

# Learning to Infer Social Ties in Large Networks

Wenbin Tang, Honglei Zhuang, Jie Tang  
Dept. of Computer Science  
Tsinghua University



# Real social networks are complex...

- Nobody exists only in one social network.
  - Public network vs. private network
  - Business network vs. family network
- However, existing networks (e.g., Facebook and Twitter) are trying to lump everyone into one big network
  - FB tries to solve this problem via **lists/groups**
  - **However...**



- Google+

**which circle?** Users do not take time to create it.

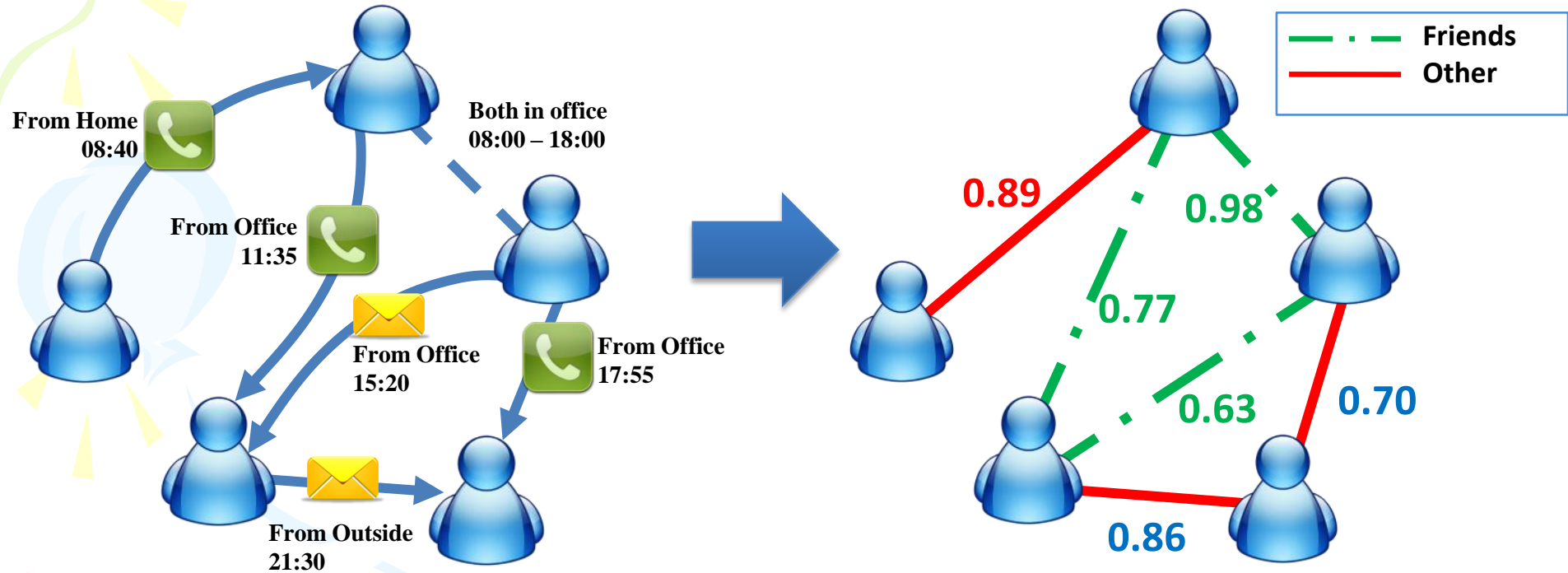




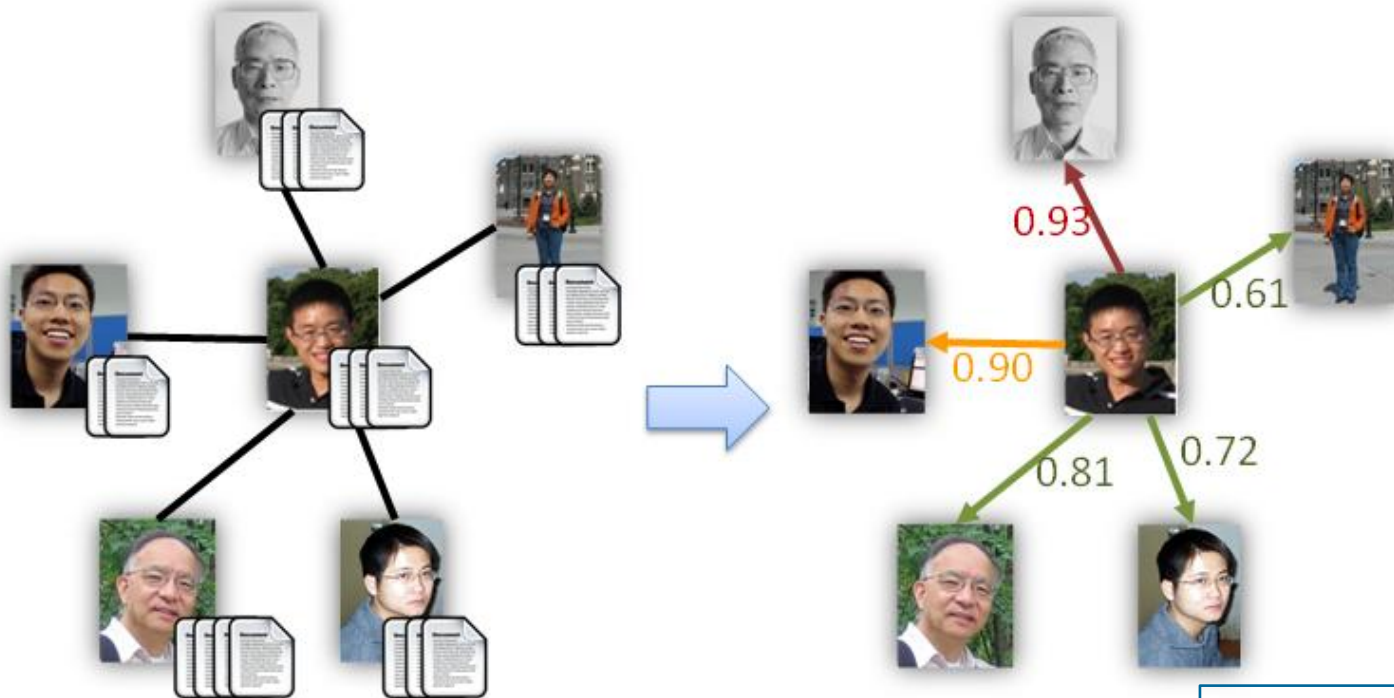
