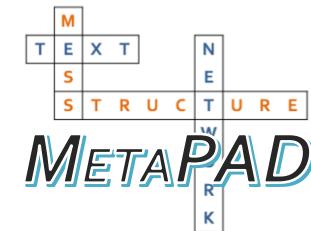
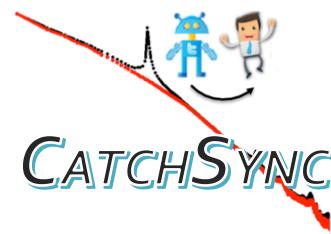


# Data-Driven Behavioral Analytics with Networks



Meng Jiang  
University of Notre Dame

<http://www.meng-jiang.com>

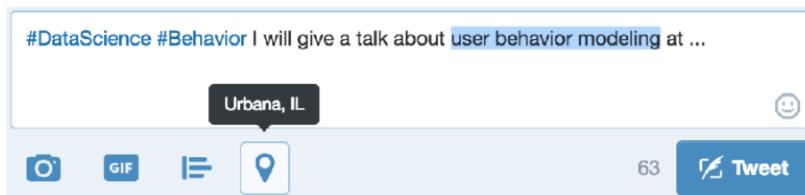
# Behavior and “Behavior Networks”

Interaction of individuals with themselves and with their environment. — Wikipedia

Social behaviors



Tweeting behaviors



Paper-publishing behaviors

**Meng Jiang, Christos Faloutsos, and Jiawei Han.** “CatchTartan: Representing and Summarizing Dynamic Behaviors.” In **SIGKDD 2016**.

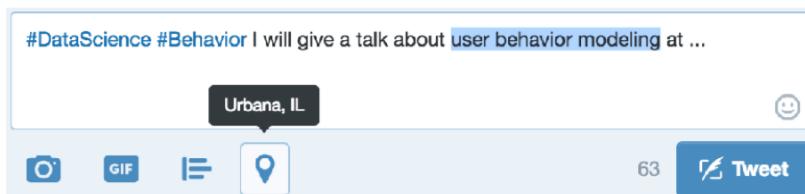
# Behavior and “Behavior Networks”

Interaction of individuals with themselves and with their environment. — Wikipedia

Social behaviors

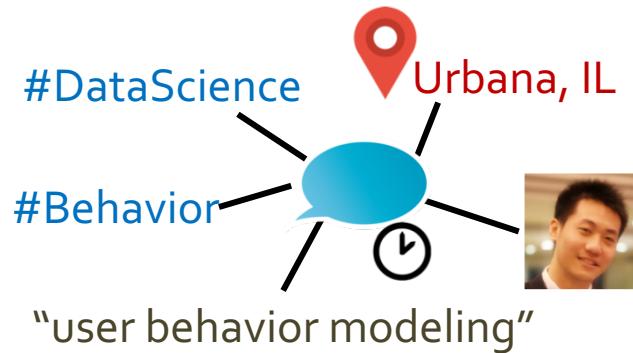
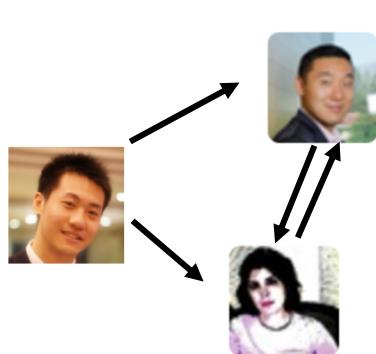


Tweeting behaviors



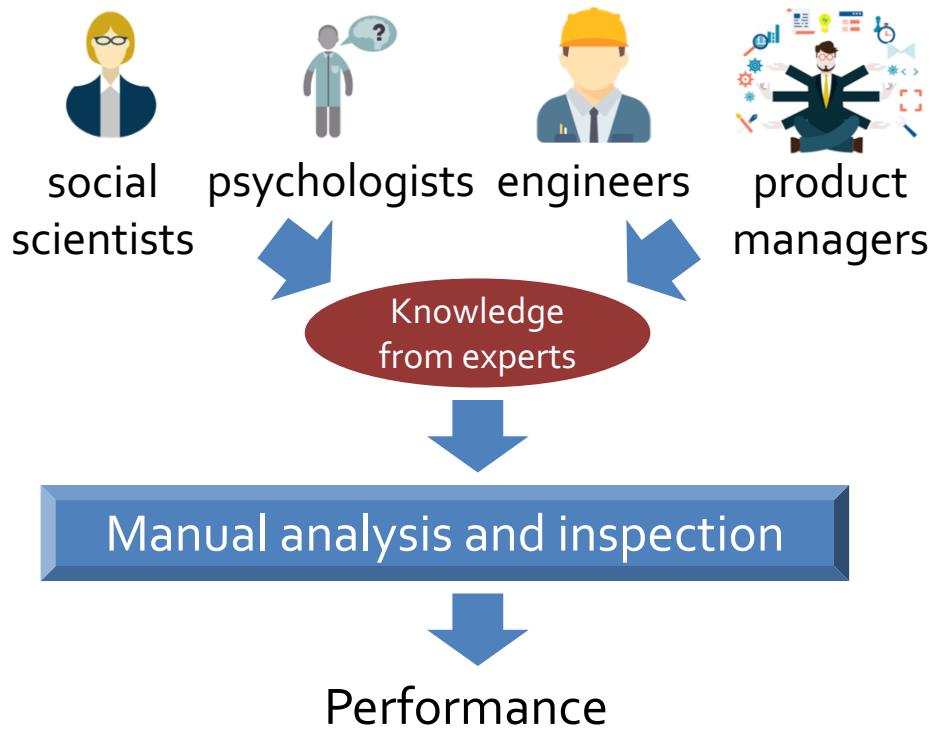
Paper-publishing behaviors

Meng Jiang, Christos Faloutsos, and Jiawei Han. “CatchTartan: Representing and Summarizing Dynamic Behaviors.” In SIGKDD 2016.

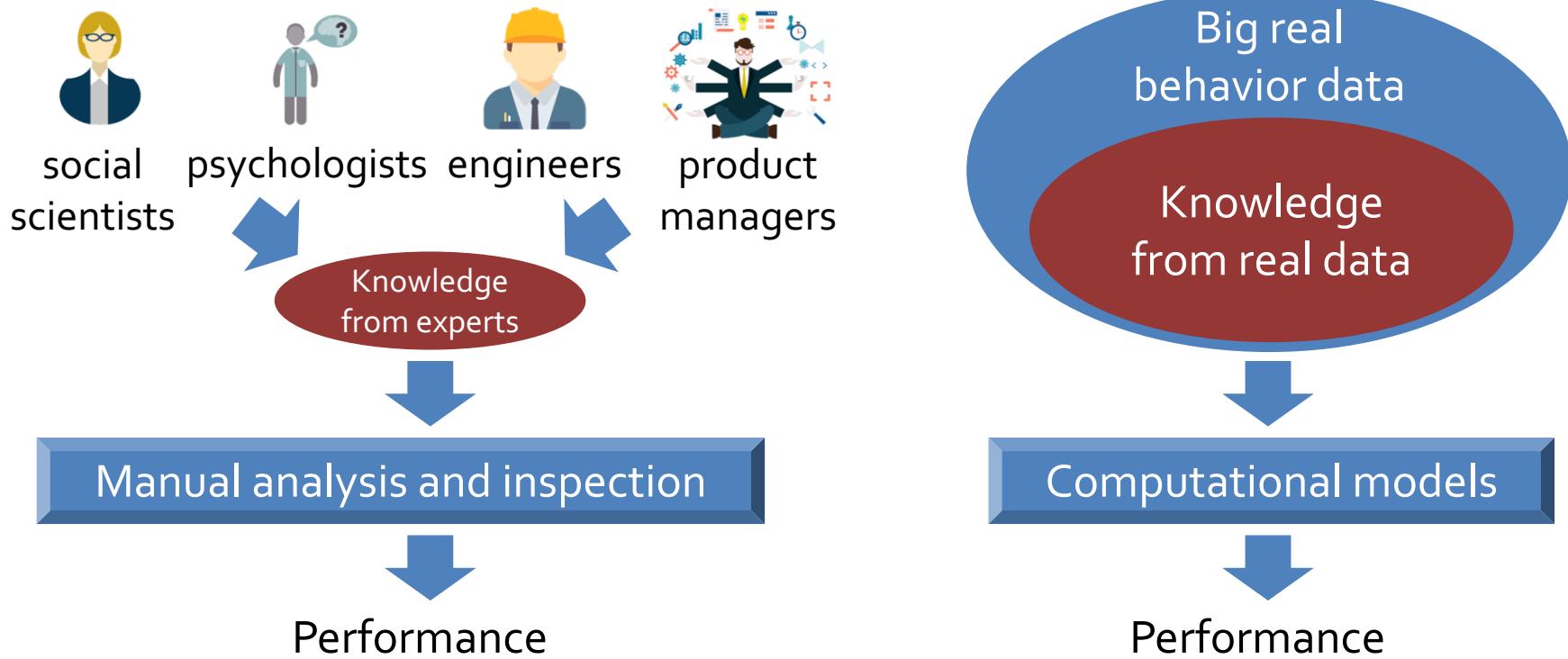


Applications: prediction, recommendation, fraud detection, spam detection...

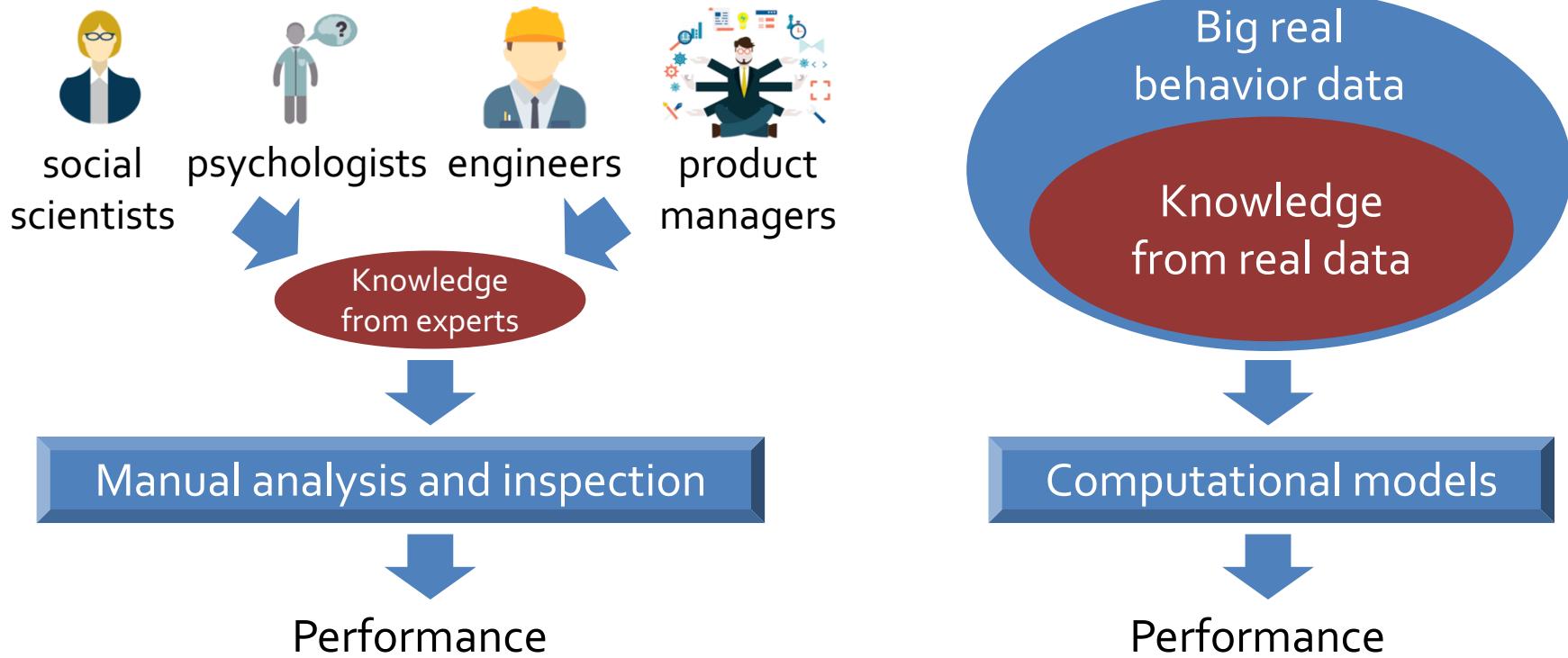
# Data-Driven Behavioral Analysis



# Data-Driven Behavioral Analysis



# Data-Driven Behavioral Analysis



Data, **knowledge**, intelligence, and trustworthiness.  
(User behavior modeling)



