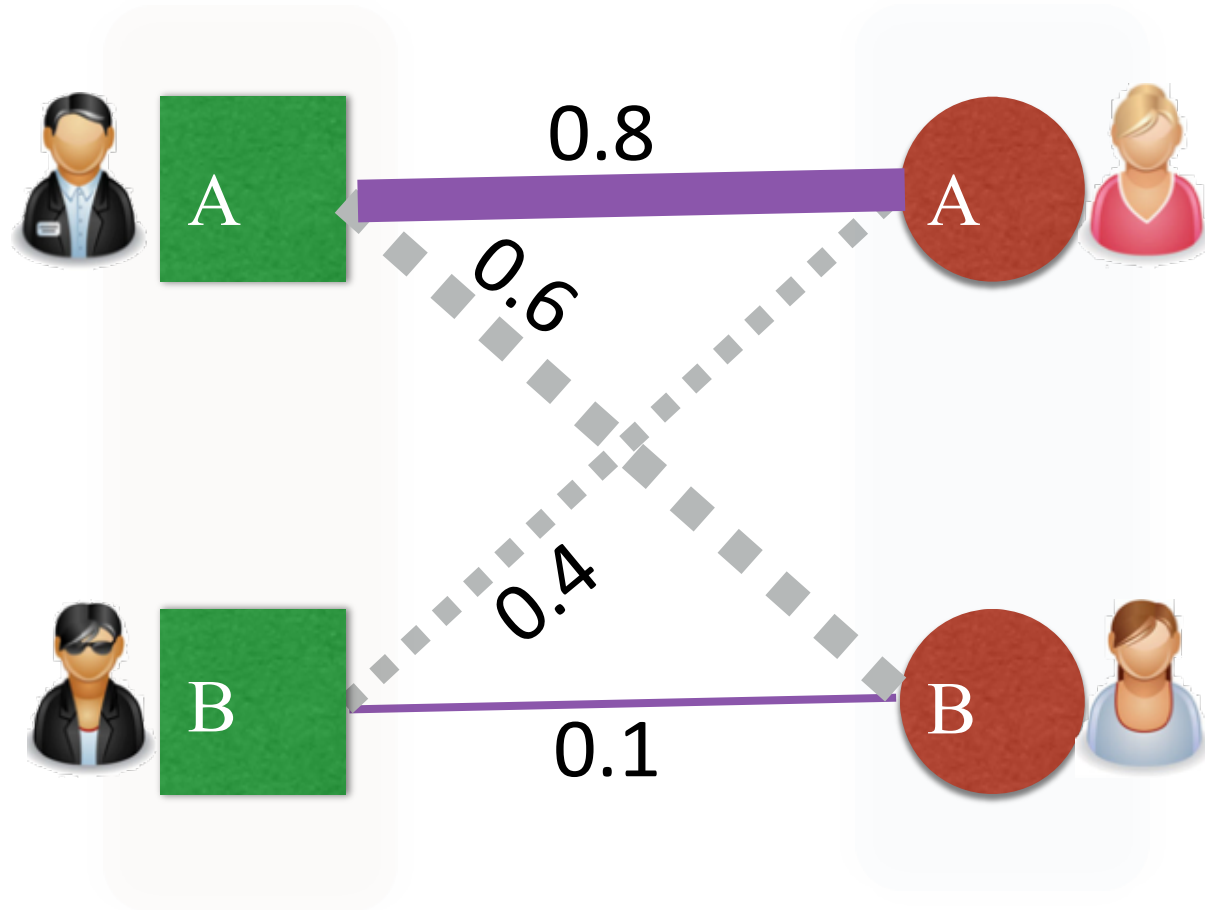
Courtesy of Xiangnan Kong for this slide

# Extract Heterogeneous Cross-Network Features

# Max Sum of Scores w.r.t. Constraints

[Kong et al. CIKM'13]



Blocking Pair

0.8

0.6

0.4

0.1

# Stable Matching/Marriage

# References

- S.Liu, S.Wang, F.Zhu, J.Zhang and R. Krishnan, "**HYDRA: Large-scale Social Identity Linkage via Heterogeneous Behavior Modeling**", SIGMOD 2014.
- R. Zafarani and H. Liu, **"Connecting users across social media sites: A behavioral-modeling approach"**.  KDD 2013.
- X. Kong, J. Zhang, and P.S. Yu, "**Inferring Anchor Links across Multiple Heterogeneous Social Networks**", ACM CIKM 2013.
- J. Zhang, X. Kong, and P.S. Yu,  "**Predicting Social Links for New Users across Aligned Heterogeneous Social Networks**" IEEE ICDM, 2013.
- J. Zhang, X. Kong, and P.S. Yu, "**Transferring Heterogeneous Links across Location-Based Social Networks**", ACM WSDM 2014.

# Conclusion

- Network structure and user behavior play important roles in network mining and analysis for *social applications*.

- In particular, we have shown how they can be employed to study
  - Network Correlation and Patterns
  - Frequent Network Patterns
  - Collaboration Patterns

- When combined, they can even shed new light into
  - Relationship Mining
  - User Identity Linkage