# Even complex than we imaged!

- Only 16% of mobile phone users in Europe have created custom contact groups
  - *users do not* take the time to create it
  - *users do not* know how to circle their friends
- The fact is that our social network is **black-white**...

# Example: Mobile network



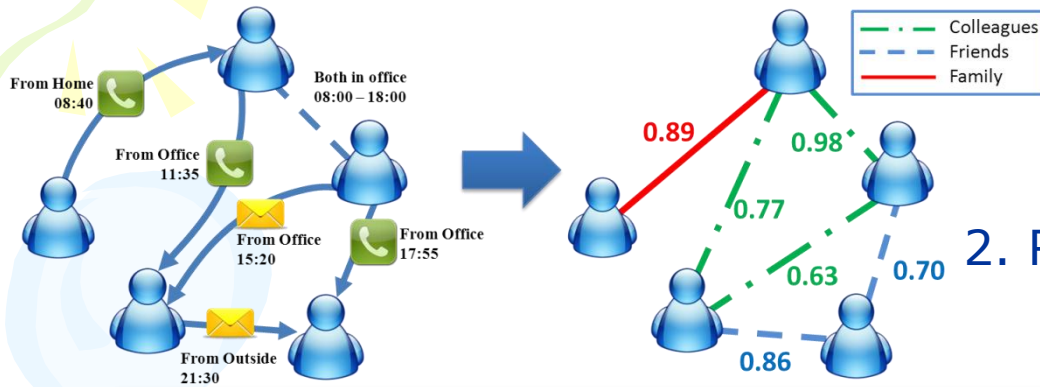
# Example: Coauthor networks



Advisor-Advisee  
Advisee-Advisor  
Coauthor

# Challenges

## 1. Relationships in Mobile Network



## 2. Relationships in Publication Network

### **Challenges:**

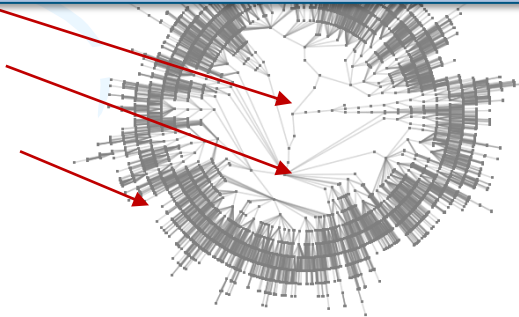
- A generalized framework for inferring social ties?
- A scalable, efficient method?

