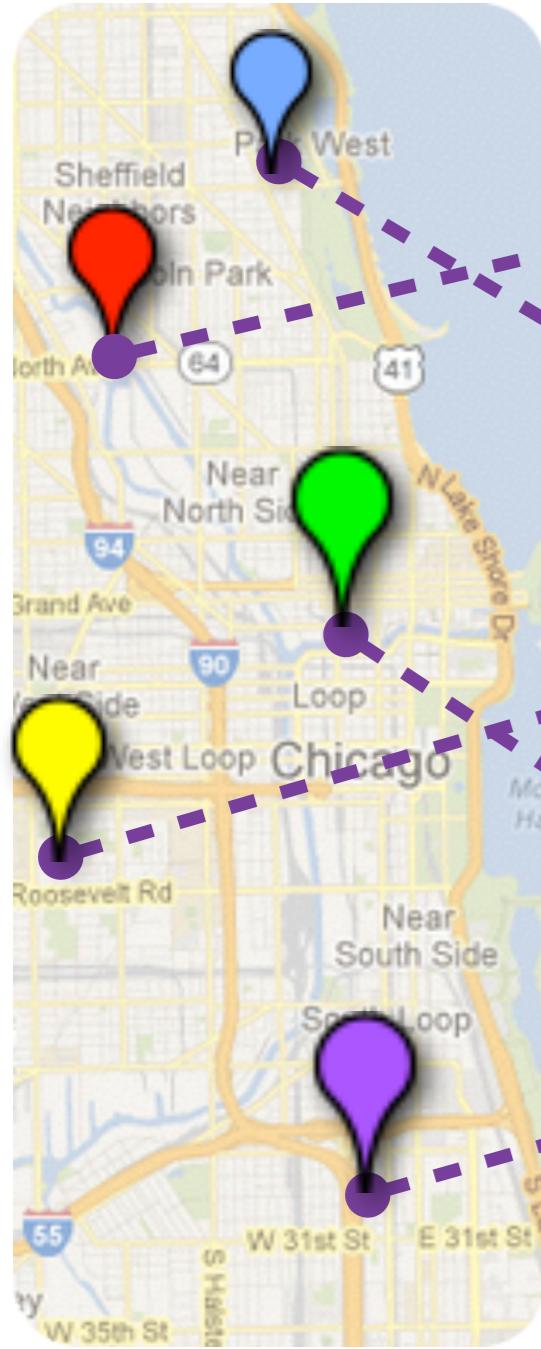


Meta-path based Multi-Network Collective Link Prediction

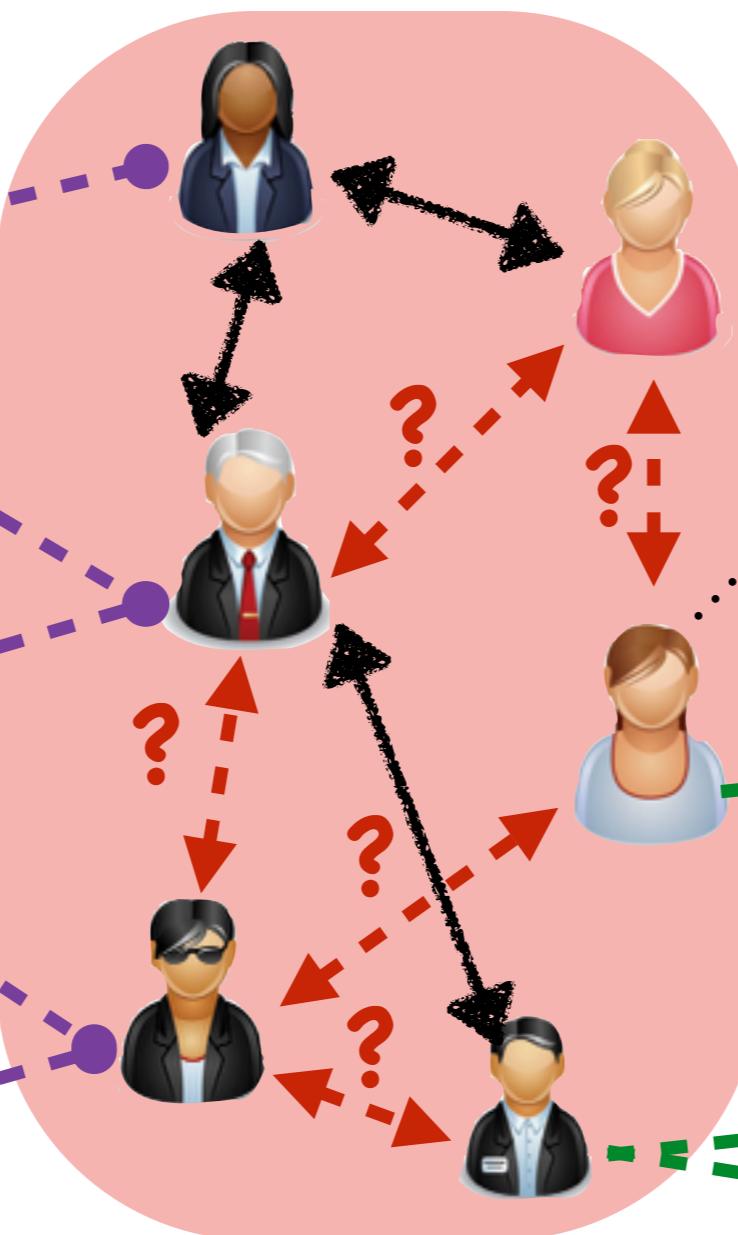
Jiawei Zhang^{1,2}, Philip S. Yu¹, Zhi-Hua Zhou²
University of Illinois at Chicago², Nanjing University²