# Basic Research Areas

- Six Disruptive Basic Research Areas
  - Engineered Materials (metamaterials and plasmonics)
  - Quantum Information and Control
  - Cognitive Neuroscience
  - Nanoscience and Nanoengineering
  - Synthetic Biology
  - Computational Modeling of Human and Social Behavior

# Research Topics in Behavior Modeling

Behavior  
Modeling

# Research Topics in Behavior Modeling

REWARDS	# TICKETS GIVEN	CONSEQUENCES	# TICKETS TAKEN AWAY
Extra Math	+5	HITTING	-3
Getting along well with others	+3	BULLYING	-4
Good Table Manners	+4	TEASING	-1
LOVE & RESPECT	+5	LYING	-2
Obeying the FIRST TIME	+3	THROWING A FIT	-3
Calm & Quiet in STORE	+3	Ignoring Parents	-4
Extra Reading	+2	SCREAMING or YELLING	-1
CLEANING up after PLAYING	+2	BAD SPORT	-2

## 1. Behavior intentions

Behavior  
Modeling

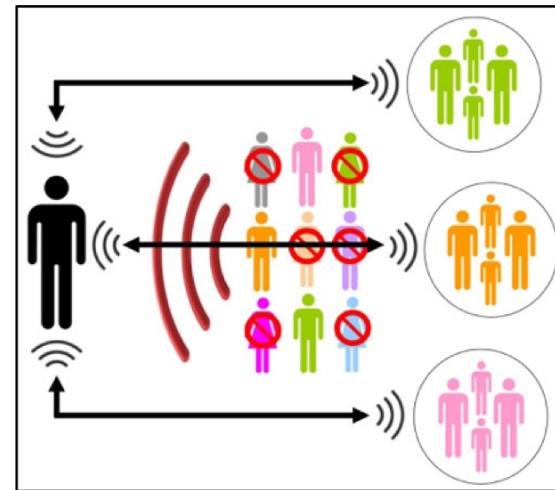
# Research Topics in Behavior Modeling

REWARDS	# TICKETS GIVEN	CONSEQUENCES	# TICKETS TAKEN AWAY
Extra Math	+5	HITTING	-3
Getting along well with others	+3	BULLYING	-4
Good Table Manners	+4	TEASING	-1
LOVE & RESPECT	+5	LYING	-2
Obeying the FIRST TIME	+3	THROWING A FIT	-3
Calm & Quiet in STORE	+3	Ignoring Parents	-4
Extra Reading	+2	SCREAMING or YELLING	-1
CLEANING up after PLAYING	+2	BAD SPORT	-2

1. Behavior intentions

2. Social contexts

Behavior  
Modeling



# Research Topics in Behavior Modeling

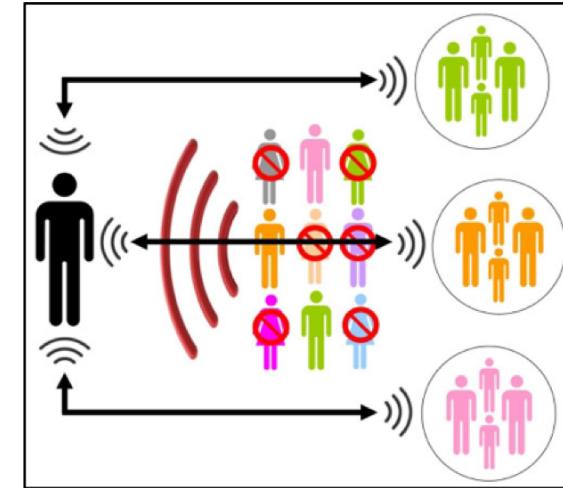
REWARDS	# TICKETS GIVEN	CONSEQUENCES	# TICKETS TAKEN AWAY
Extra Math	+5	HITTING	-3
Getting along well with others	+3	BULLYING	-4
Good Table Manners	+4	TEASING	-1
LOVE & RESPECT	+5	LYING	-2
Obeying the FIRST TIME	+3	THROWING A FIT	-3
Calm & Quiet in STORE	+3	Ignoring Parents	-4
Extra Reading	+2	SCREAMING or TELLING	-1
CLEANING up after PLAYING	+2	BAD SPORT	-2

1. Behavior intentions

2. Social contexts

3. Spatiotemporal contexts

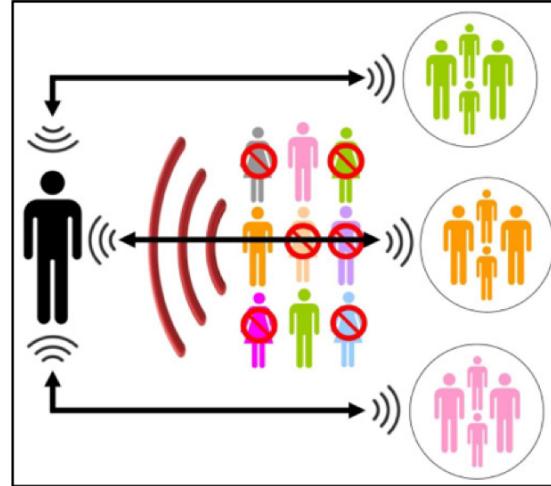
Behavior  
Modeling



# Research Topics in Behavior Modeling

REWARDS	# TICKETS GIVEN	CONSEQUENCES	# TICKETS TAKEN AWAY
Extra Math	+5	HITTING	-3
Getting along well with others	+3	BULLYING	-4
Good Table Manners	+4	TEASING	-1
LOVE & RESPECT	+5	LYING	-2
Obeying the FIRST TIME	+3	THROWING A FIT	-3
Calm & Quiet in STORE	+3	Ignoring Parents	-4
Extra Reading	+2	SCREAMING or YELLING	-1
CLEANING up after PLAYING	+2	BAD SPORT	-2

1. Behavior intentions

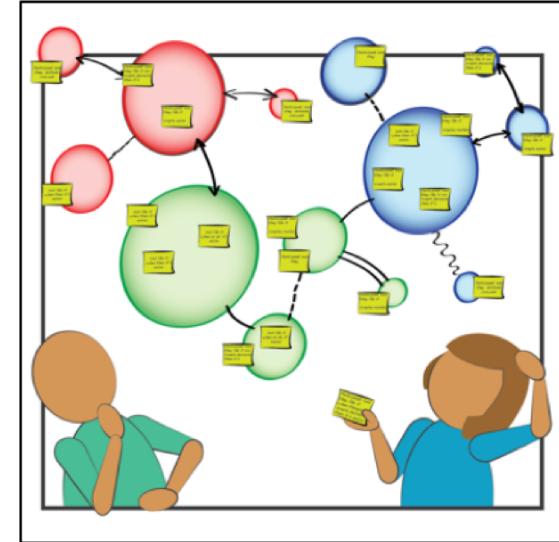


2. Social contexts

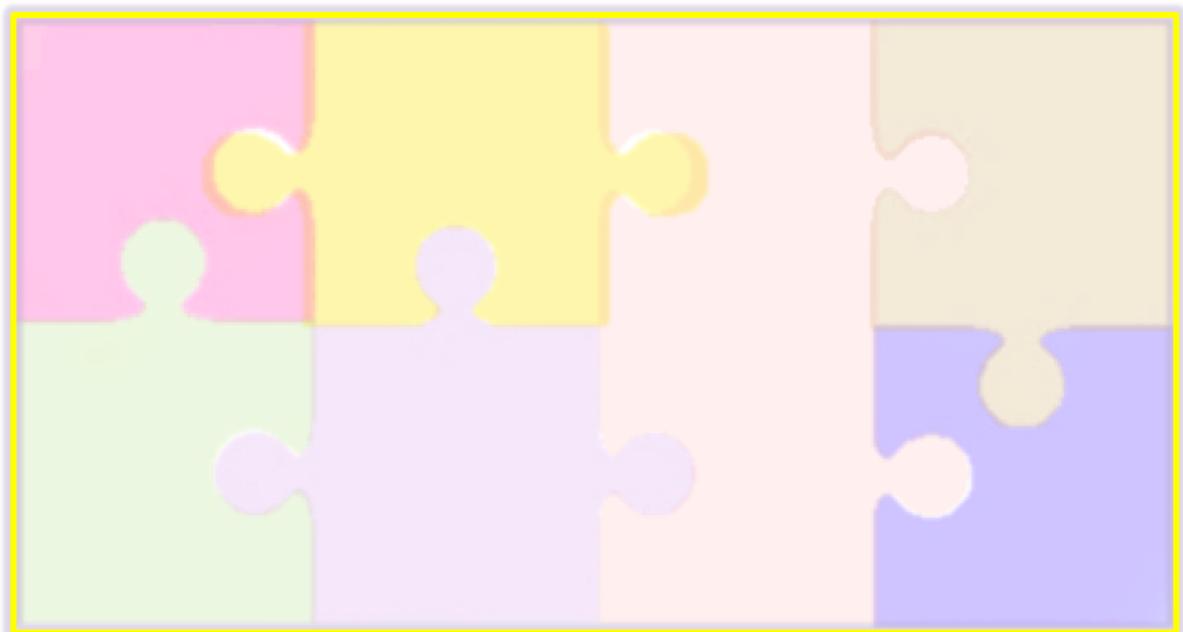
3. Spatiotemporal contexts



4. Behavior content



# My Research “Area”



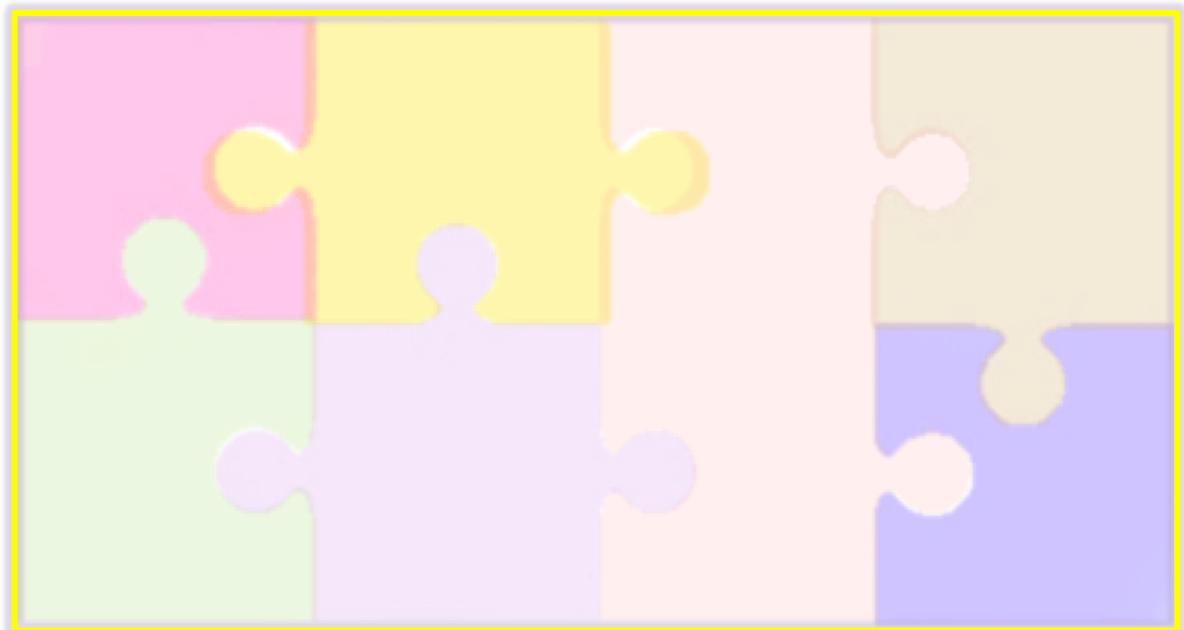
# My Research “Area”