3.  
Co

How to infer

Manager

Employee

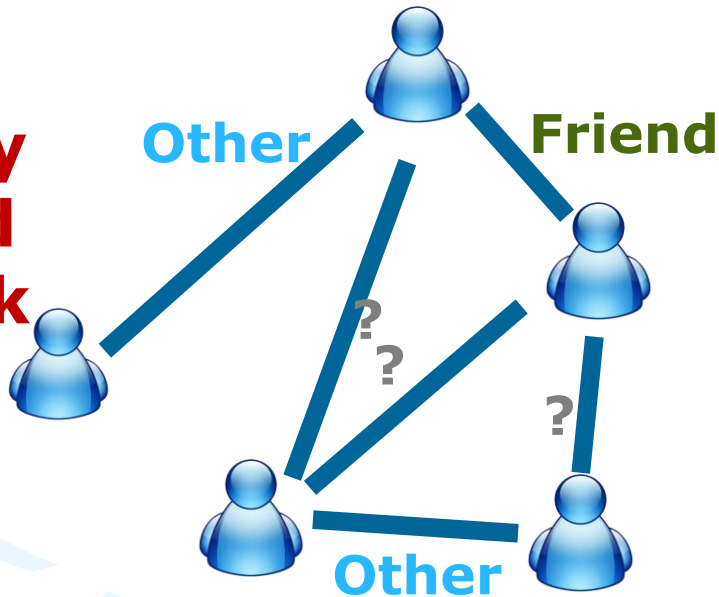


Advisor-Advisee  
Advisee-Advisor  
Coauthor

# Problem Formulation

Input:  $G = (V, E^L, E^U, R^L, W)$

**Partially  
Labeled  
Network**



$V$ : Set of Users

$E^L, R^L$ : Labeled relationships

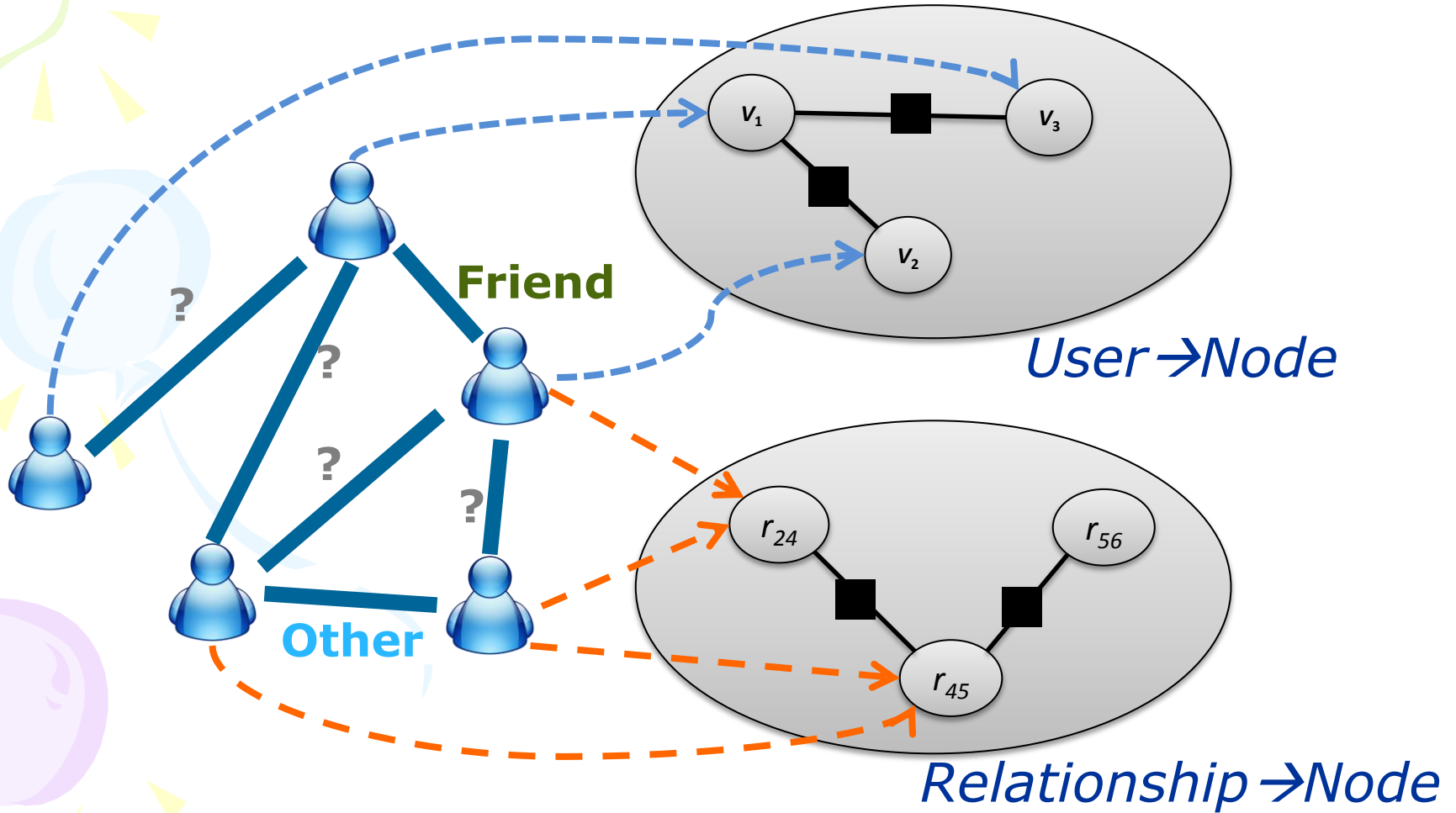
$E^U$ : Unlabeled relationships

Input:  
 $G = (V, E^L, E^U, R^L, W)$



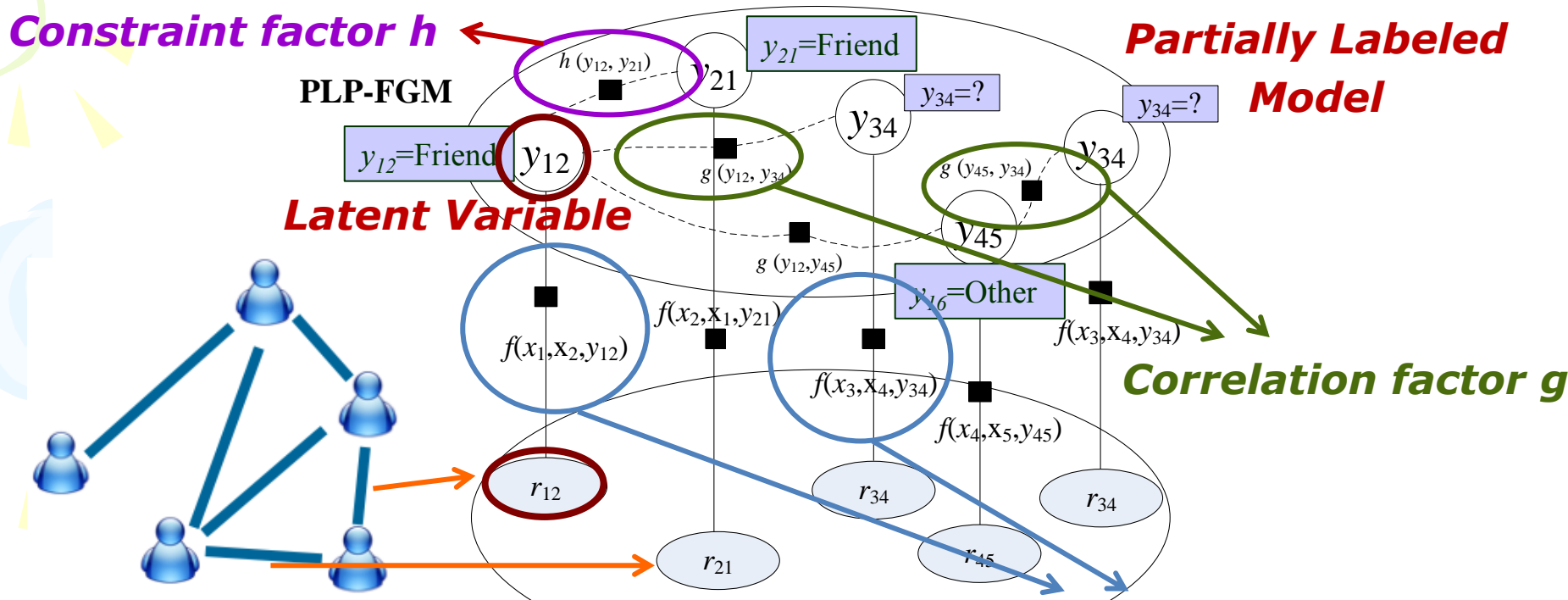
Output:  
 $f: G \rightarrow R$

# Basic Idea





# Partially Labeled Pairwise Factor Graph Model (PLP-FGM)



## Problem:

Ma For each relationship, identify which type has the highest probability?

# Solutions<sub>(con't)</sub>

- Different ways to instantiate factors
  - We use exponential-linear functions
    - Attribute Factor:

$$f(y_i, \mathbf{x}_i) = \frac{1}{Z_\lambda} \exp\{\lambda^T \Phi(y_i, \mathbf{x}_i)\}$$

- Correlation / Constraint Factor:

$$g(y_i, G(y_i)) = \frac{1}{Z_\alpha} \exp\left\{ \sum_{y_j \in G(y_i)} \alpha^T \mathbf{g}(y_i, y_j) \right\}$$

$$h(y_i, H(y_i)) = \frac{1}{Z_\beta} \exp\left\{ \sum_{y_j \in H(y_i)} \beta^T \mathbf{h}(y_i, y_j) \right\}$$

- $\theta = [\lambda, \alpha, \beta], s = [\Phi^T, g^T, h^T]^T$
- Log-Likelihood of labeled Data:

$$\mathcal{O}(\theta) = \log \sum_{Y|Y^L} \exp\{\theta^T \mathbf{S}\} - \log \sum_Y \exp\{\theta^T \mathbf{S}\}$$

# Learning Algorithm

- Maximize the log-likelihood of labeled relationships

**Input:** learning rate  $\eta$

**Output:** learned parameters  $\theta$

Initialize  $\theta$ ;

repeat

    Calculate  $\mathbb{E}_{p_{\theta}(Y|Y^L,G)}\mathbf{S}$  using LBP ;

    Calculate  $\mathbb{E}_{p_{\theta}(Y|G)}\mathbf{S}$  using LBP ;

    Calculate the gradient of  $\theta$  according to Eq. 7:

$$\nabla_{\theta} = \mathbb{E}_{p_{\theta}(Y|Y^L,G)}\mathbf{S} - \mathbb{E}_{p_{\theta}(Y|G)}\mathbf{S}$$