Traditional social link prediction in one single social network

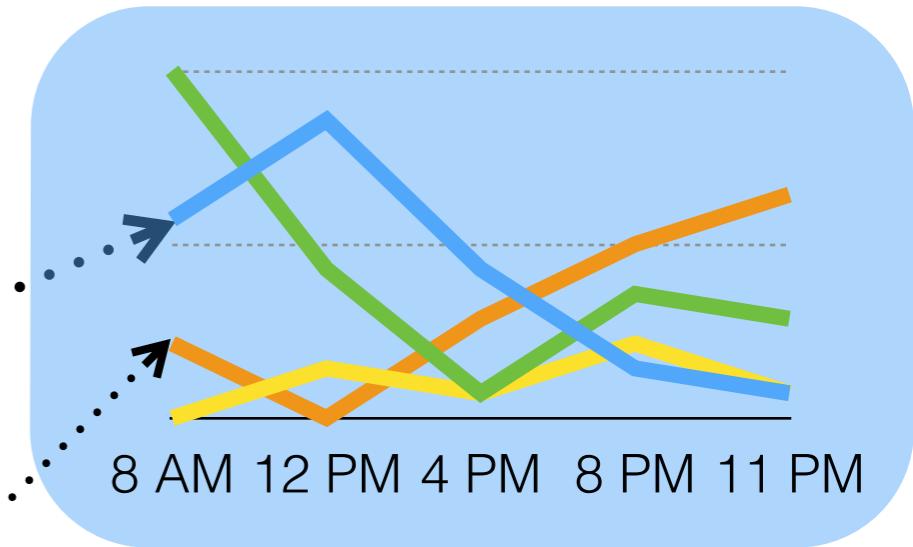
Locations



Social Links



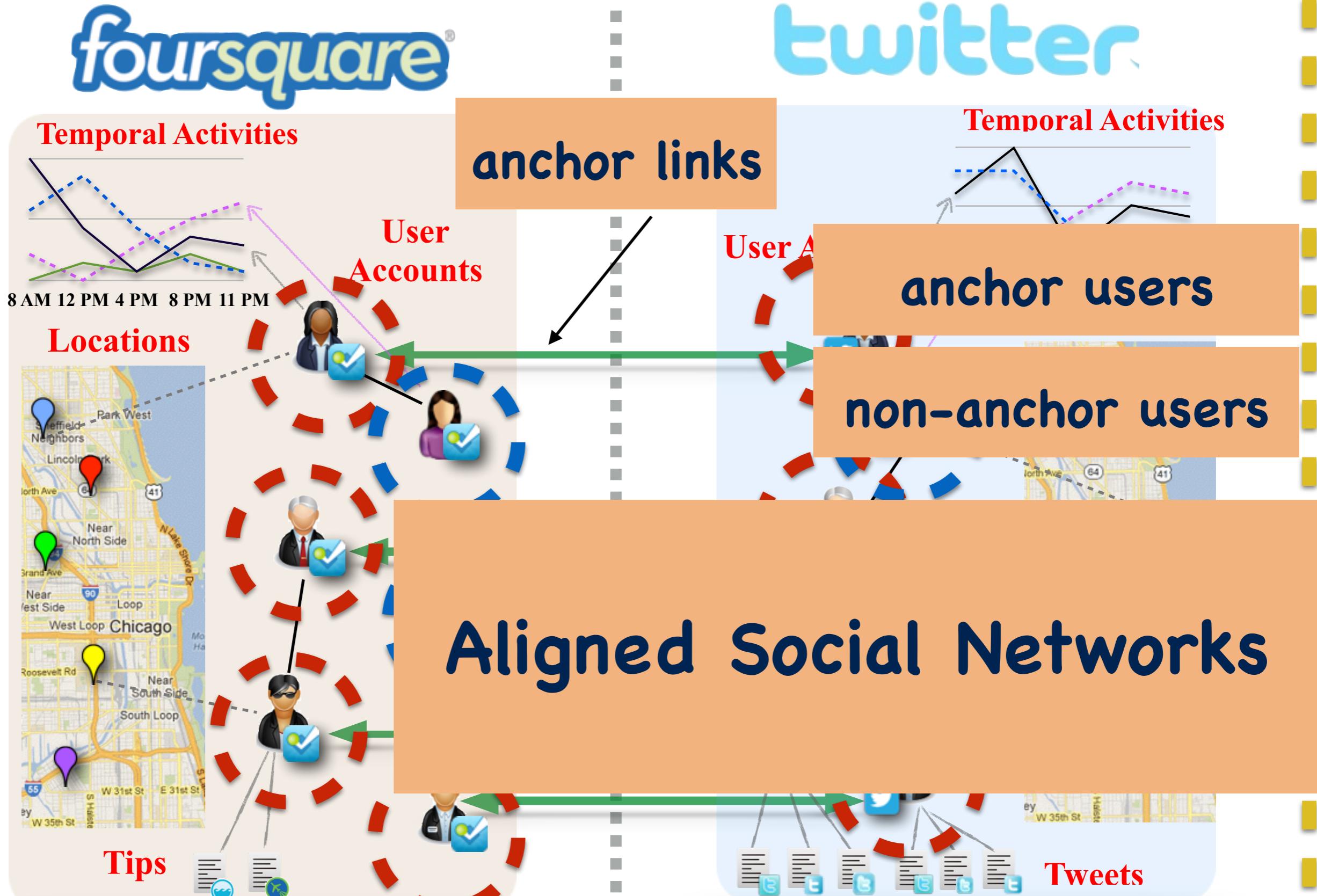
Temporal Activities



Contents: Tweets

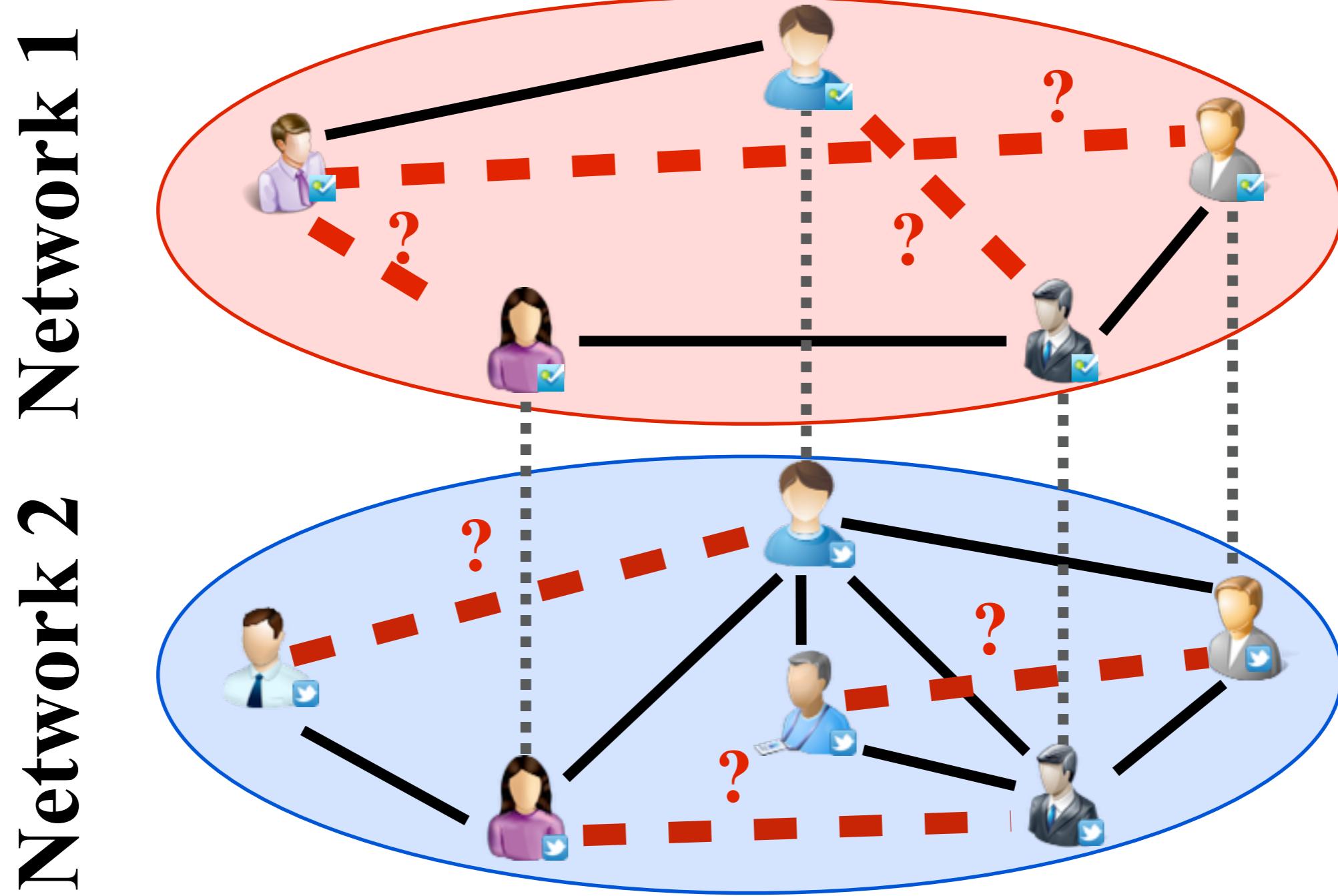


Users use multiple social networks simultaneously



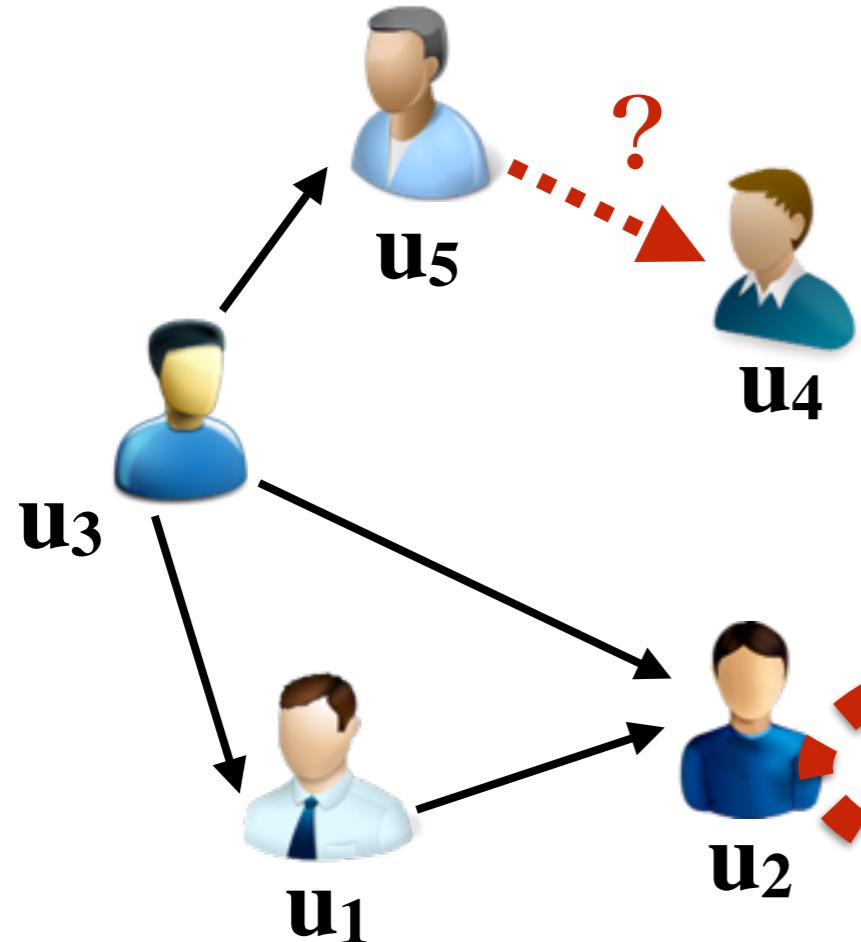
Predicting social links in multiple aligned networks simultaneously