*Intelligence:*

Behavior prediction  
and recommendation

*Trustworthiness:*

Suspicious behavior  
detection

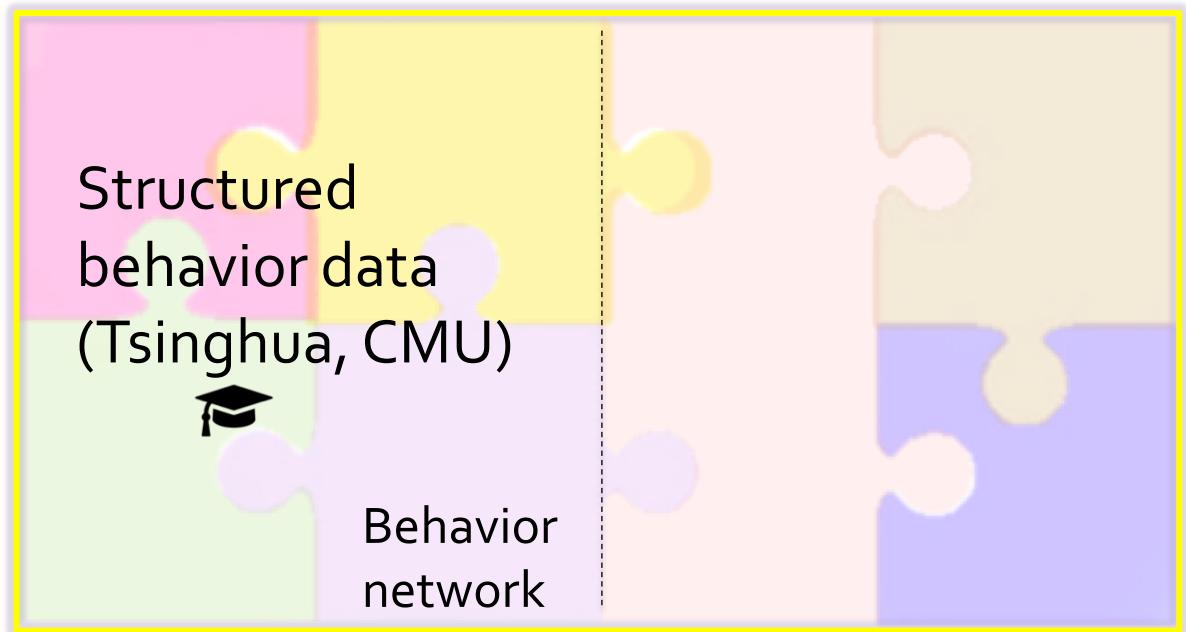


# My Research “Area”

Social      Spatiotemporal  
contexts    contexts

*Intelligence:*  
Behavior prediction  
and recommendation

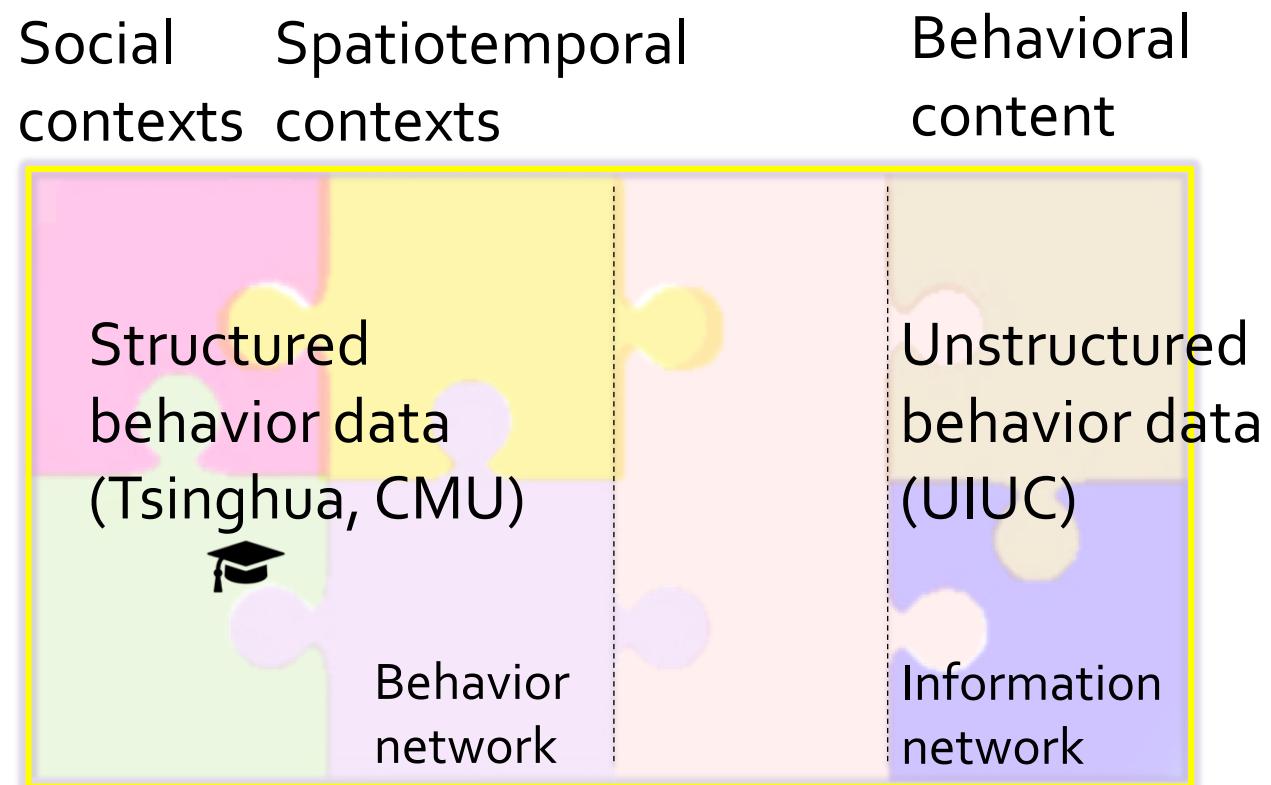
*Trustworthiness:*  
Suspicious behavior  
detection



# My Research “Area”

*Intelligence:*  
Behavior prediction  
and recommendation

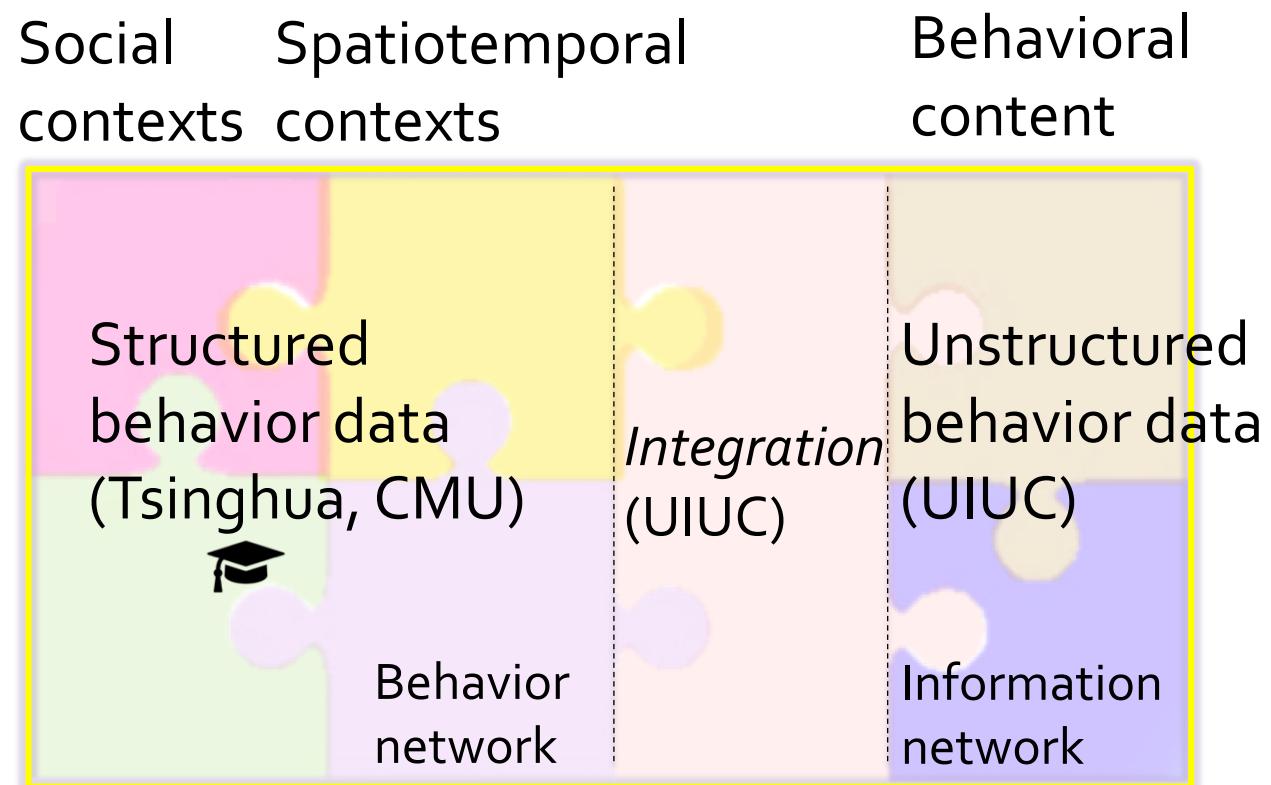
*Trustworthiness:*  
Suspicious behavior  
detection



# My Research “Area”

*Intelligence:*  
Behavior prediction  
and recommendation

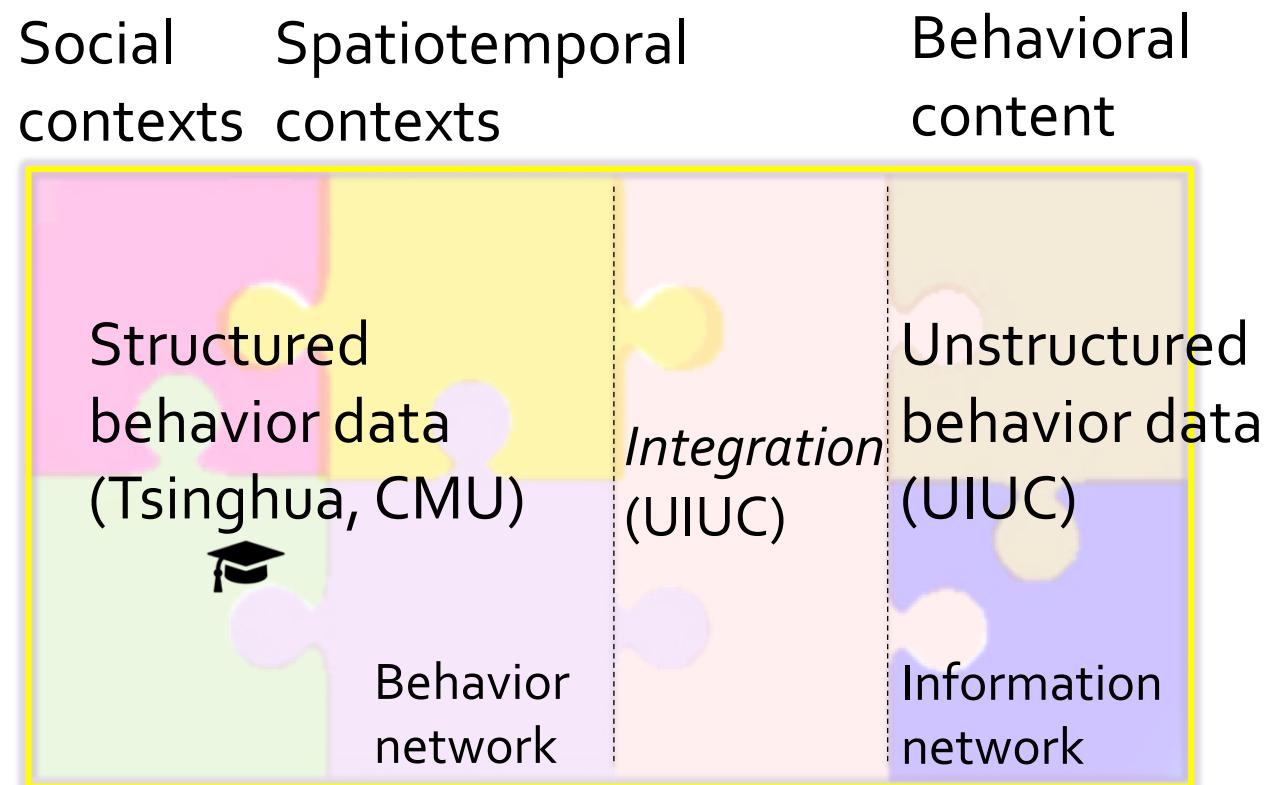
*Trustworthiness:*  
Suspicious behavior  
detection



# My Research “Area”

*Intelligence:*  
Behavior prediction  
and recommendation

*Trustworthiness:*  
Suspicious behavior  
detection



Ask good **Questions**.  
Find good Data-Driven **Methodologies**.  
Propose good **Solutions**.