    Update parameter  $\theta$  with the learning rate  $\eta$ : **Expectation Computing**

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \cdot \nabla_{\theta} \quad \text{Loopy Belief Propagation}$$

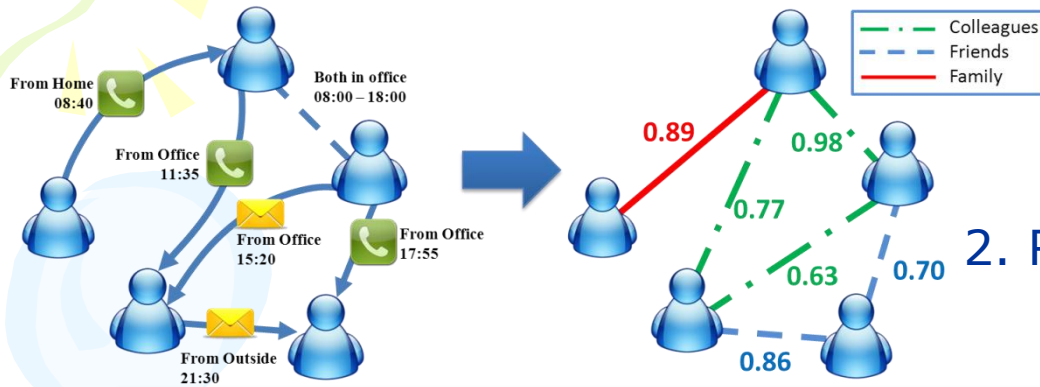
until *Convergence*;

Algorithm 1: Learning PLP-FGM.

**Gradient Decent Method**

# Challenges

## 1. Relationships in Mobile Network



## 2. Relationships in Publication Network

### **Challenges:**

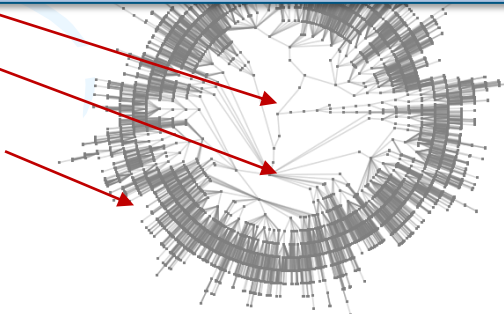
- A generalized framework for inferring social ties?
- A scalable, efficient method?

## 3. Co

How to infer

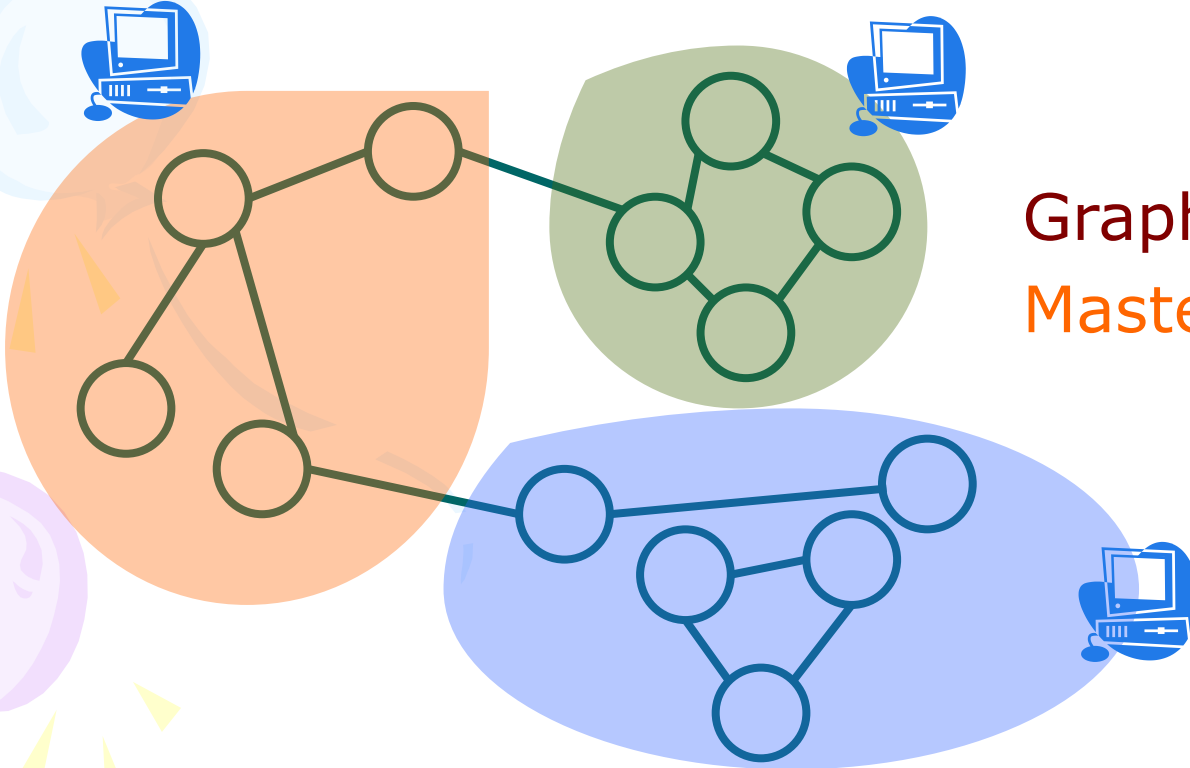
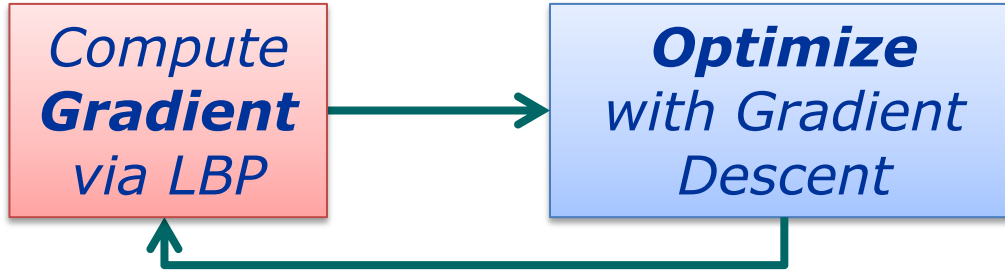
Manager