..... anchor link — existing social links - ? - - social links to be predicted



class imbalance problem
negative instances >>
positive instances

network structure



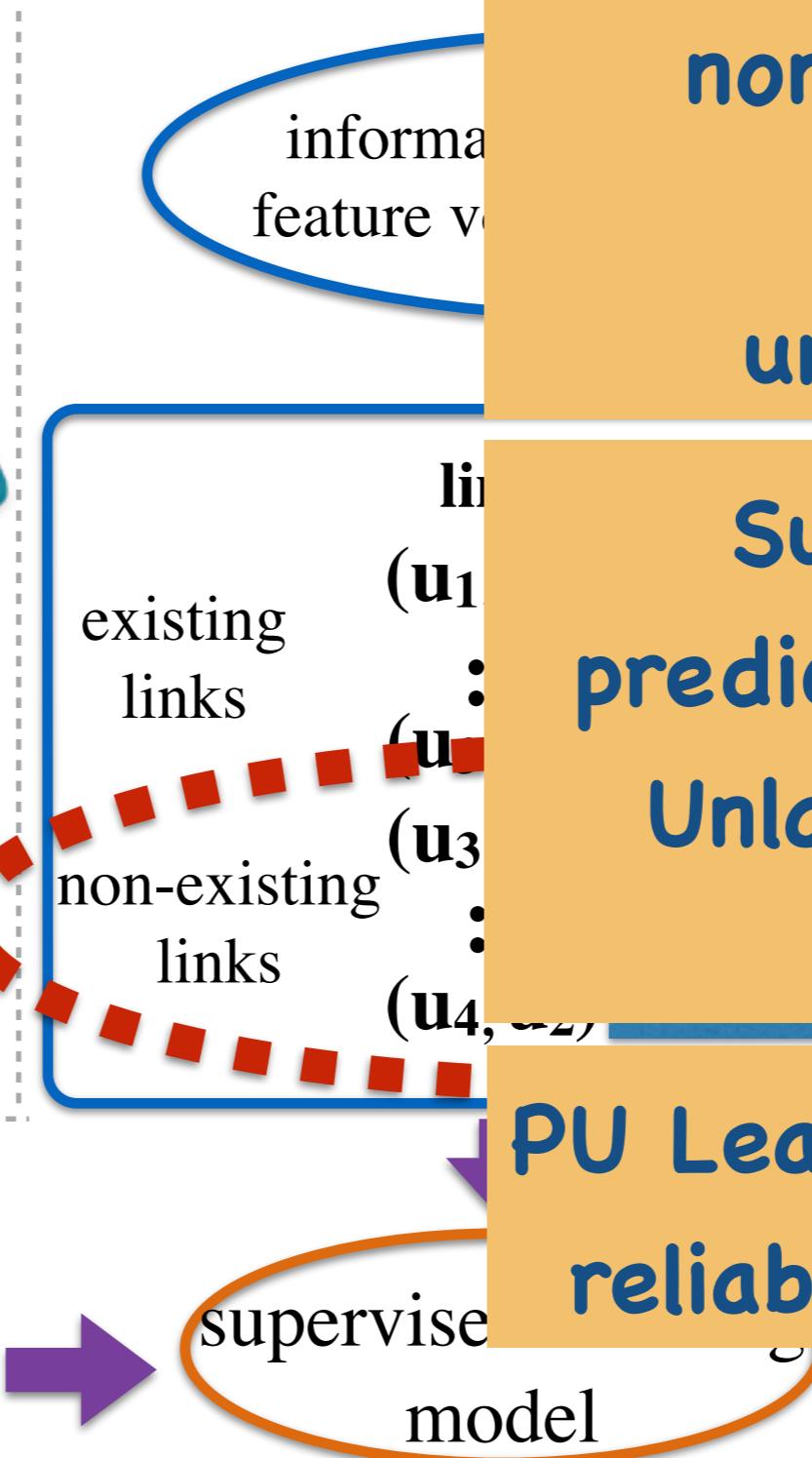
link to be predicted
 (u_5, u_4)

non-existing links
!=
negative links

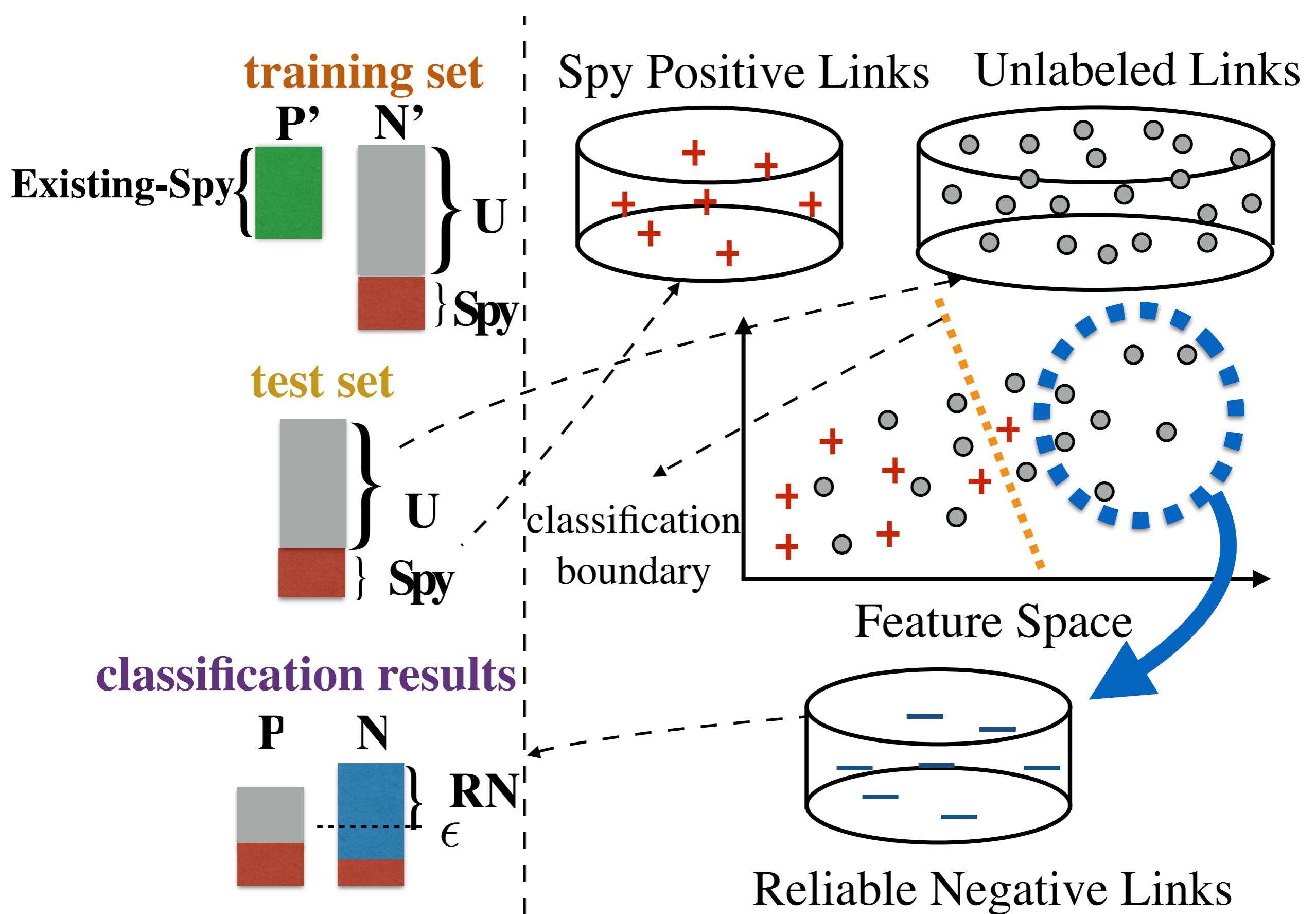
non-existing links
should be
unlabeled links

Supervised link
prediction ==> Positive
Unlabeled (PU) link
prediction

PU Learning: How to find
reliable negative links?
label/score

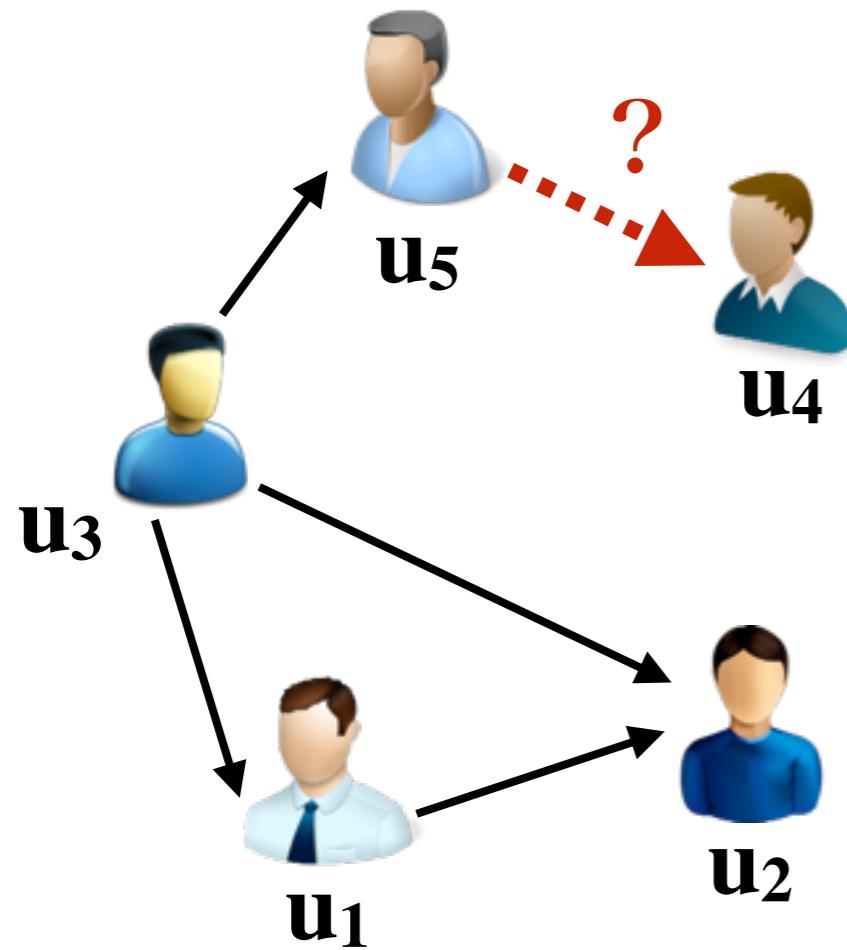


Reliable Negative Links Extraction



PU Link Prediction Setting

network structure



what kind of information
are there in the network?