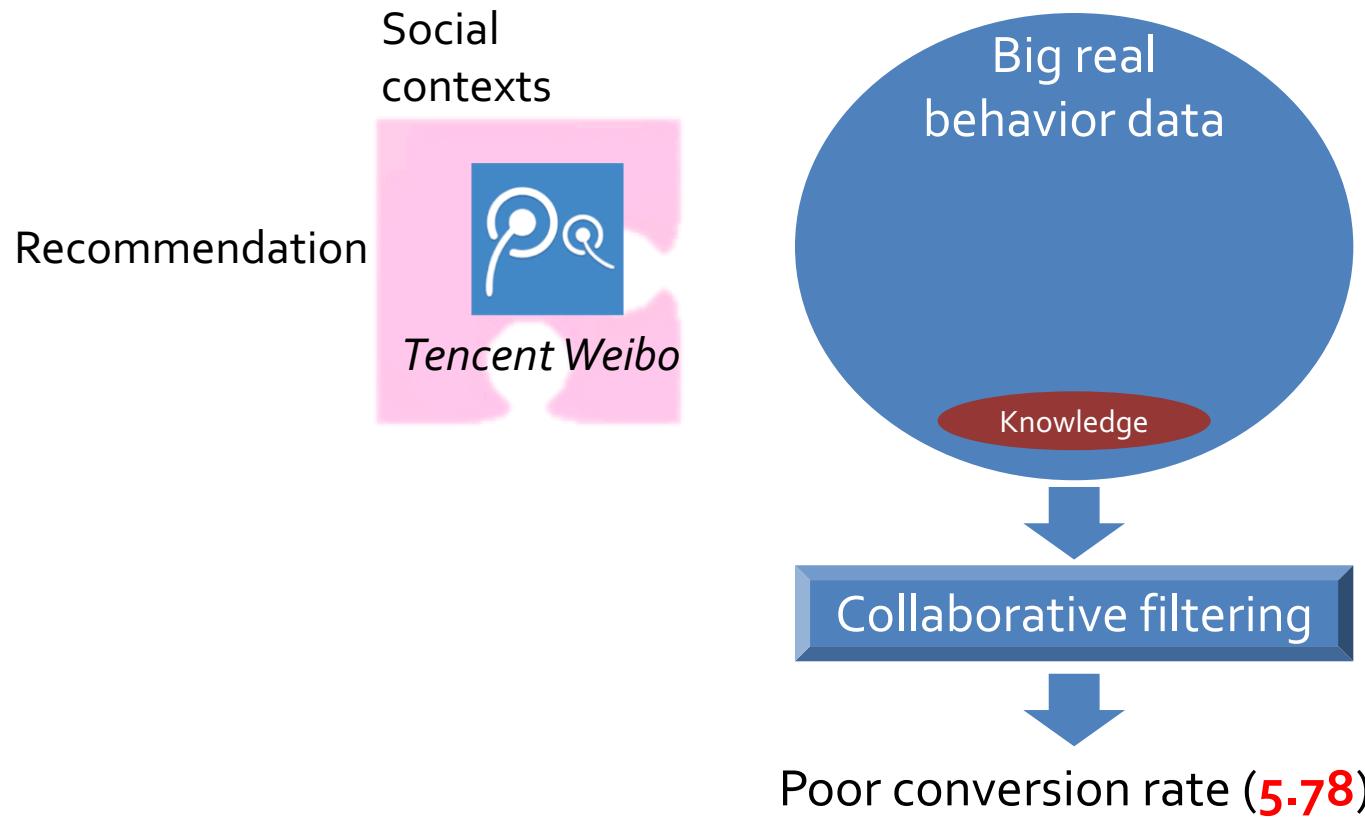
# Q1: A Social Recommender System

Social  
contexts

Recommendation

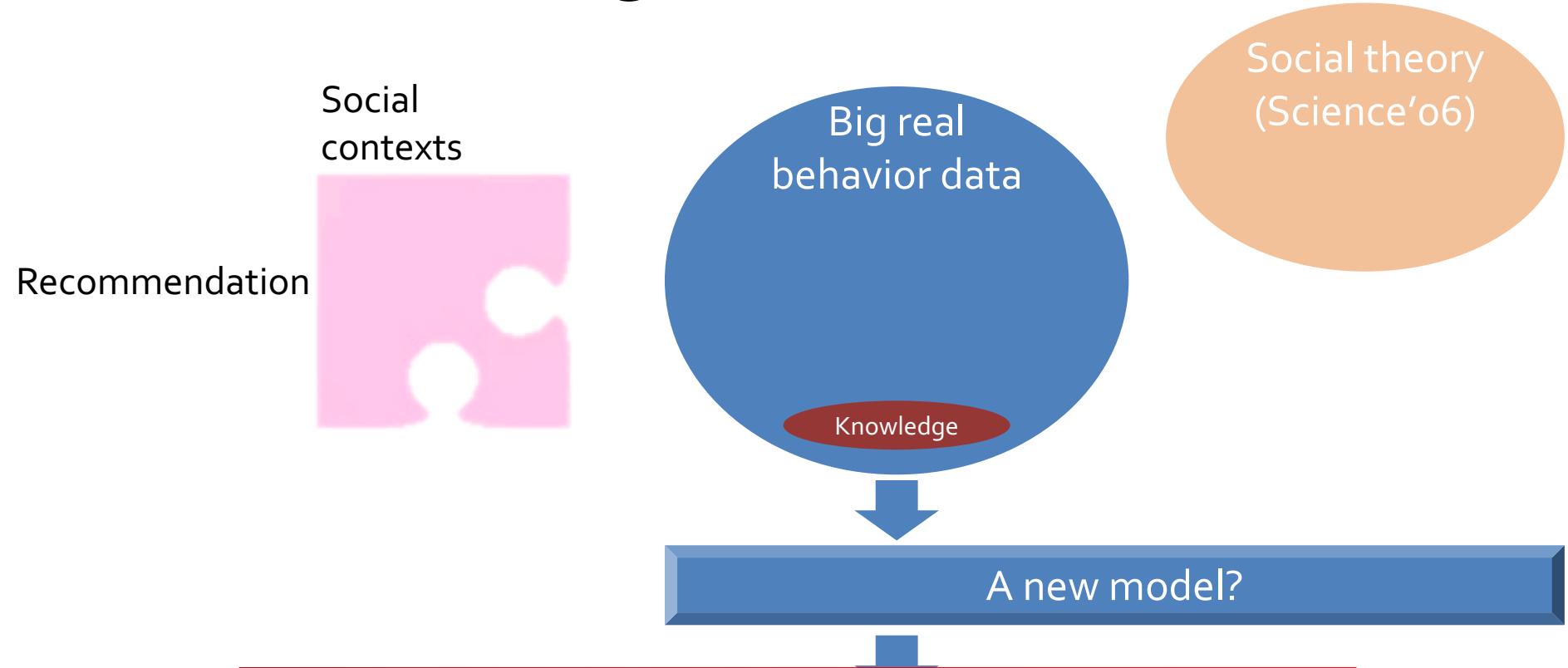


# Q1: A Social Recommender System



Collaborative filtering for recommenders [Breese et al. UAI'98; Getoor and Sahami WEBKDD'99; Herlocker et al. CSCW'00, TOIS'04; Koren et al. KDD'08 Computer'09; Liu et al. SIGIR'08]

# M1: Knowledge from Social Theories

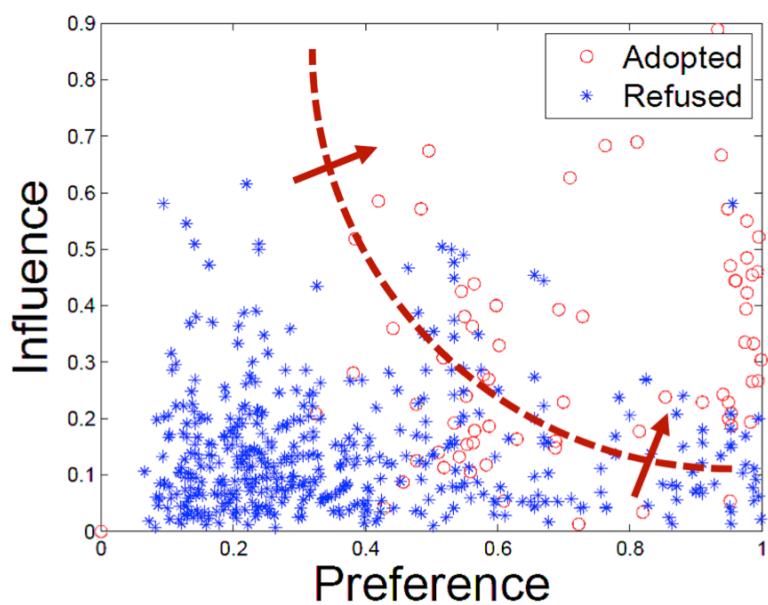
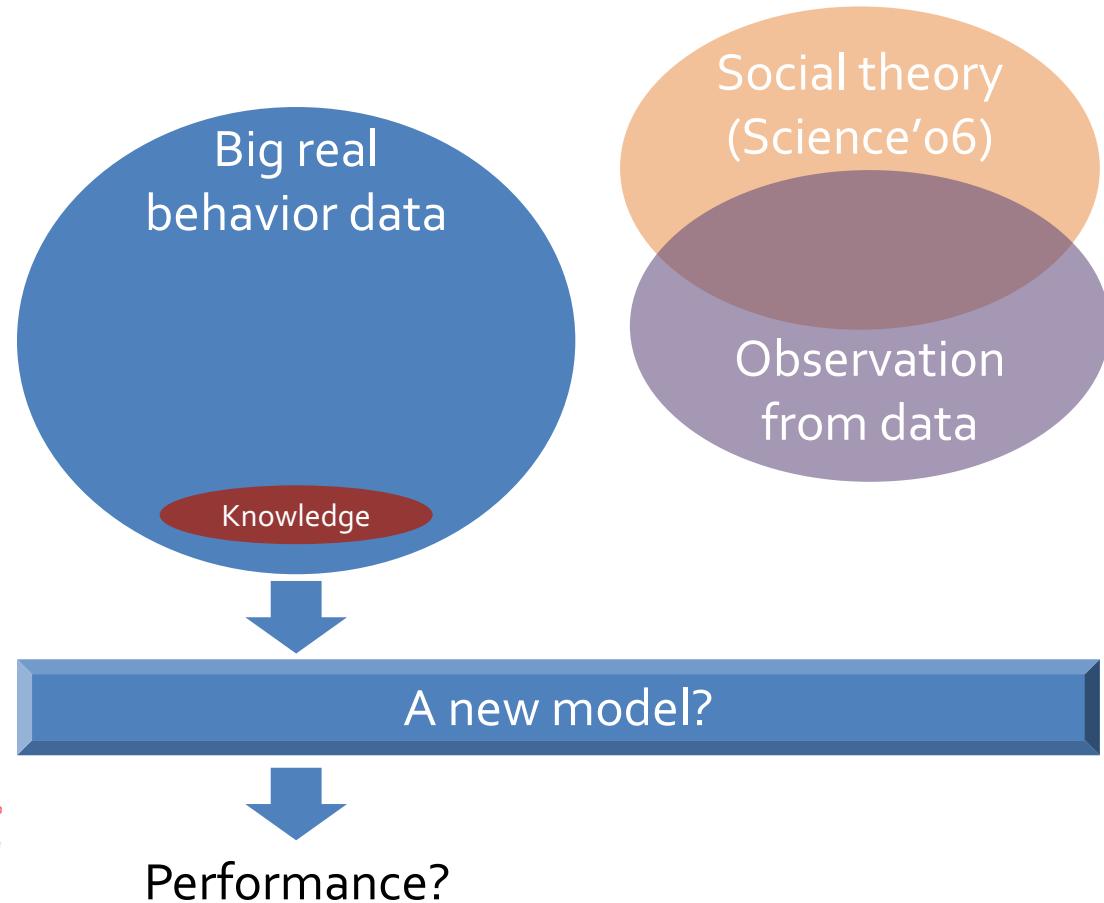


“... Differences between independent and social influence conditions are significant ...”