Employee



Advisor-Advisee  
Advisee-Advisor  
Coauthor

# Distributed Learning



Graph Partition  
Master-Slave Computing

# Data Sets

- Coauthor Network (Publication)
  - To infer **Advisor-Advisee** relationship
  - Papers from DBLP
- Email Network (Email)
  - To infer **Manager-Subordinate** relationship
  - Using Enron Email Dataset
- Mobile Network (Mobile)
  - To infer **Friendship**
  - 107 users (ten-month). Published by MIT

<b>Data Set</b>	<b>Users</b>	<b>Unlabeled Relationships</b>	<b>Labeled Relationships</b>
Publication	1,036,990	1,984,164	6,096
Email	151	3,424	148
Mobile	107	5,122	314



# Baselines

- **Baselines:**

- SVM:

- Use the same feature defined in our model to train a classification model

- TCFG:

- An unsupervised method to identify advisor-advisee relationships

- **PLP-FGM-S**

- Do not use partially-labeled property
    - Train parameters on **the labeled sub-graph**

# Performance Analysis

Data Set	Method	Precision	Recall	F <sub>1</sub> -score
Publication	SVM	72.5	54.9	62.1
	TPFG	82.8	<b>89.4</b>	86.0
	PLP-FGM-S	77.1	78.4	77.7
	PLP-FGM	<b>91.4</b>	87.7	<b>89.5</b>
Email	SVM	79.1	<b>88.6</b>	83.6
	PLP-FGM-S	85.8	85.6	85.7
	PLP-FGM	<b>88.6</b>	87.2	<b>87.9</b>
Mobile	SVM	<b>92.7</b>	64.9	76.4
	PLP-FGM-S	88.1	71.3	78.8
	PLP-FGM	89.4	<b>75.2</b>	<b>81.6</b>