	link	features	label
existing links	(u ₁ , u ₂)	[blue bar]	+1
⋮	(u ₃ , u ₅)	[blue bar]	+1
reliable negative links	(u _x , u _y)	[blue bar]	-1
⋮	(u _x , u _y)	[blue bar]	-1

link to be predicted

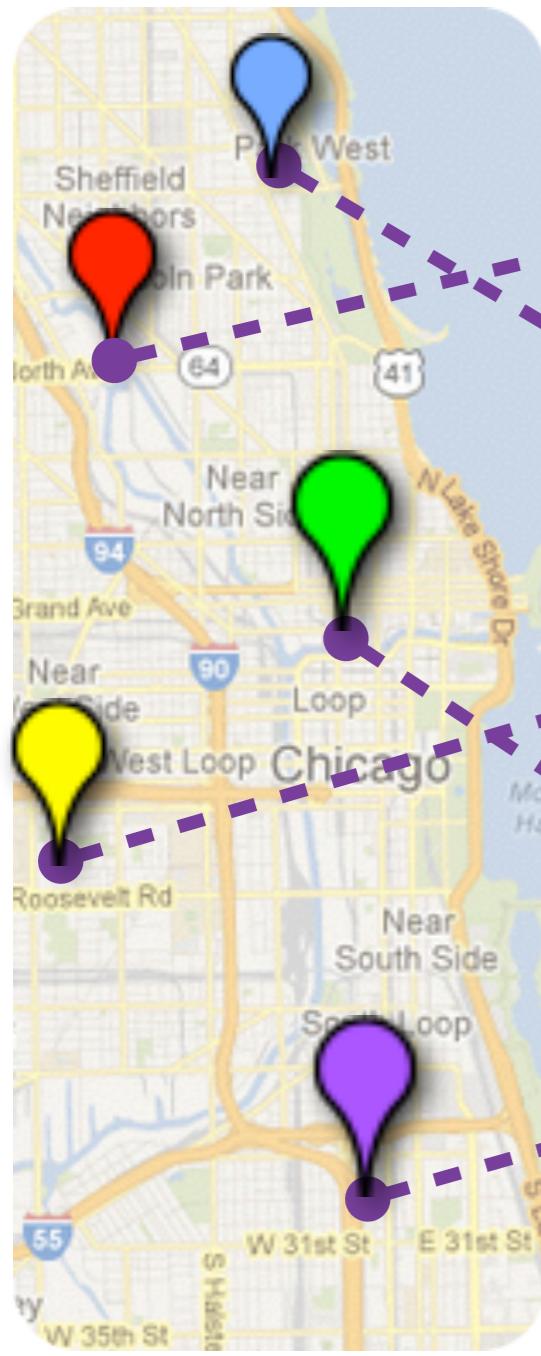
(u₅, u₄)

supervised learning
model

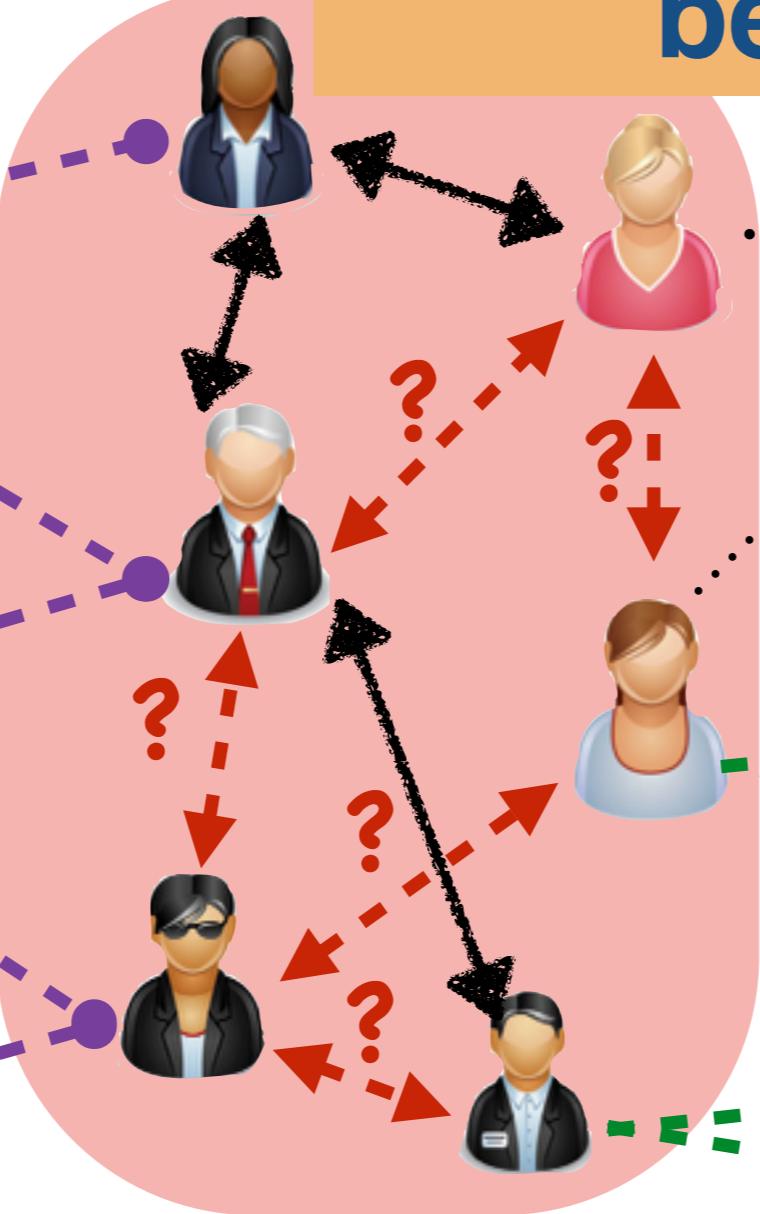
scores

Heterogeneous Information

Locations



Social what kind of features can
be extracted ?