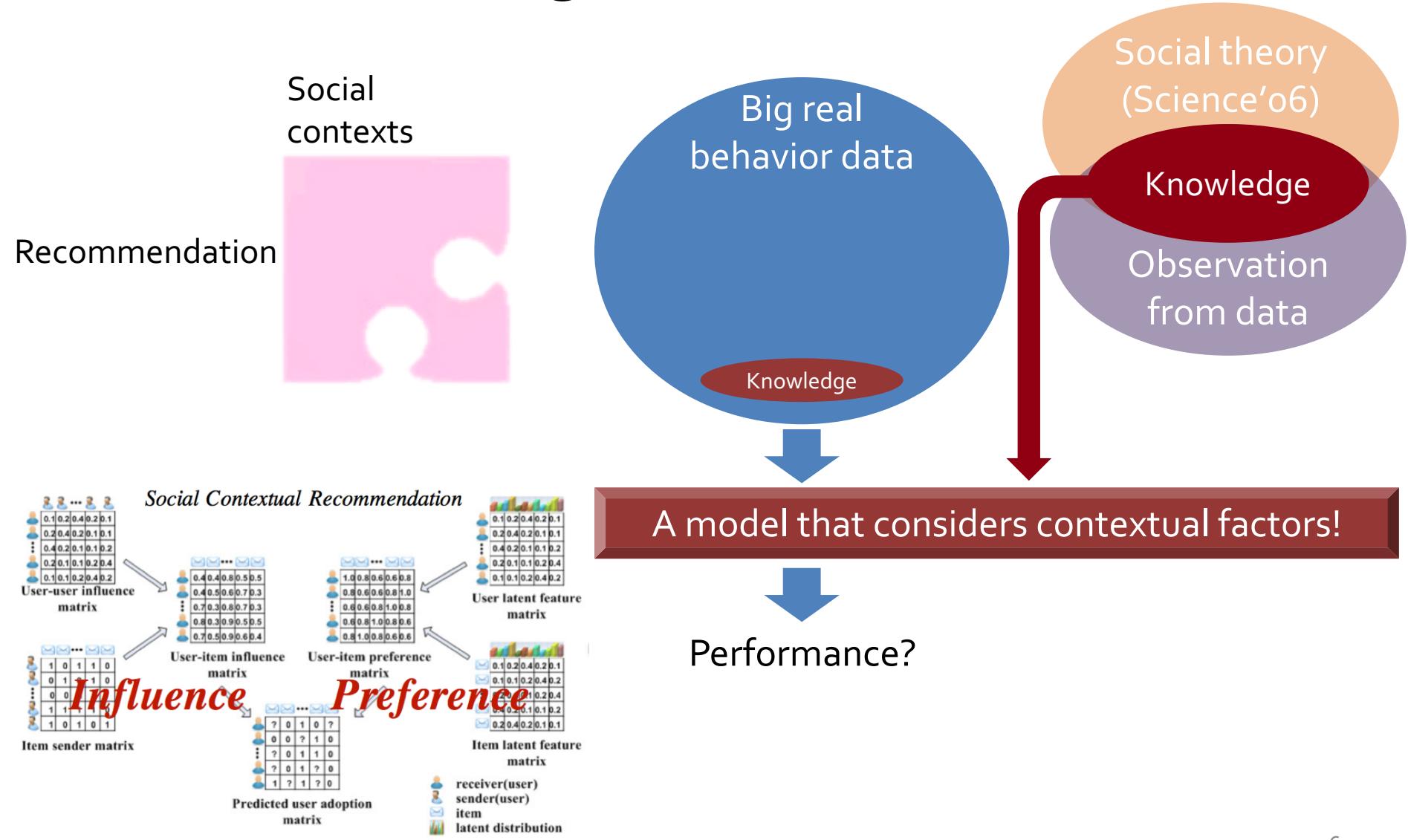
Salganik, Dodds, and Watts. “Experimental Study of Inequality and Unpredictability in an Artificial Cultural Market.” *Science*, Vol. 311, 2006.

# M1: Knowledge from Social Theories

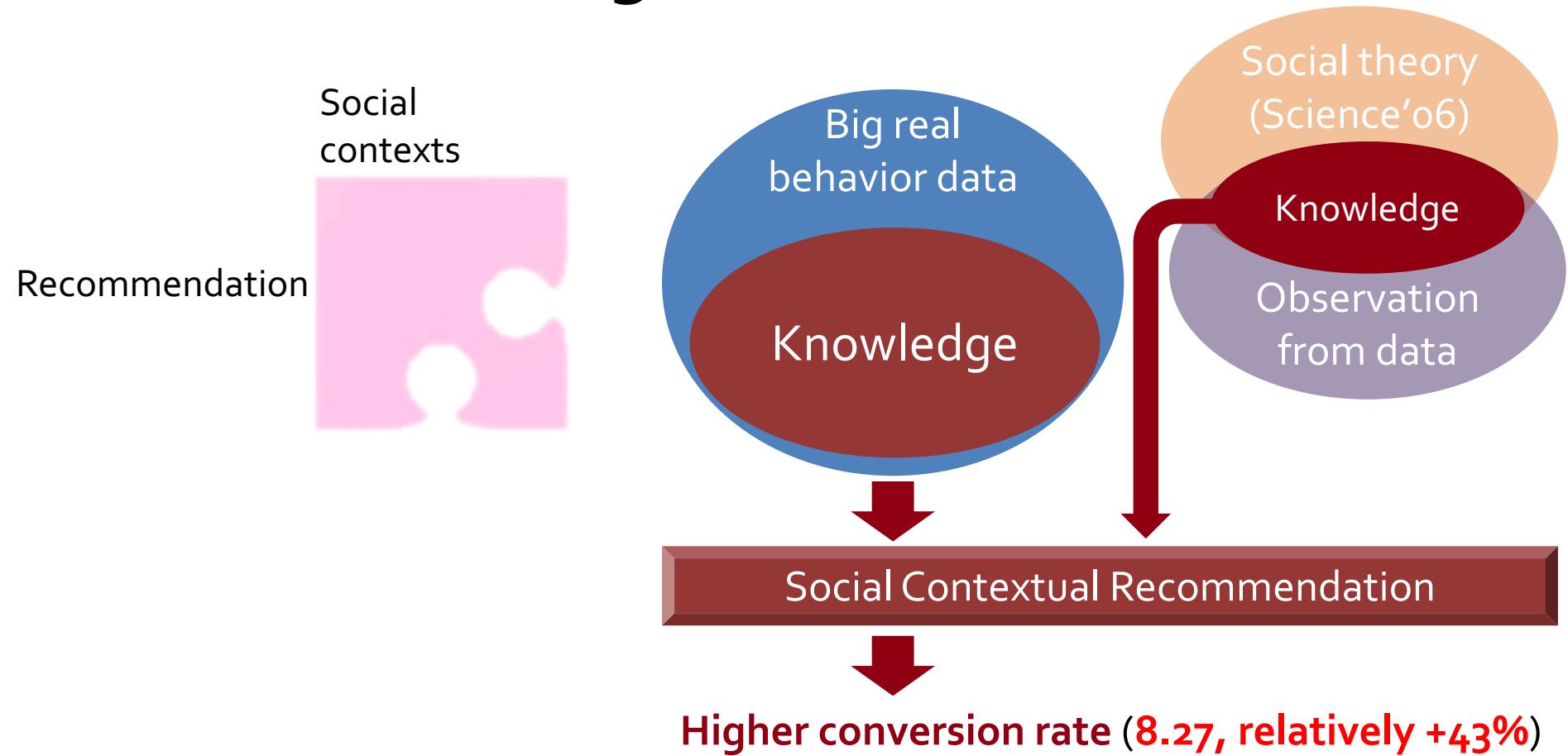
Social contexts  
Recommendation



# M1: Knowledge from Social Theories



# M1: Knowledge from Social Theories



S1: CIKM'12: cited by 211; TKDE'14: cited by 82.

# S1: Social Contextual Recommendation

- Optimization
- Gradient descent method

$$P(\mathbf{R}|\mathbf{S}, \mathbf{U}, \mathbf{V}, \sigma_R^2) = \prod_{i=1}^M \prod_{j=1}^N \mathcal{N}(\mathbf{R}_{ij} | \mathbf{S}_i \mathbf{G}_j^\top \odot \mathbf{U}_i^\top \mathbf{V}_j, \sigma_R^2)$$

behavior      influence      preference

$$\mathcal{J} = \|\mathbf{R} - \mathbf{SG}^\top \odot \mathbf{U}^\top \mathbf{V}\|_F^2 + \alpha \|\mathbf{W} - \mathbf{U}^\top \mathbf{U}\|_F^2$$

$$+ \beta \|\mathbf{C} - \mathbf{V}^\top \mathbf{V}\|_F^2 + \gamma \|\mathbf{S} - \mathbf{F}\|_F^2$$

$$+ \delta \|\mathbf{S}\|_F^2 + \eta \|\mathbf{U}\|_F^2 + \lambda \|\mathbf{V}\|_F^2$$

behavior      user-user interaction

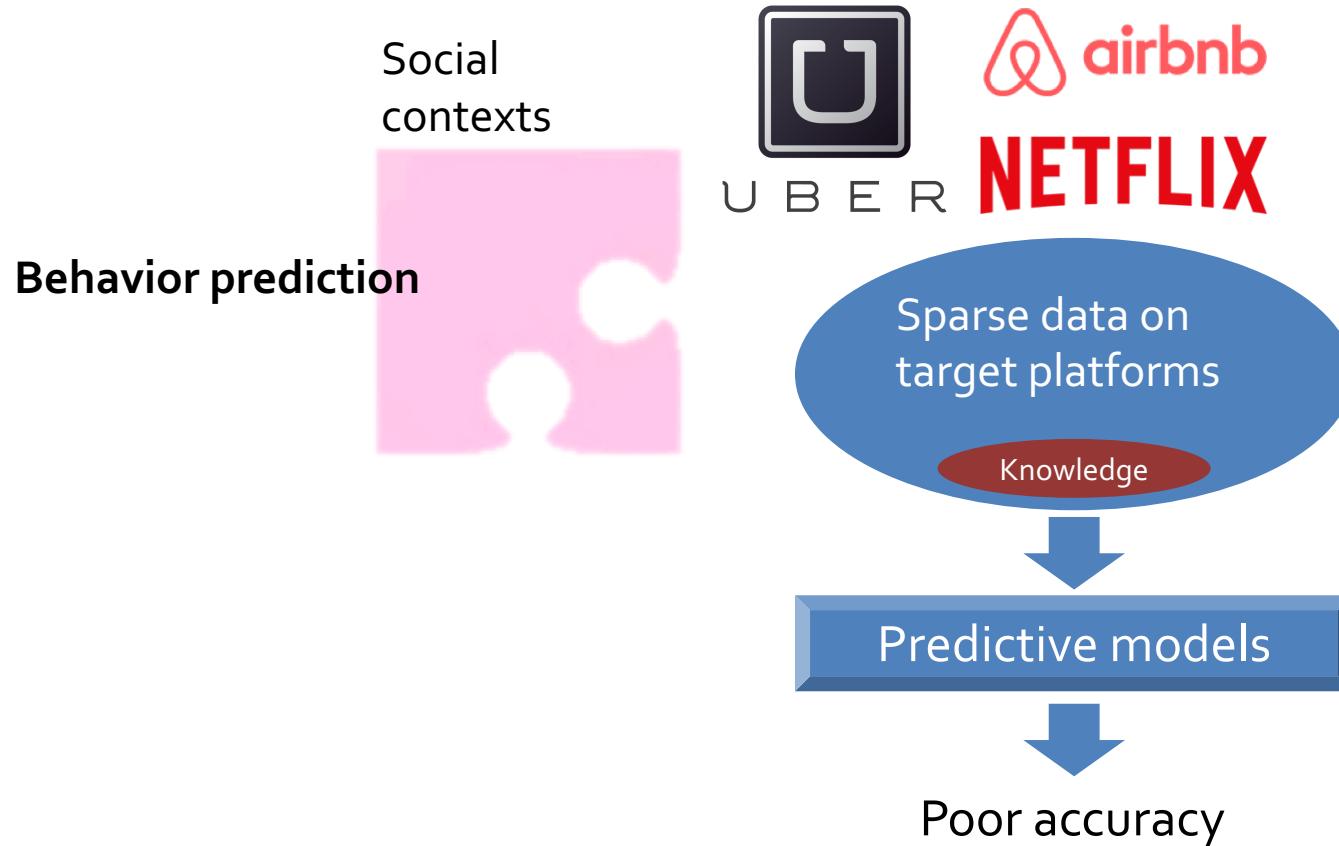
item content      avoid overfitting      social relation

$$\frac{\partial \mathcal{J}}{\partial \mathbf{S}} = 2 \left( -\mathbf{R}(\mathbf{G} \odot \mathbf{V}^\top \mathbf{U}) + (\mathbf{S}\mathbf{G}^\top \odot \mathbf{U}^\top \mathbf{V})\mathbf{G} + \gamma(\mathbf{S} - \mathbf{F}) + \delta\mathbf{S} \right)$$

$$\frac{\partial \mathcal{J}}{\partial \mathbf{U}} = 2 \left( -\mathbf{VR}^\top + \mathbf{V}(\mathbf{GS}^\top \odot \mathbf{V}^\top \mathbf{U}) - 2\alpha\mathbf{UW} + 2\alpha\mathbf{UU}^\top \mathbf{U} + \eta\mathbf{U} \right)$$