**SVM:** Use the same feature to train a classification model

**TPFG:** An unsupervised method to identify advisor-advisee relationships

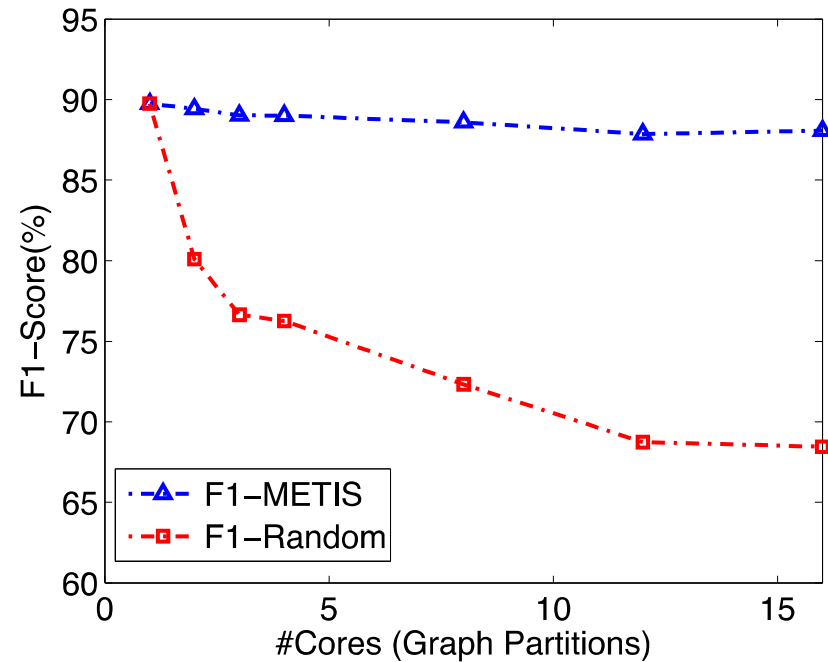
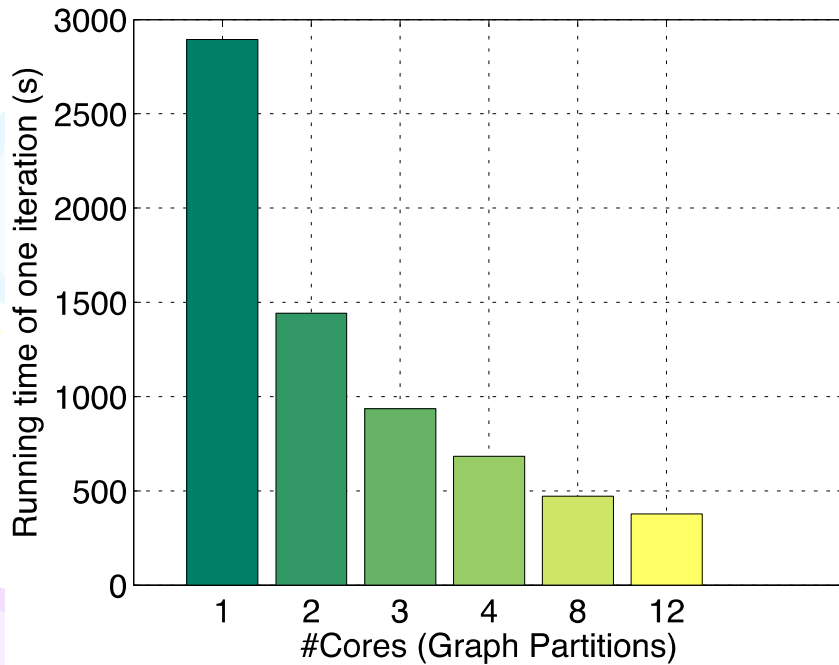
**PLP-FGM-S:** Train PLP-FGM model on the labeled sub-graph



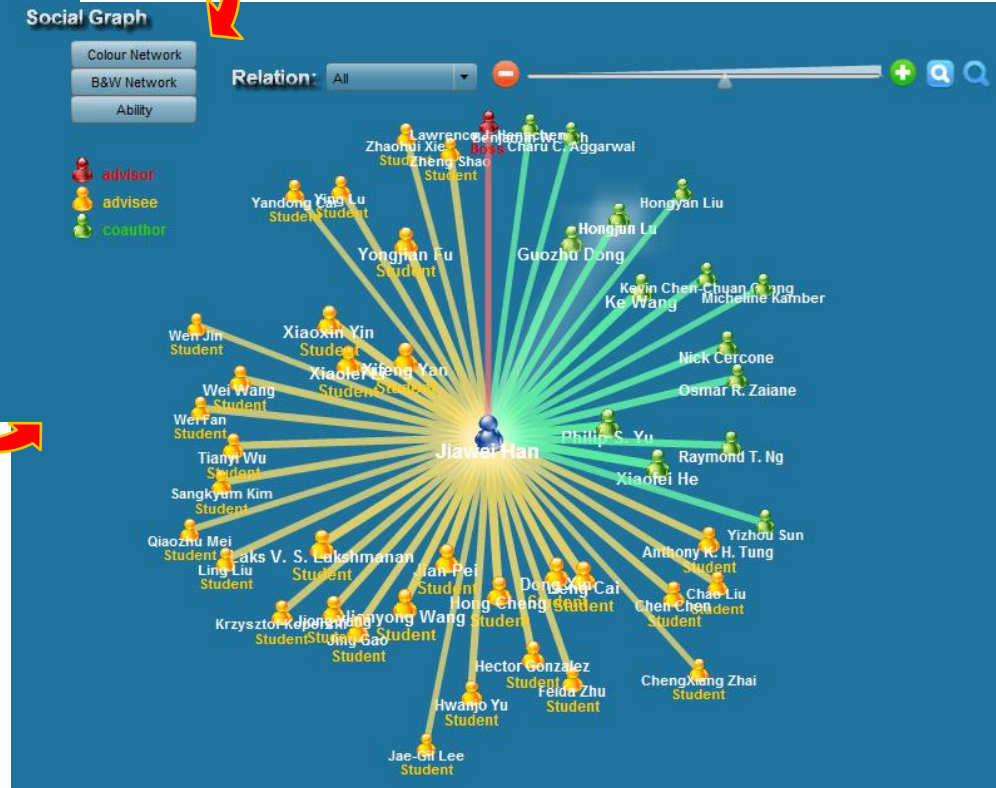
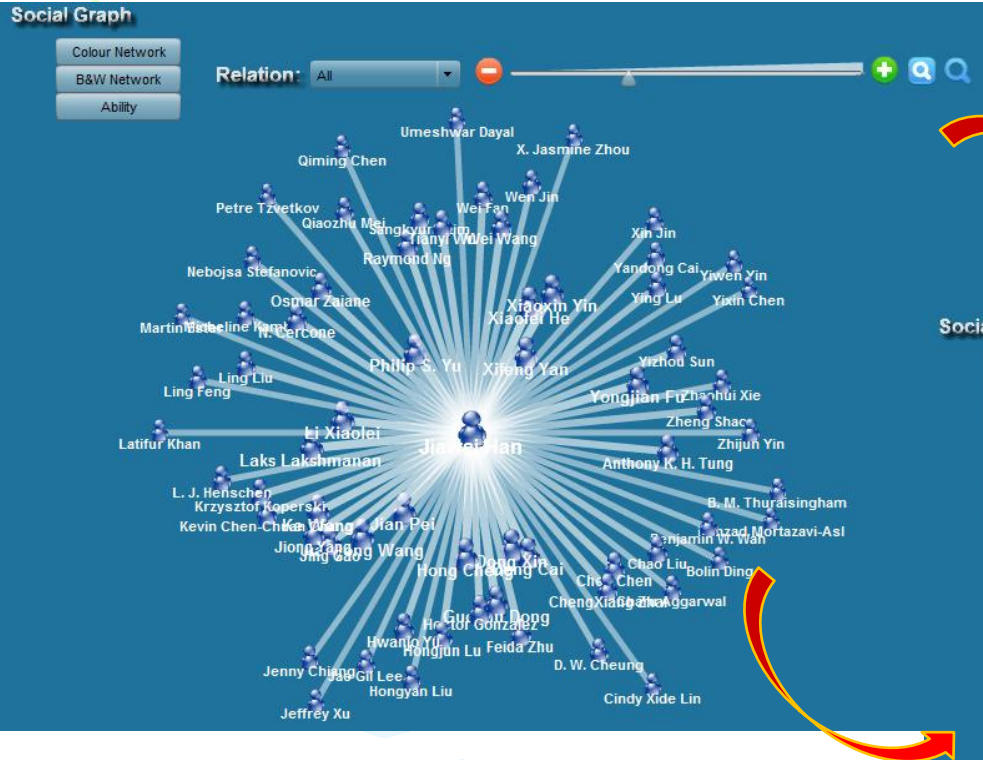
# Factor Contribution Analysis