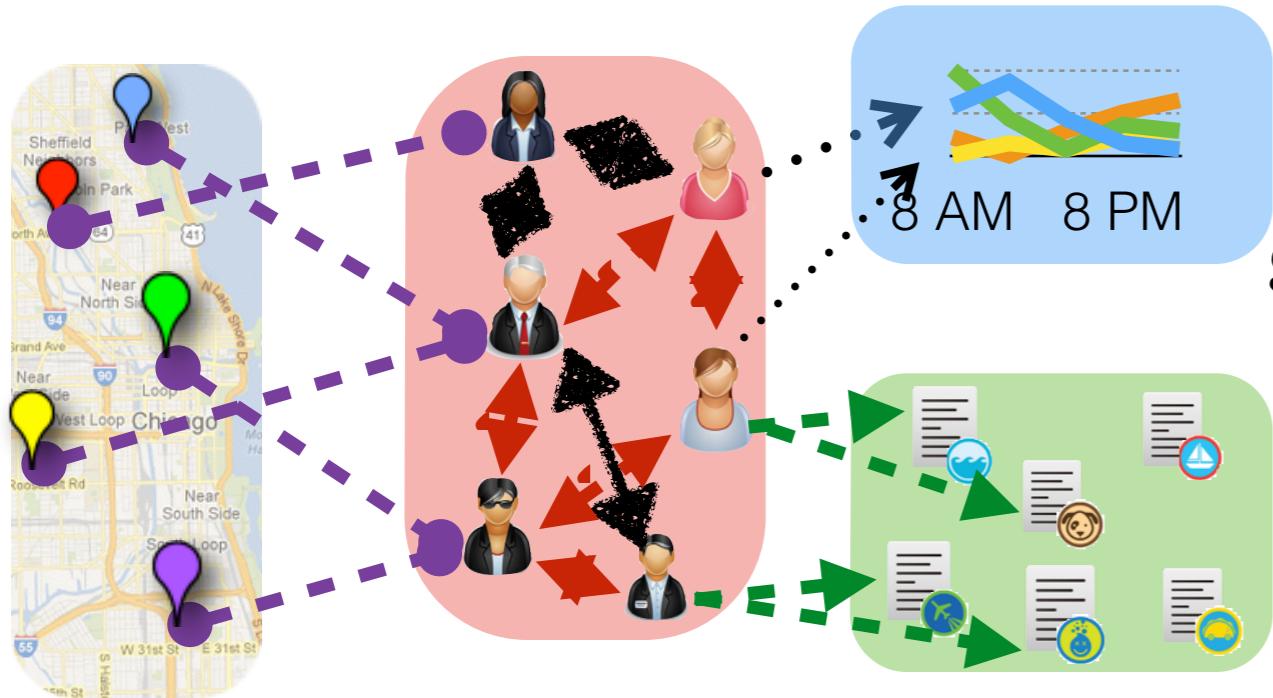


8 AM 12 PM 4 PM 8 PM 11 PM

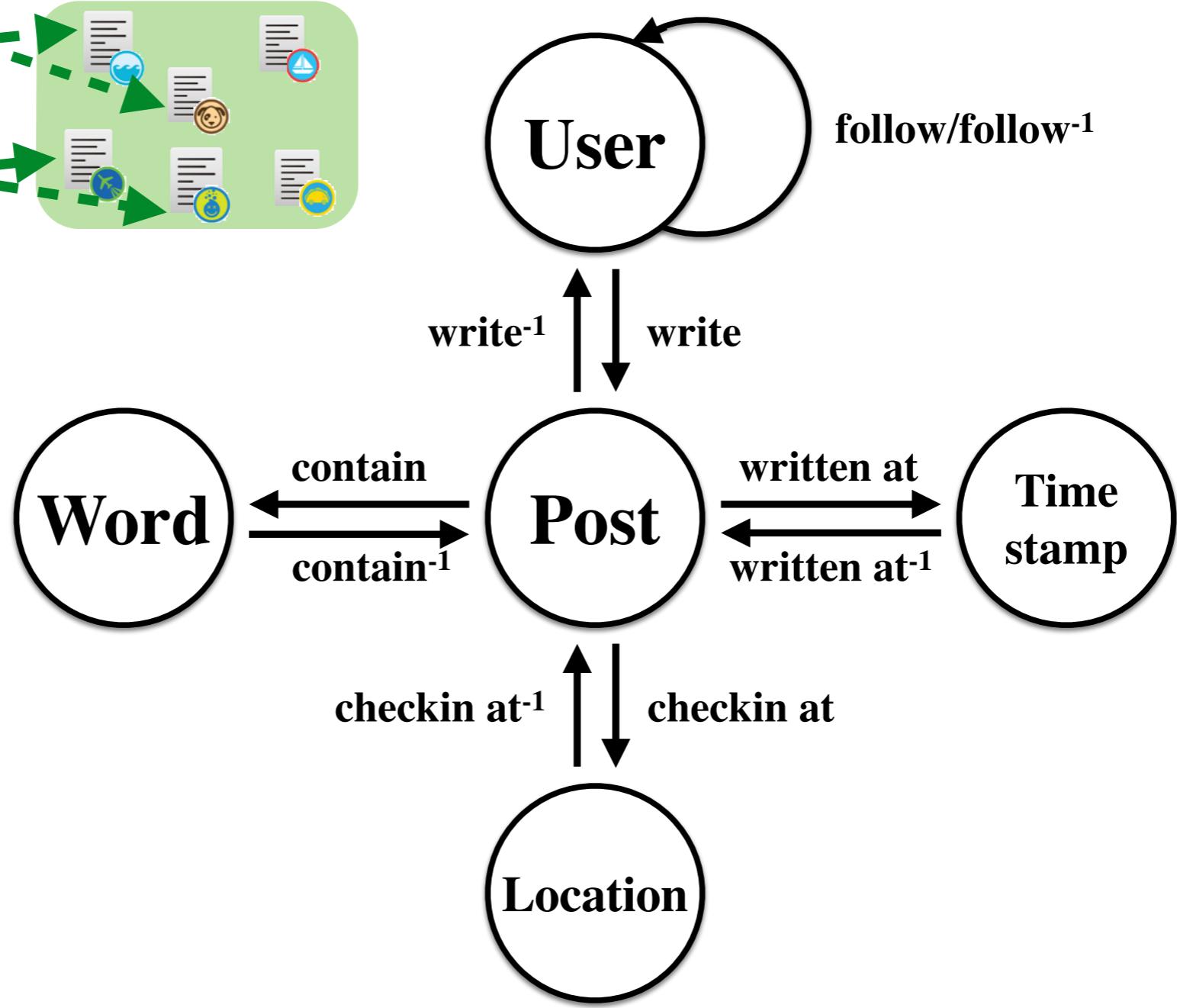
Contents: Tweets



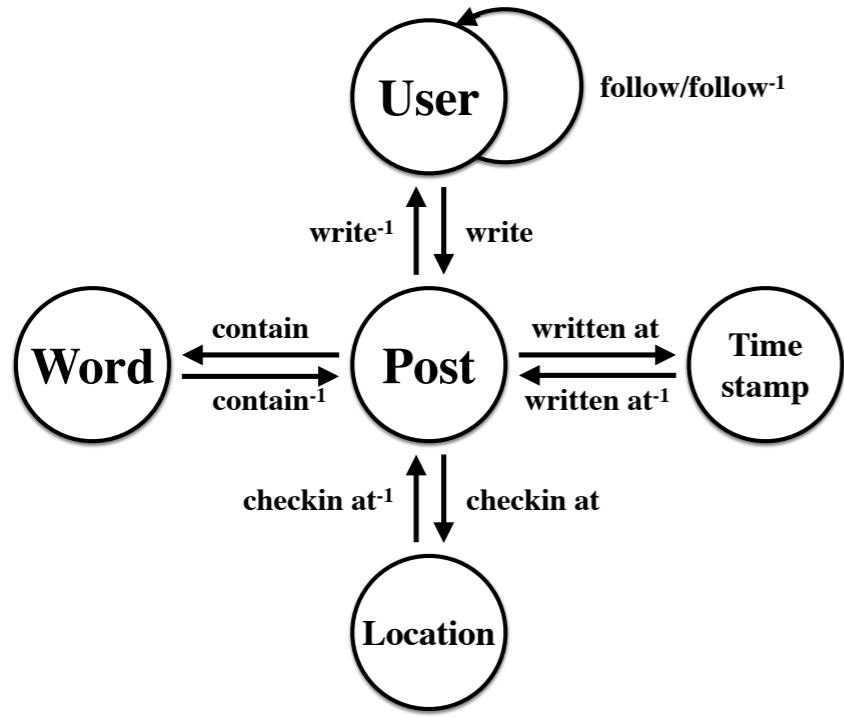
Network Schema



social network schema



Intra-network social meta paths



Definition 10 (Intra-Network Social Meta Path): For a given meta path $\Phi = T_1 \xrightarrow{R_1} T_2 \xrightarrow{R_2} \dots \xrightarrow{R_{k-1}} T_k$ defined based on S_G , if T_1 and T_k are both the “User” node type, then P is defined as a *social meta path*. Depending on whether T_1, \dots, T_k and R_1, \dots, R_{k-1} are the same or not, P can be divided into two categories: *homogeneous intra-network social meta path* and *heterogeneous intra-network social meta path*.

Homogeneous Intra-Network Social Meta Path

- **ID 0. Follow:** User $\xrightarrow{\text{follow}}$ User, whose notation is “ $U \rightarrow U$ ” or $\Phi_0(U, U)$.
- **ID 1. Follower of Follower:** User $\xrightarrow{\text{follow}}$ User $\xrightarrow{\text{follow}}$ User, whose notation is “ $U \rightarrow U \rightarrow U$ ” or $\Phi_1(U, U)$.
- **ID 2. Common Out Neighbor:** User $\xrightarrow{\text{follow}}$ User $\xrightarrow{\text{follow}^{-1}}$ User, whose notation is “ $U \rightarrow U \leftarrow U$ ” or $\Phi_2(U, U)$.
- **ID 3. Common In Neighbor:** User $\xrightarrow{\text{follow}^{-1}}$ User $\xrightarrow{\text{follow}}$ User, whose notation is “ $U \leftarrow U \rightarrow U$ ” or $\Phi_3(U, U)$.

Heterogeneous Intra-Network Social Meta Path

- **ID 4. Common Words:** User $\xrightarrow{\text{write}}$ Post $\xrightarrow{\text{contain}}$ Word $\xrightarrow{\text{contain}^{-1}}$ Post $\xrightarrow{\text{write}^{-1}}$ User, whose notation is “ $U \rightarrow P \rightarrow W \leftarrow P \leftarrow U$ ” or $\Phi_4(U, U)$.
- **ID 5. Common Timestamps:** User $\xrightarrow{\text{write}}$ Post $\xrightarrow{\text{contain}}$ Time $\xrightarrow{\text{contain}^{-1}}$ Post $\xrightarrow{\text{write}^{-1}}$ User, whose notation is “ $U \rightarrow P \rightarrow T \leftarrow P \leftarrow U$ ” or $\Phi_5(U, U)$.
- **ID 6. Common Location Checkins:** User $\xrightarrow{\text{write}}$ Post $\xrightarrow{\text{attach}}$ Location $\xrightarrow{\text{attach}^{-1}}$ Post $\xrightarrow{\text{write}^{-1}}$ User, whose notation is “ $U \rightarrow P \rightarrow L \leftarrow P \leftarrow U$ ” or $\Phi_6(U, U)$.