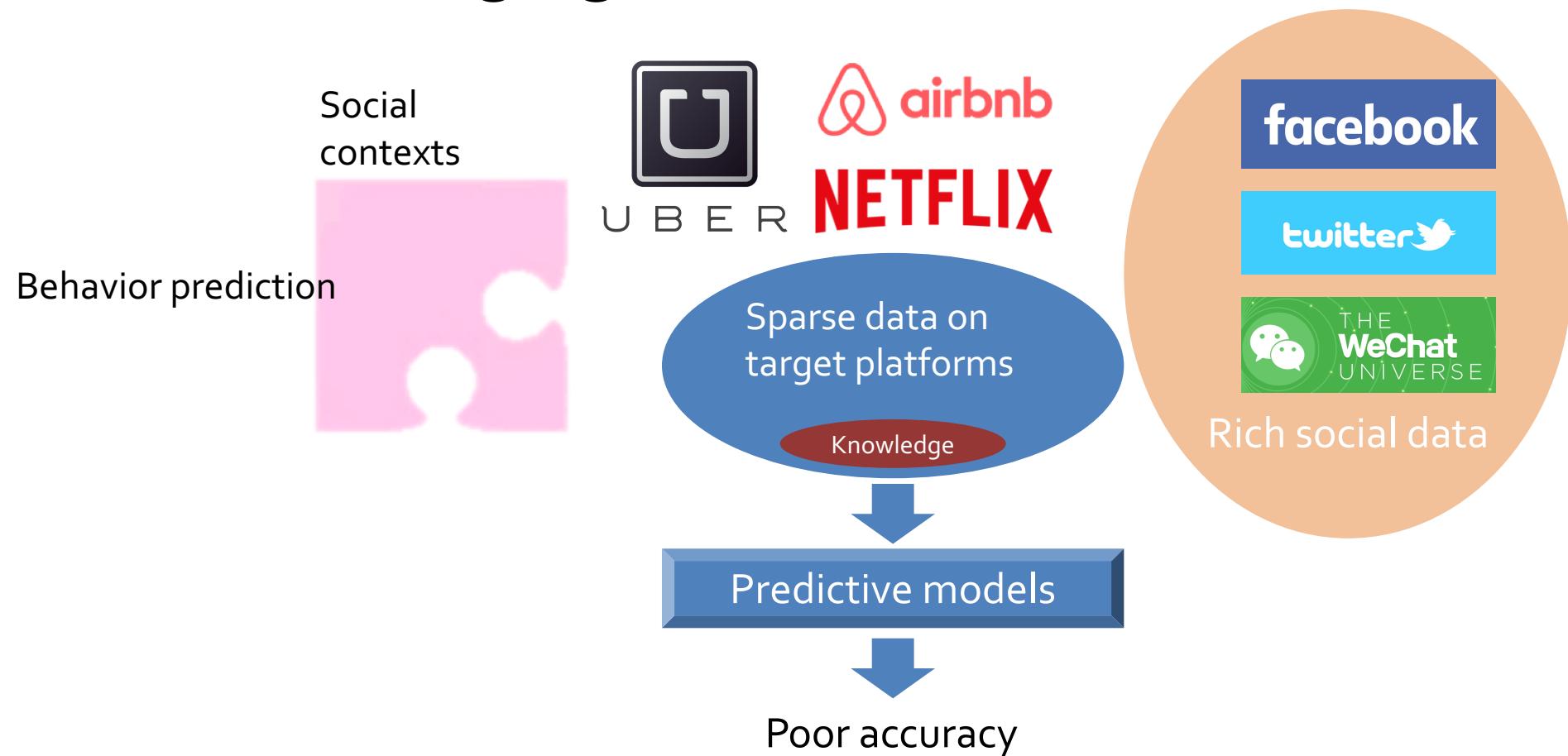
$$\frac{\partial \mathcal{J}}{\partial \mathbf{V}} = 2 \left( -\mathbf{UR} + \mathbf{U}(\mathbf{SG}^\top \odot \mathbf{U}^\top \mathbf{V}) - 2\beta\mathbf{VC} + 2\beta\mathbf{VV}^\top \mathbf{V} + \lambda\mathbf{V} \right)$$

# Q2: Leveraging Social Data for Prediction

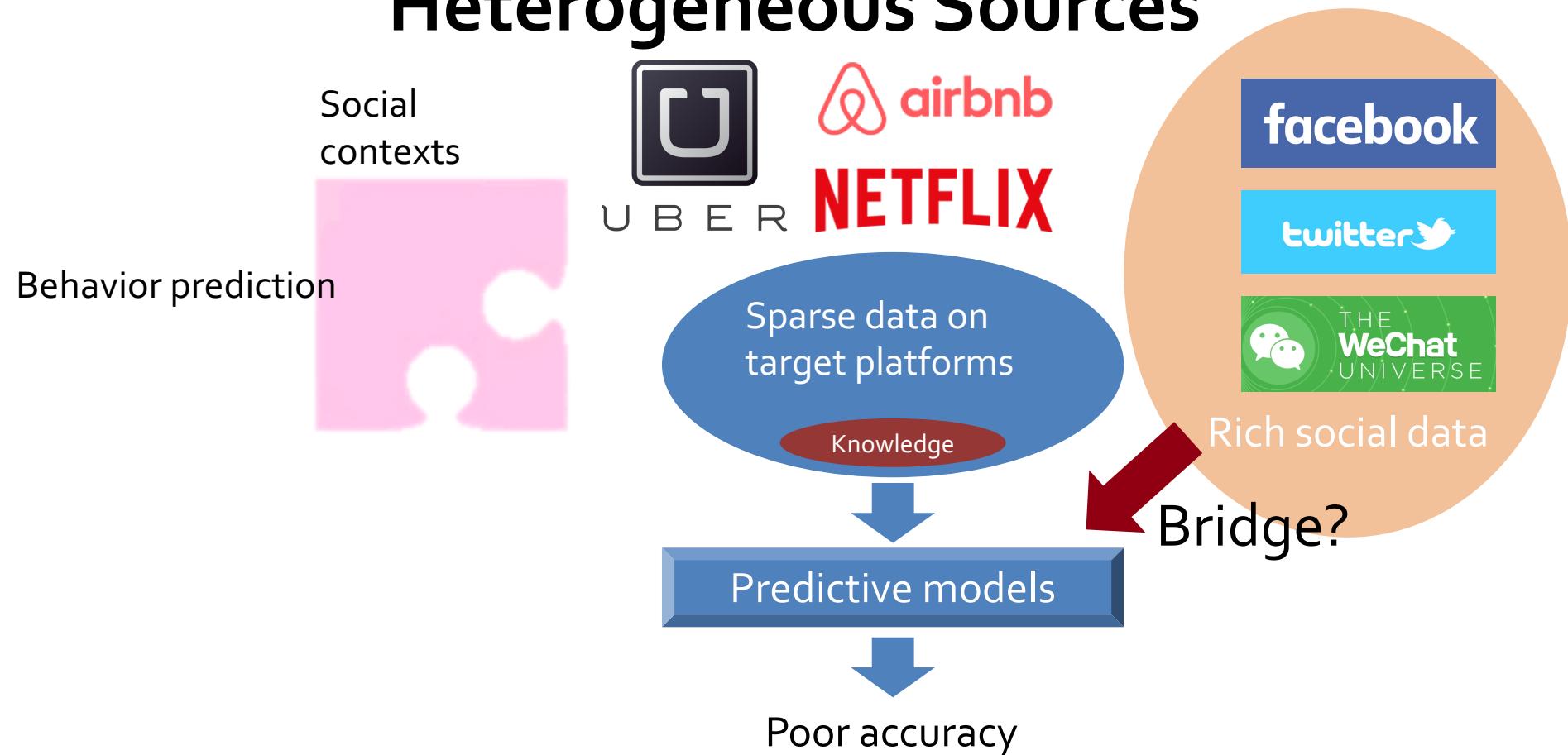


Data sparseness [Herlocker et al. CSCW'00; Sarwar et al. WWW'01; Burke UM&UAI'02; Ma et al. TOIS'11 TIST'11; Tang et al. Soc. Netw. Anal. Min.; Xue et al. VLDB'15; Han and Obradovic et al. SDM'16]

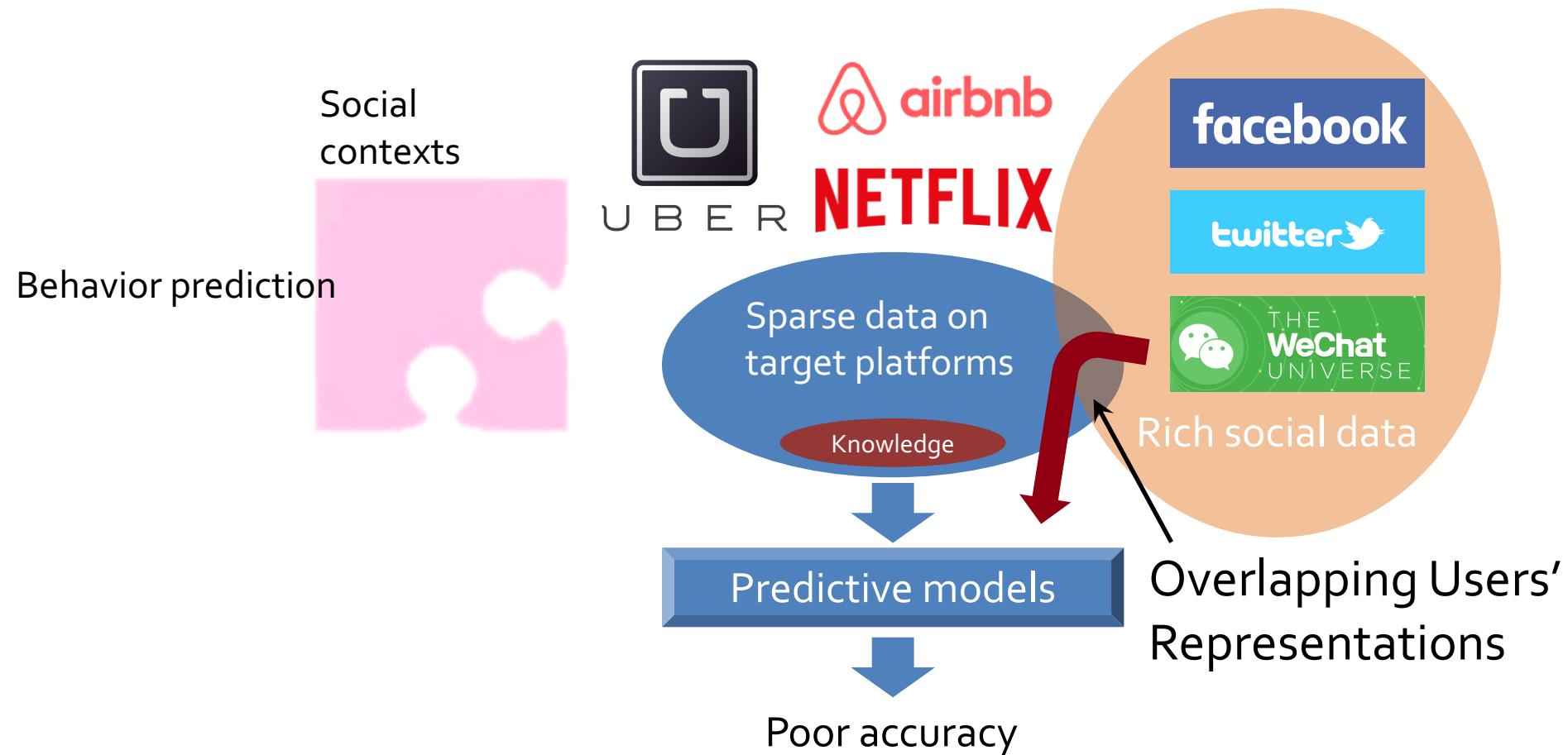
# Q2: Leveraging Social Data for Prediction