Data Set	Factor used	F <sub>1</sub> -score
Publication	Attributes	64.9
	+Co-advisor	75.0(+10.1%)
	+Co-advisee	74.7(+9.8%)
	All	89.5(+24.6%)
Email	Attributes	80.3
	+Co-recipient	80.6(+0.3%)
	+Co-manager	83.2(+2.9%)
	+Co-subordinate	85.0(+4.7%)
	All	87.9(+7.6%)
Mobile	Attributes	80.2
	+Co-location	80.4(+0.2%)
	+Related-call	80.2(+0.0%)
	All	81.6(+1.4%)

# Distributed Learning Performance



# System Arnetminer





# Conclusion

- Formulate the problem of inferring the types of social ties
- Propose the PLP-FGM model to solve this problem, and present a distributed learning algorithm
- Validate the approach in different real data sets



# Future work

- Make online social networks *colorful*
  - How to involve user into learning process?
  - Connect with social theories?

The background features three large, stylized swirls in purple, green, and blue. Interspersed among these swirls are several yellow starburst shapes, each composed of multiple triangular points radiating from a central point. The overall aesthetic is bright and celebratory.

# Thank you!

Any Questions?



# Correlation Definition

- Mobile Dataset:
  - Co-location
    - 3 users in the same location.
  - Related-call
    - A Make a call to B&C at the same place/time
- For more information, please refer to the paper😊



# Feature Definition

Data set	Factor	Description	
Publication	Paper count	$ P_i ,  P_j $	
	Paper ratio	$ P_i / P_j $	
	Coauthor ratio	$ P_i \cap P_j / P_i ,  P_i \cap P_j / P_j $	
	Conference coverage	The proportion of the conferences which both $v_i$ and $v_j$ attended among conferences $v_j$ attended.	
	First-paper-year-diff	The difference in year of the earliest publication of $v_i$ and $v_j$ .	
Email	Traffics	Sender	Recipients Include
		$v_i$	$v_j$
		$v_j$	$v_i$
		$v_i$	$v_k$ and not $v_j$
		$v_j$	$v_k$ and not $v_i$
		$v_k$	$v_i$ and not $v_j$
		$v_k$	$v_j$ and not $v_i$
Mobile	#voice calls	The total number of voice call logs between two users.	
	#messages	Number of messages between two users.	
	Night-call ratio	The proportion of calls at night (8pm to 8am).	
	Call duration	The total duration time of calls between two users.	
	#proximity	The total number of proximity logs between two users.	
	In-role proximity ratio	The proportion of proximity logs in “working place” and in working hours (8am to 8pm).	





# Existing Methods...

- [Diehl:07] try to identify the relationships by learning a ranking function in **Email network**.
- Wang et al. [Wang:10] propose an unsupervised algorithm for mining the advisor-advisee relationships from the **Publication network**.
- Both algorithms focus on a specific domain
  - not easy to extend to other problems.