New network problem

..... anchor

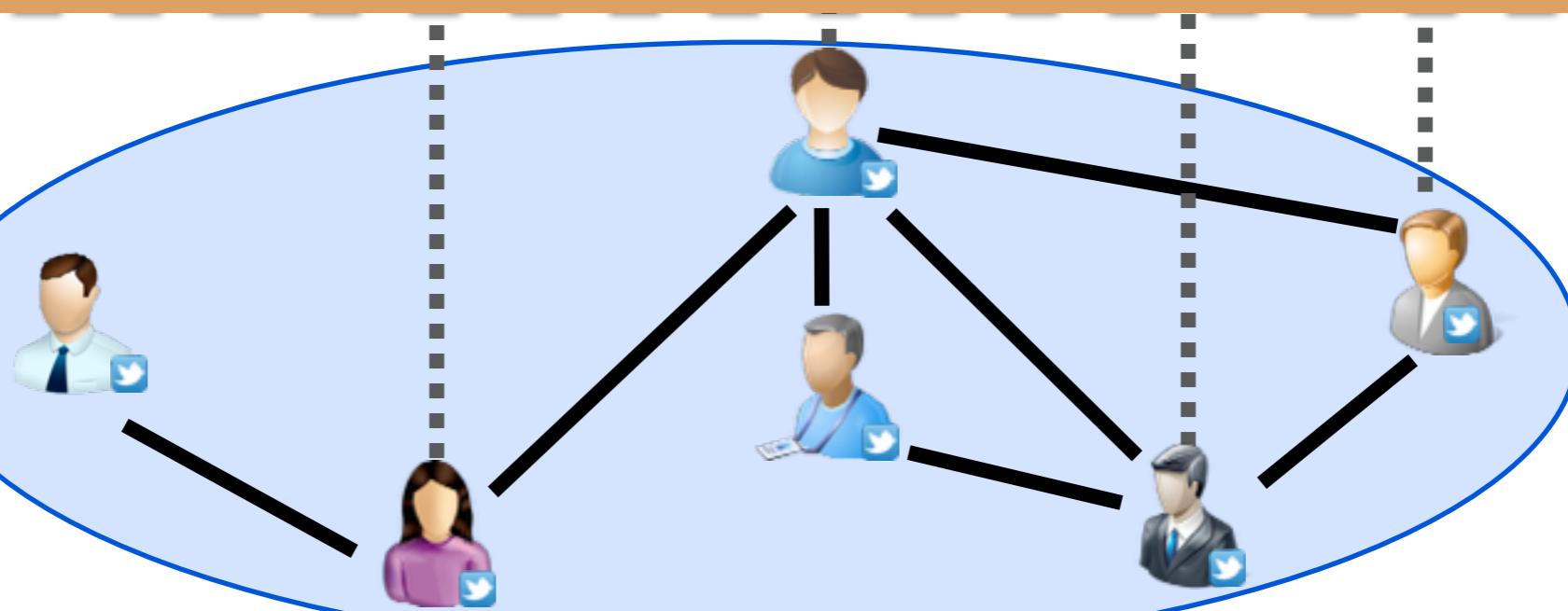
New network

sparse information ==> sparse feature

Network 1

information in other aligned networks can
be transferred to the new network or not?

Network 2



Anchor Meta path & Inter-network social meta paths

Definition 12 (Anchor Meta Path): Let U^i, U^j be the user nodes of G^i and G^j respectively and $A^{i,j}$ be the anchor links between G^i and G^j . Meta path $\Upsilon = T_1 \xleftrightarrow{R_1} T_2$ is an *anchor meta path* between network G^i and G^j iff $T_1 = U^i$ and $T_2 = U^j$ and $R_1 = A^{i,j}$. The notation of *anchor meta path* from G^i to G^j is $\Upsilon(U^i, U^j)$ and the length of $\Upsilon(U^i, U^j)$ is 1.

Definition 13 (Inter-Network Meta Path): Meta path $\Psi = T_1 \xrightarrow{R_1} T_2 \xrightarrow{R_2} \dots \xrightarrow{R_{k-1}} T_k$ is an *inter-network meta path* across G^i and G^j iff $\exists m \in \{1, 2, \dots, k-1\}$, $T_m \xleftrightarrow{R_m} T_{m+1} = \Upsilon(U^i, U^j)$.

Category 1: $\Upsilon(U^i, U^j) \circ (\Phi(U^j, U^j) \cup \Phi_0(U^j, U^j)) \circ \Upsilon(U^j, U^i)$, whose notation is $\Psi_1(U^i, U^i)$;

Category 2.: $(\Phi(U^i, U^i) \cup \Phi_0(U^i, U^i)) \circ \Upsilon(U^i, U^j) \circ (\Phi(U^j, U^j) \cup \Phi_0(U^j, U^j)) \circ \Upsilon(U^j, U^i)$, whose notation is $\Psi_2(U^i, U^i)$;

Category 3.: $\Upsilon(U^i, U^j) \circ (\Phi(U^j, U^j) \cup \Phi_0(U^j, U^j)) \circ \Upsilon(U^j, U^i) \circ (\Phi(U^i, U^i) \cup \Phi_0(U^i, U^i))$, whose notation is $\Psi_3(U^i, U^i)$;

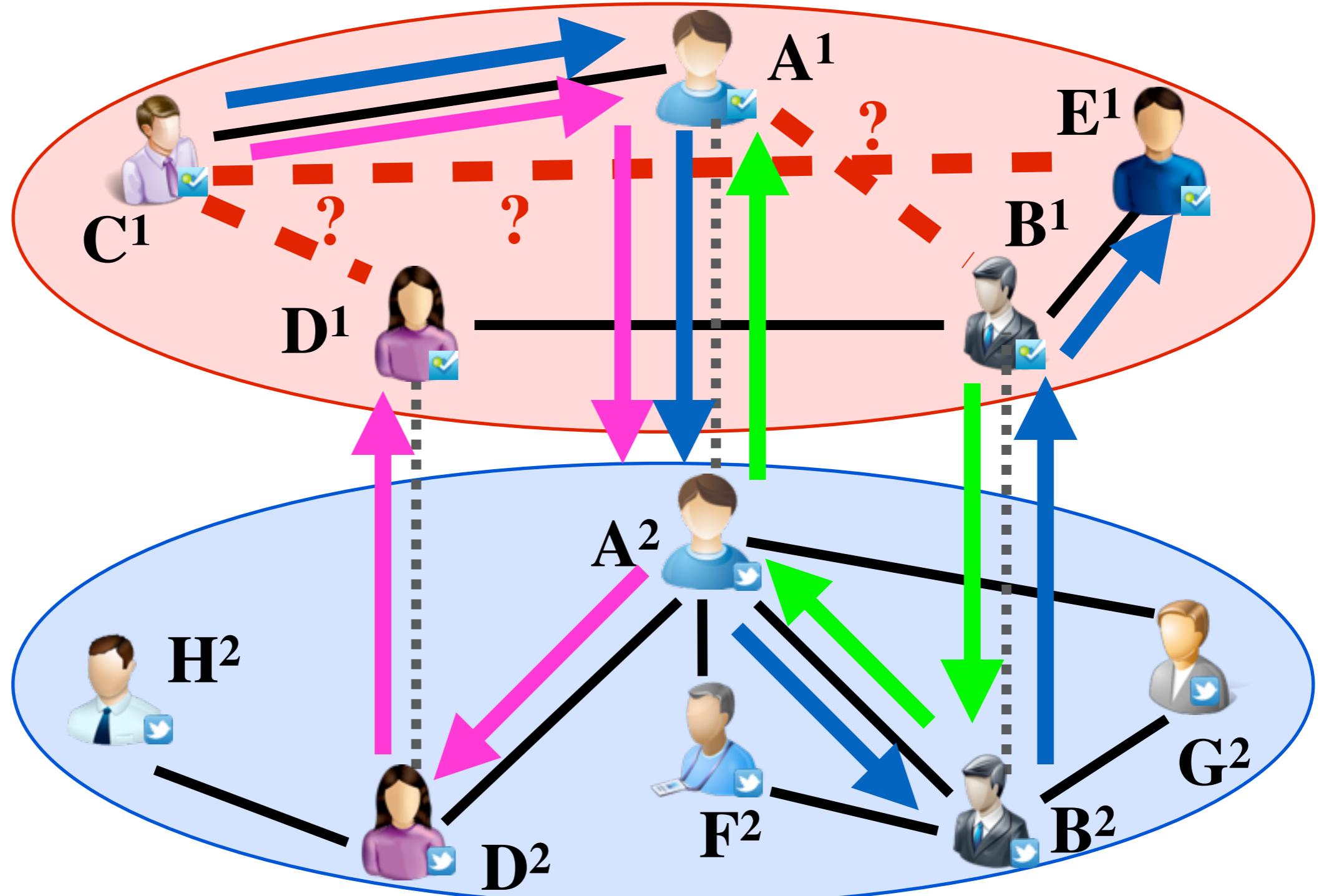
Category 4.: $(\Phi(U^i, U^i) \cup \Phi_0(U^i, U^i)) \circ \Upsilon(U^i, U^j) \circ (\Phi(U^j, U^j) \cup \Phi_0(U^j, U^j)) \circ \Upsilon(U^j, U^i) \circ (\Phi(U^i, U^i) \cup \Phi_0(U^i, U^i))$, whose notation is $\Psi_4(U^i, U^i)$;