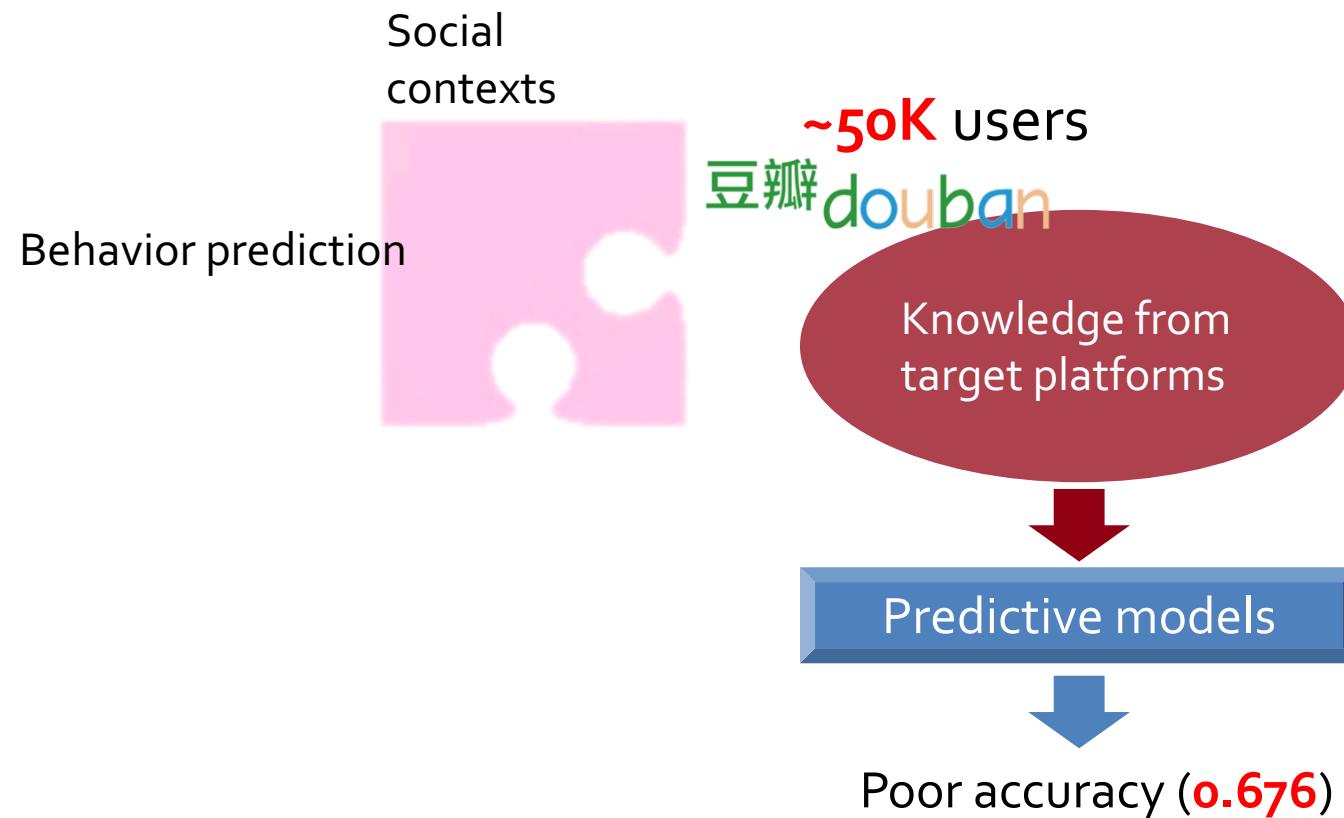
# M2: Knowledge Transfer from Heterogeneous Sources



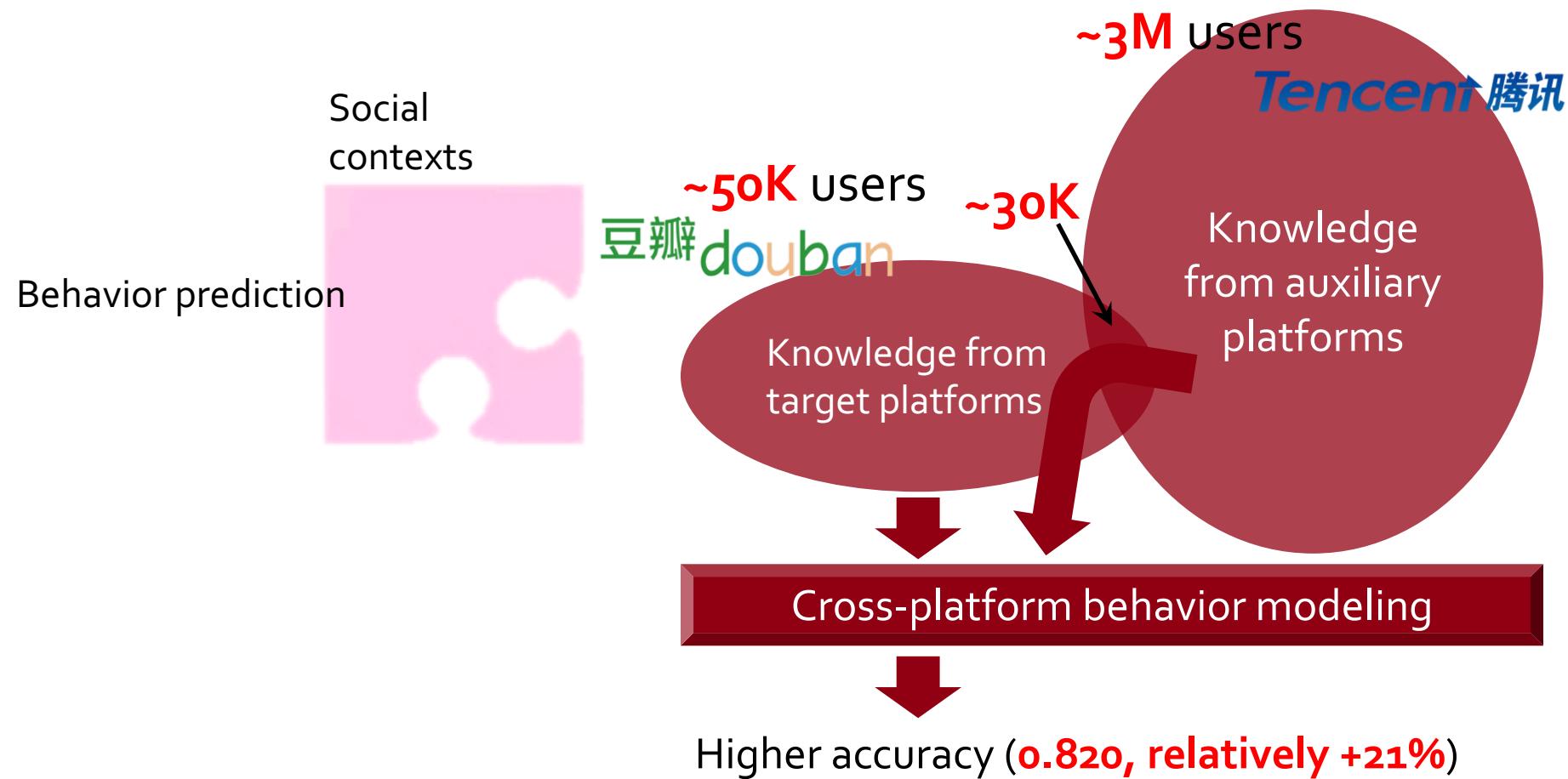
## M2: Knowledge Transfer from Heterogeneous Sources



## M2: Knowledge Transfer from Heterogeneous Sources

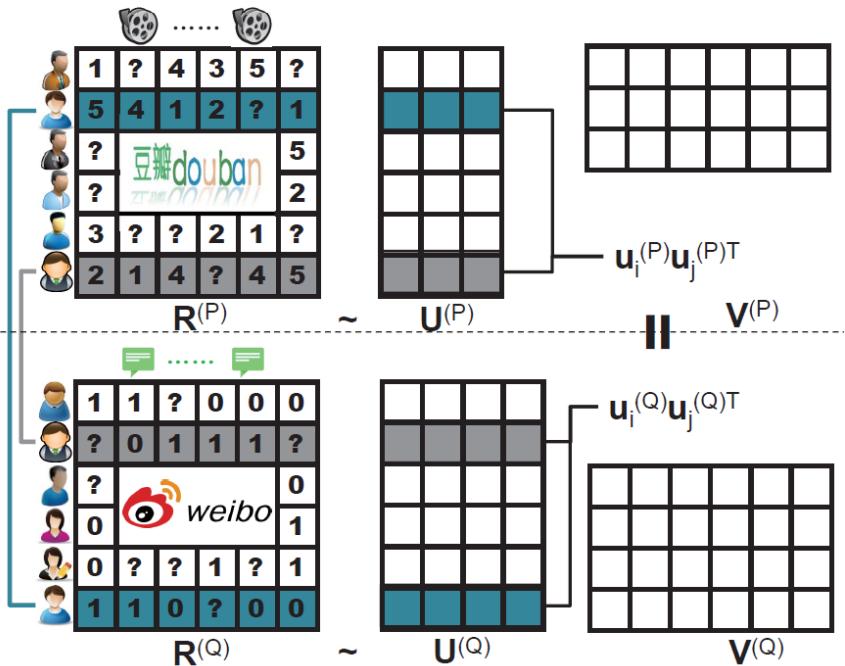


## M2: Knowledge Transfer from Heterogeneous Sources



S2: CIKM'12: cited by 72; TKDE'15: cited by 44; AAAI'16: cited by 13.

# S2: Cross-Platform Behavior Modeling



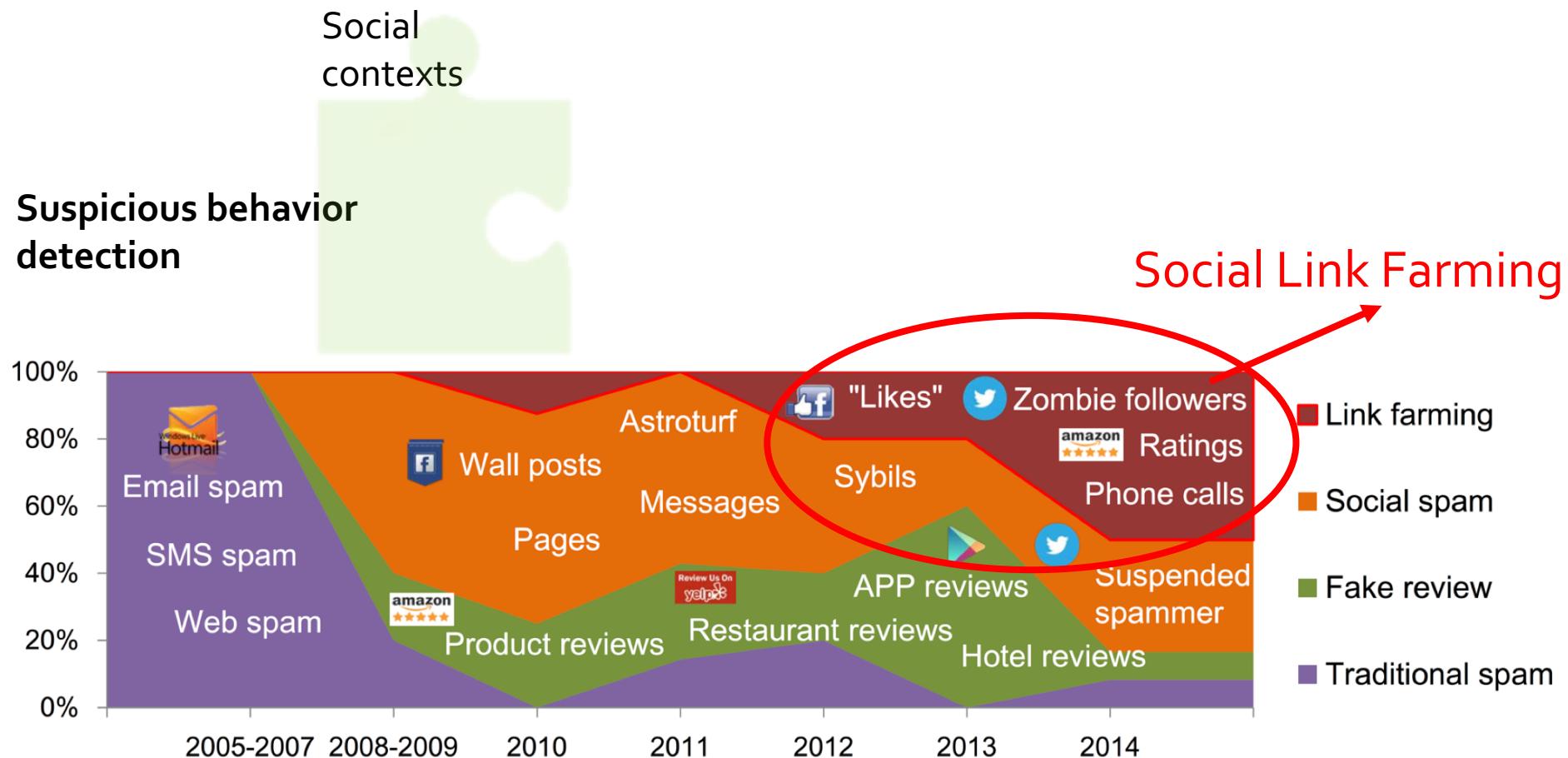
$$\begin{aligned}
 \mathcal{J} = & \sum_{i,j} W_{i,j}^{(P)} \left( R_{i,j}^{(P)} - \sum_r U_{i,r}^{(P)} V_{r,j}^{(P)} \right)^2 \\
 & + \lambda \sum_{i,j} W_{i,j}^{(Q)} \left( R_{i,j}^{(Q)} - \sum_r U_{i,r}^{(Q)} V_{r,j}^{(Q)} \right)^2 \\
 & + \mu \sum_{i_1,j_1,i_2,j_2} W_{i_1,j_1}^{(P,Q)} W_{i_2,j_2}^{(P,Q)} \left( A_{i_1,i_2}^{(P)} - A_{j_1,j_2}^{(Q)} \right)^2
 \end{aligned}$$

Overlapping user similarity  
(Pair-wise regularization)

Target platform      Auxiliary platform

Overlapping user similarity  
(Pair-wise regularization)

# Q3: Catching Social Link Farming



Meng Jiang, Peng Cui, and Christos Faloutsos. "Suspicious behavior detection: current trends and future directions." *IEEE Intelligent Systems*, 2016. (Survey paper)