Inter-network social meta path instances

Links: anchor link — social link - ? - potential social link

Paths: B¹ to A¹ C¹ to D¹ C¹ to E¹

Network 1 Network 2



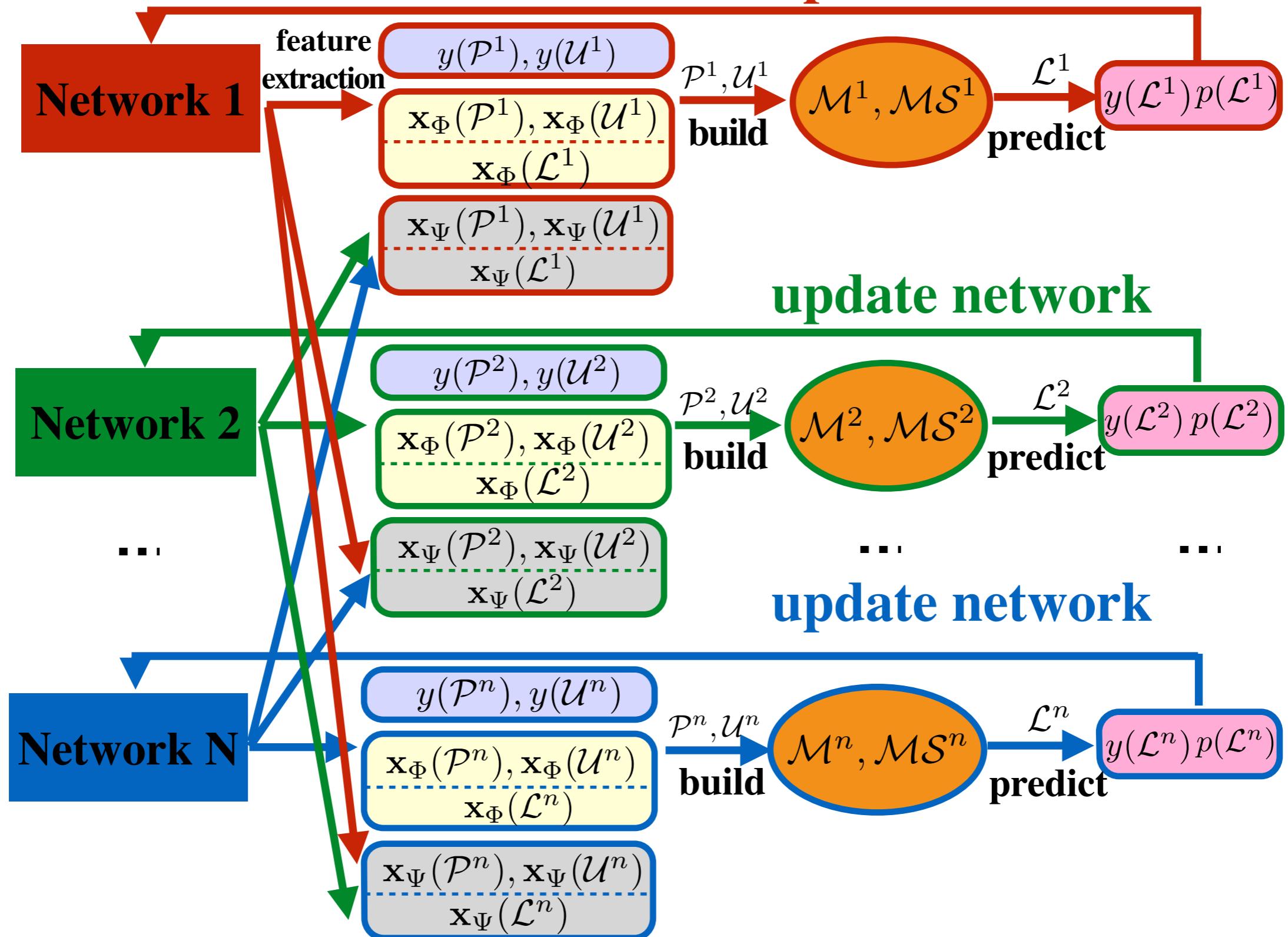
Meta path selection

Let variable $X_i \in [\mathbf{x}_\Phi^T, \mathbf{x}_\Psi^T]^T$ be a feature extracted based on a meta path in $\{\Phi, \Psi\}$ and variable Y be the *label*. $P(Y = y)$ denotes the *prior probability* that links in the training set having label y and $P(X_i = x)$ represents the *frequency* that feature X_i has value x . Information theory related measure *mutual information* (mi) is used as the ranking criteria:

$$mi(X_i) = \sum_x \sum_y P(X_i = x, Y = y) \log \frac{P(X_i = x, Y = y)}{P(X_i = x)P(Y = y)}$$

Multi-network collective link prediction framework

update network



Dataset

- Foursquare and Twitter

Table 2: Properties of the Heterogeneous Networks

		network	
property		Twitter	Foursquare
# node	user	5,223	5,392
	tweet/tip	9,490,707	48,756
	location	297,182	38,921
# link	friend/follow	164,920	76,972
	write	9,490,707	48,756
	locate	615,515	48,756

Experiment Settings

- Ground truth: existing social link among users
 - hide part of the existing links in the test set
 - build model to discover these links
- Comparison Methods
 - MLI (Multi-network Link Ientifier)
 - LI (Link Ientifier): predict links in each network independently
 - SCAN(Supervised Cross-Aligned-Network link prediction): supervised link prediction, no meta path selection,
 - SCAN_s (SCAN with source network): features are extracted based on inter-network meta paths
 - SCAN_t (SCAN with target network): features are extracted based on intra-network meta paths
- Evaluation Metrics
 - AUC, Accuracy, F1

collective link prediction is better than independent link prediction

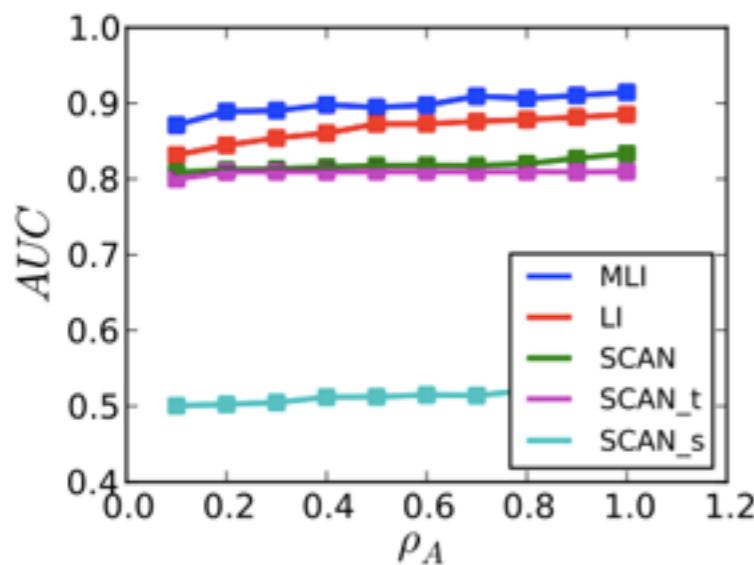
Experiment Results

PU link prediction setting and meta path selection can improve the results

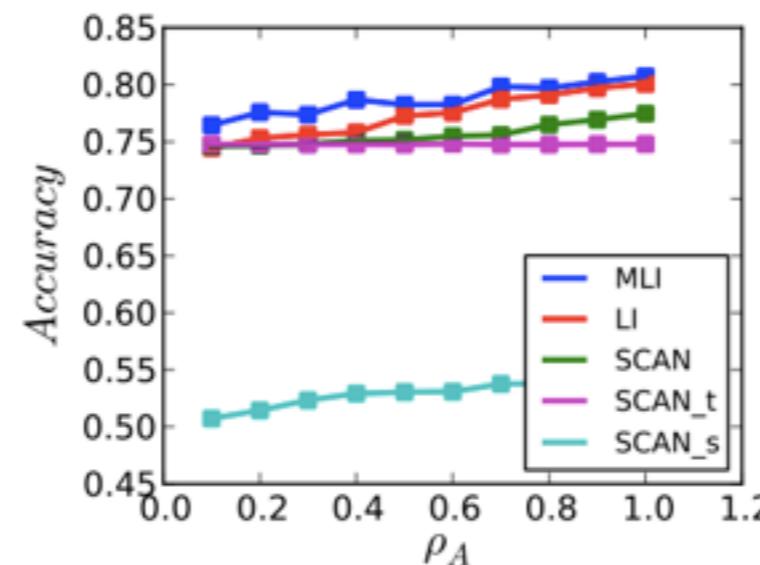
network	measure	methods	MLI	LI	SCAN	SCANT	SCANs	MLI	LI	SCAN	SCANT	SCANs	MLI	LI	SCAN	SCANT	SCANs	MLI	LI	SCAN	SCANT	SCANs	
Foursquare	AUC																						
		MLI	0.524±0.013	0.524±0.017	0.524±0.012	0.524±0.005	0.524±0.002	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	
		LI	0.568±0.013	0.624±0.053	0.699±0.004	0.722±0.006	0.761±0.01	0.568±0.013	0.624±0.053	0.699±0.004	0.722±0.006	0.761±0.01	0.568±0.013	0.624±0.053	0.699±0.004	0.722±0.006	0.761±0.01	0.568±0.013	0.624±0.053	0.699±0.004	0.722±0.006	0.761±0.01	
		SCAN	0.558±0.007	0.6±0.006	0.683±0.071	0.714±0.009	0.721±0.007	0.491±0.019	0.568±0.004	0.66±0.008	0.685±0.007	0.711±0.007	0.491±0.019	0.568±0.004	0.66±0.008	0.685±0.007	0.711±0.007	0.491±0.019	0.568±0.004	0.66±0.008	0.685±0.007	0.711±0.007	
		SCANT	0.491±0.019	0.568±0.004	0.66±0.008	0.685±0.007	0.711±0.007	0.548±0.011	0.548±0.055	0.548±0.007	0.548±0.008	0.548±0.007	0.548±0.011	0.548±0.055	0.548±0.007	0.548±0.008	0.548±0.007	0.548±0.011	0.548±0.055	0.548±0.007	0.548±0.008	0.548±0.007	
		SCANs	0.524±0.013	0.524±0.017	0.524±0.012	0.524±0.005	0.524±0.002	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	
	Accuracy																						
		MLI	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	
		LI	0.568±0.013	0.624±0.053	0.699±0.004	0.722±0.006	0.761±0.01	0.568±0.013	0.624±0.053	0.699±0.004	0.722±0.006	0.761±0.01	0.568±0.013	0.624±0.053	0.699±0.004	0.722±0.006	0.761±0.01	0.568±0.013	0.624±0.053	0.699±0.004	0.722±0.006	0.761±0.01	
		SCAN	0.558±0.007	0.6±0.006	0.683±0.071	0.714±0.009	0.721±0.007	0.491±0.019	0.568±0.004	0.66±0.008	0.685±0.007	0.711±0.007	0.491±0.019	0.568±0.004	0.66±0.008	0.685±0.007	0.711±0.007	0.491±0.019	0.568±0.004	0.66±0.008	0.685±0.007	0.711±0.007	
		SCANT	0.491±0.019	0.568±0.004	0.66±0.008	0.685±0.007	0.711±0.007	0.548±0.011	0.548±0.055	0.548±0.007	0.548±0.008	0.548±0.007	0.548±0.011	0.548±0.055	0.548±0.007	0.548±0.008	0.548±0.007	0.548±0.011	0.548±0.055	0.548±0.007	0.548±0.008	0.548±0.007	
		SCANs	0.524±0.013	0.524±0.017	0.524±0.012	0.524±0.005	0.524±0.002	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	0.622±0.01	0.692±0.00	0.70±0.005	0.769±0.004	0.779±0.000	
	F1																						
		MLI	0.622±0.01	0.695±0.02	0.722±0.018	0.742±0.005	0.771±0.005	0.622±0.01	0.695±0.02	0.722±0.018	0.742±0.005	0.771±0.005	0.622±0.01	0.695±0.02	0.722±0.018	0.742±0.005	0.771±0.005	0.622±0.01	0.695±0.02	0.722±0.018	0.742±0.005	0.771±0.005	
		LI	0.63±0.017	0.635±0.015	0.66±0.007	0.684±0.01	0.715±0.016	0.531±0.096	0.559±0.002	0.55±0.016	0.584±0.002	0.6±0.011	0.531±0.096	0.559±0.002	0.55±0.016	0.584±0.002	0.6±0.011	0.531±0.096	0.559±0.002	0.55±0.016	0.584±0.002	0.6±0.011	
		SCAN	0.6±0.02	0.609±0.006	0.614±0.031	0.632±0.018	0.645±0.018	0.531±0.096	0.559±0.002	0.55±0.016	0.584±0.002	0.6±0.011	0.531±0.096	0.559±0.002	0.55±0.016	0.584±0.002	0.6±0.011	0.531±0.096	0.559±0.002	0.55±0.016	0.584±0.002	0.6±0.011	
		SCANT	0.531±0.096	0.559±0.002	0.55±0.016	0.584±0.002	0.6±0.011	0.56±0.041	0.56±0.015	0.56±0.015	0.56±0.015	0.56±0.013	0.531±0.096	0.559±0.002	0.55±0.016	0.584±0.002	0.6±0.011	0.531±0.096	0.559±0.002	0.55±0.016	0.584±0.002	0.6±0.011	
		SCANs	0.56±0.016	0.56±0.041	0.56±0.015	0.56±0.015	0.56±0.013	0.56±0.016	0.56±0.015	0.56±0.015	0.56±0.015	0.56±0.013	0.56±0.016	0.56±0.015	0.56±0.015	0.56±0.015	0.56±0.013	0.56±0.016	0.56±0.015	0.56±0.015	0.56±0.015	0.56±0.013	

Parameter Analysis

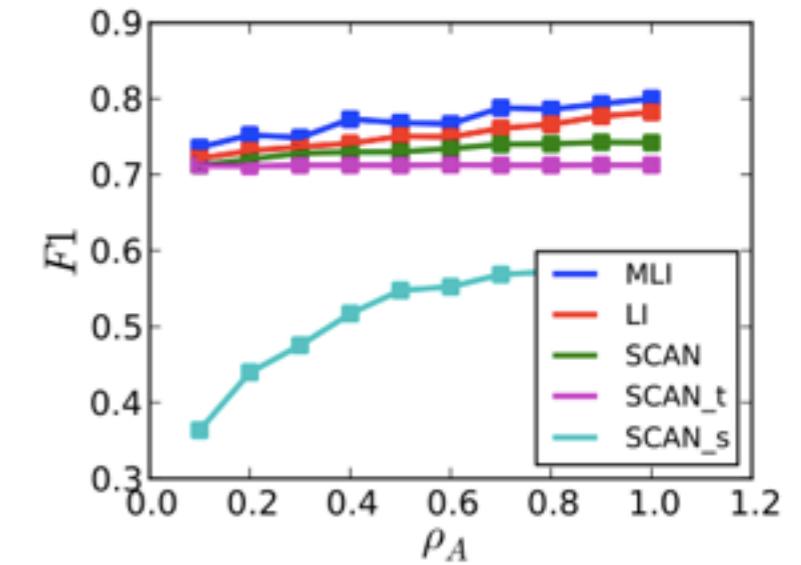
- ratio of anchor links



(a) Foursquare-AUC



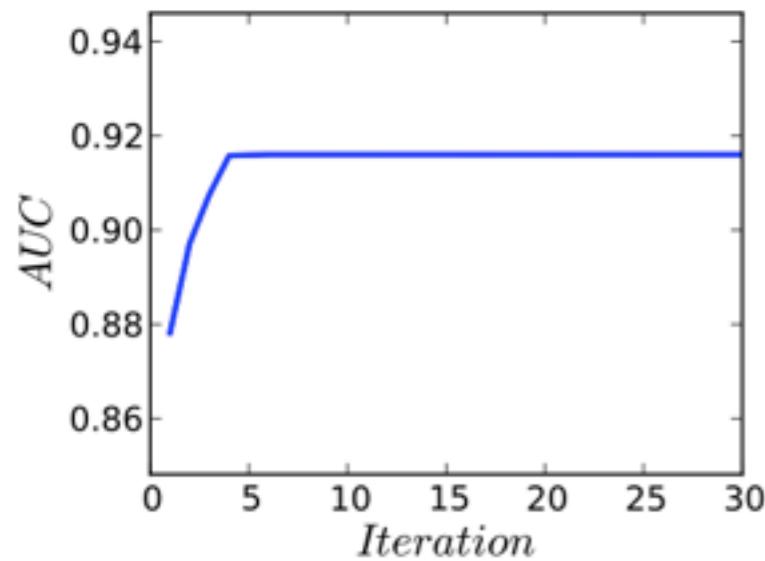
(b) Foursquare-Acc.



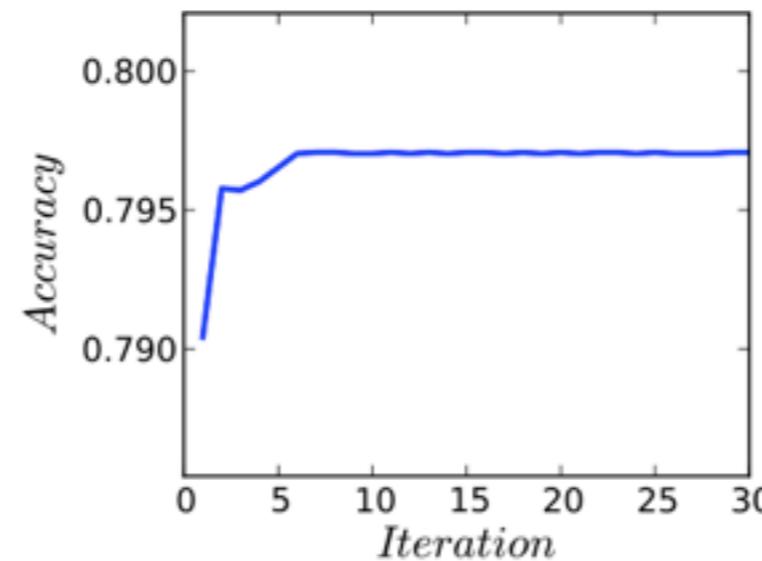
(c) Foursquare-F1

the more anchor links we have, the better performance we can achieve

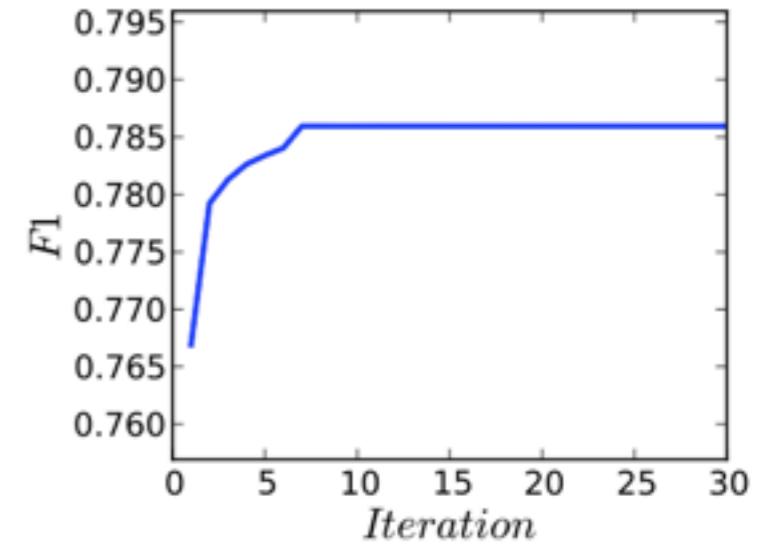
Convergence Analysis



(a) Foursquare-AUC



(b) Foursquare-Acc.



(c) Foursquare-F1

converge quickly, in less than 10 iterations

Conclusions

- Problem studied: collective link prediction across multiple aligned social networks
- Proposed Method:
 - PU Link Prediction Setting
 - Intra-network & Inter-network Meta Path based Feature Extraction
 - Meta path selection
 - Multi-network Collective PU Link Prediction Framework
- Experiment Results:
 - Collective Link Prediction is better than Independent Link Prediction
 - PU Link Prediction & Meta Path Selection can improve the results
 - Using information across networks can achieve better results
 - MLI can perform well consistently for different anchor link ratios & can converge quickly

Q&A