# Q3: Catching Social Link Farming

Social contexts

Suspicious behavior detection

**CATCH SYNC**

**5,000 FOLLOWERS** 100 FREE  
\$69.99  
Delivery within 3-4 days  
Buy Now  
Save + 3%  
VISA

**2,000 FOLLOWERS** 200 FREE  
\$29.99  
Delivery within 2-3 days  
Buy Now  
Save + 2%  
VISA

**1,000 FOLLOWERS** 200 FREE  
\$15.99  
Delivery within 1-2 days  
Buy Now  
VISA

**10,000 FOLLOWERS** 300 FREE  
\$119.99  
Delivery within 4-5 days  
Buy Now  
Save + 14%  
VISA

**20,000 FOLLOWERS** 1000 FREE  
\$229.99  
Delivery within 5-8 days  
Buy Now  
Save + 34%  
VISA

**25,000 Facebook Likes**  
\$265  
Lifetime Replacement Warranty  
Dedicated 24/7 Customer Service  
100% Risk Free, Try Us Today  
Order starts within 24 - 48 hours  
Order completed within 22 days

**50,000 Facebook Likes**  
\$525  
Lifetime Replacement Warranty  
Dedicated 24/7 Customer Service  
100% Risk Free, Try Us Today  
Order starts within 24 - 48 hours  
Order completed within 35 days

**100,000 Facebook Likes**  
\$1,000  
Lifetime Replacement Warranty  
Dedicated 24/7 Customer Service  
100% Risk Free, Try Us Today  
Order starts within 24 - 48 hours  
Order completed within 35 days

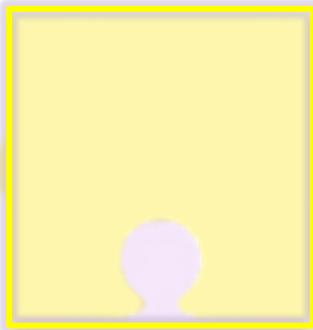
**200,000 Facebook Likes**  
\$1,750  
Lifetime Replacement Warranty  
Dedicated 24/7 Customer Service  
100% Risk Free, Try Us Today  
Order starts within 24 - 48 hours  
Order completed within 35 days

S3: CATCH SYNC

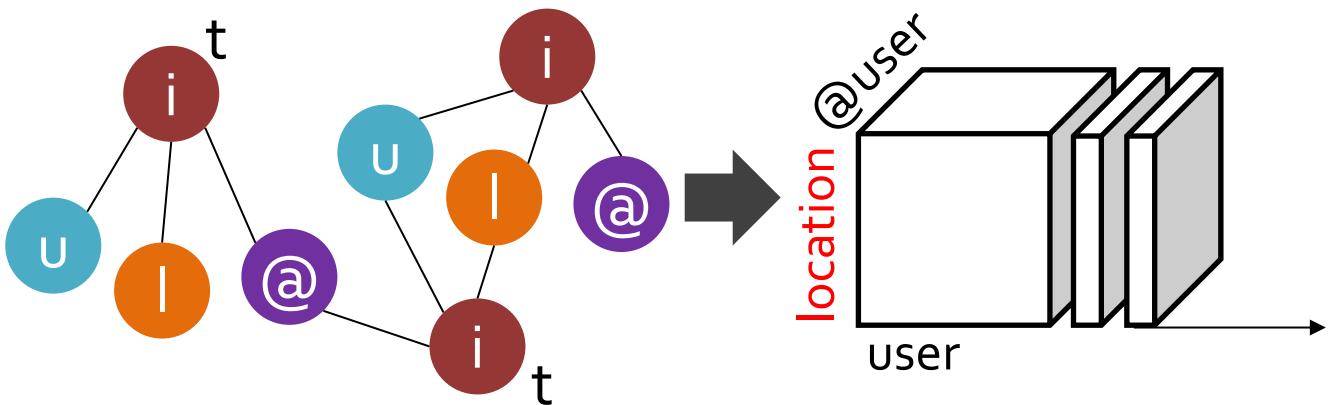
KDD'14 best paper finalist

## Q4: Knowledge from Spatiotemporal Information M4: Tensor Methods for Modeling Multiple Dimensions

Spatiotemporal contexts



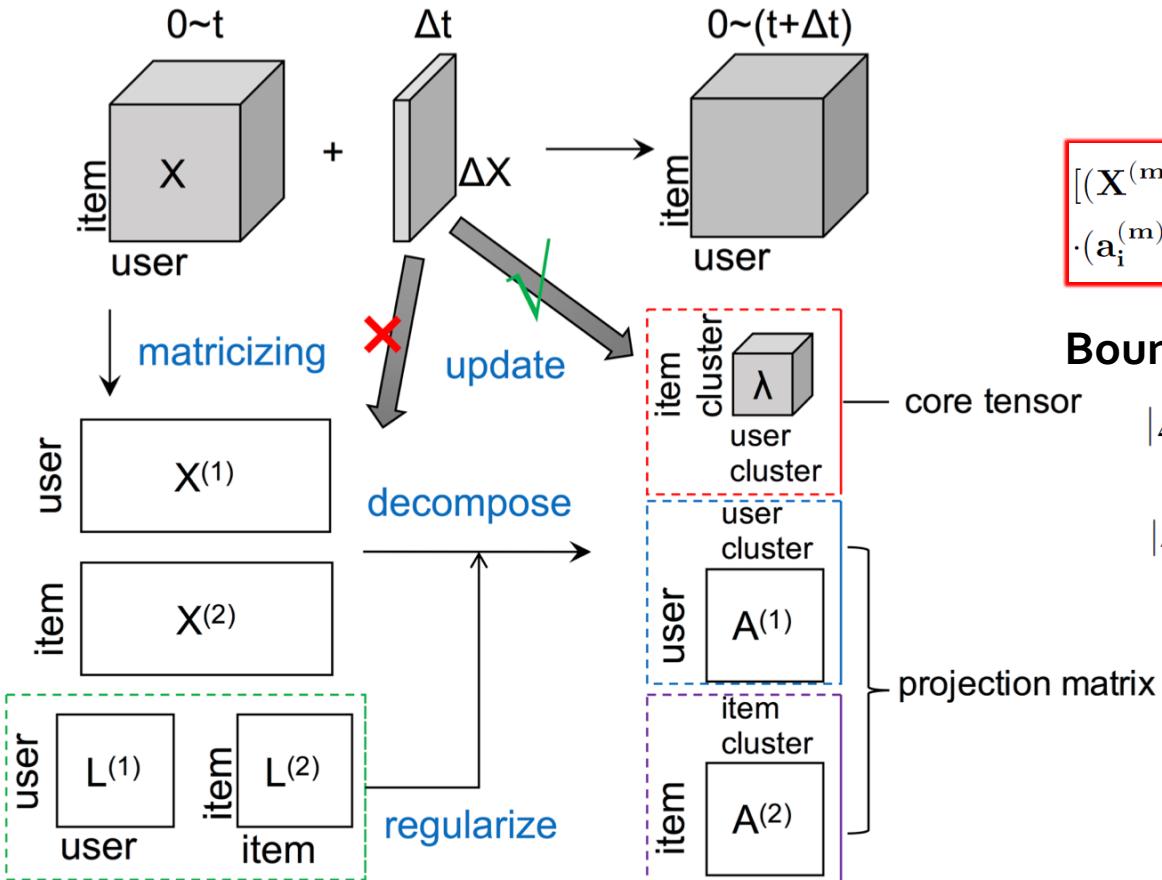
Q4-1: Who-@-whom prediction: High **complexity!**



M4-1: Incremental tensor decomposition:  
Approximation based on tensor perturbation theory. We proved bound guarantees of errors.

S4-1: RMSE reduced from **1.120** to **0.894**  
**(relatively -20.2%)**; running time reduced from **25 hours** to **51 minutes**. KDD'14: cited by **28**.

# S4-1: Flexible Evolutionary Multifaceted Analysis



**Tensor perturbation theory:**

$$[(\mathbf{X}^{(m)} + \Delta \mathbf{X}^{(m)}) (\mathbf{X}^{(m)} + \Delta \mathbf{X}^{(m)})^\top + \mu^{(m)} \mathbf{L}^{(m)}] \cdot (\mathbf{a}_i^{(m)} + \Delta \mathbf{a}_i^{(m)}) = (\lambda_i^{(m)} + \Delta \lambda_i^{(m)}) (\mathbf{a}_i^{(m)} + \Delta \mathbf{a}_i^{(m)})$$

**Bounds (guarantee for approximation):**

$$|\Delta \lambda_i^{(m)}| \leq 2(\lambda_{\mathbf{X}^{(m)\top} \mathbf{X}^{(m)}}^{\max})^{\frac{1}{2}} \|\Delta \mathbf{X}^{(m)}\|_2$$

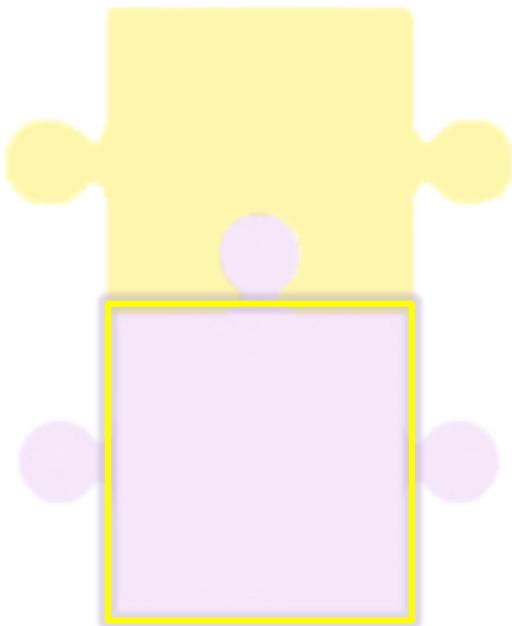
$$|\Delta \mathbf{a}_i^{(m)}| \leq 2 \|\Delta \mathbf{X}^{(m)}\|_2 \sum_{j \neq i} \frac{(\lambda_{\mathbf{X}^{(m)\top} \mathbf{X}^{(m)}}^{\max})^{\frac{1}{2}}}{|\lambda_i^{(m)} - \lambda_j^{(m)}|}$$

projection matrix

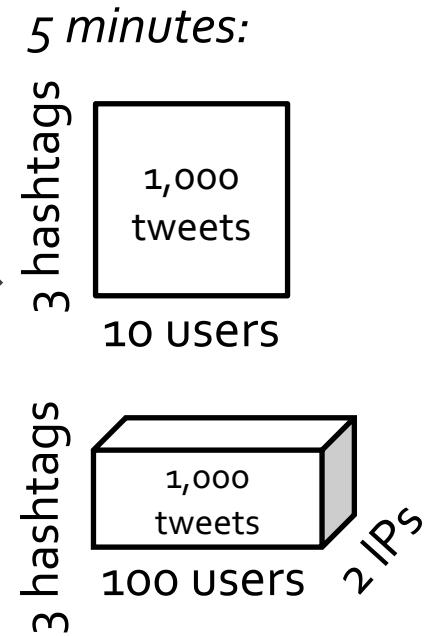
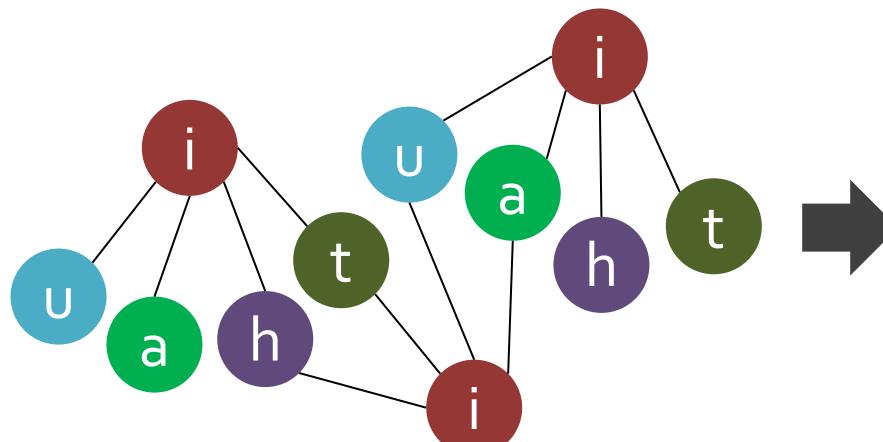
# Q4: Knowledge from Spatiotemporal Information

# M4: Tensor Methods for Modeling Multiple Dimensions

Spatiotemporal  
contexts



Q4-2: Spam detection: Evaluating suspiciousness?

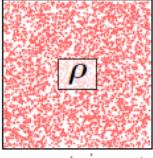
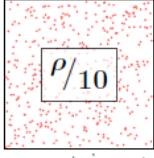
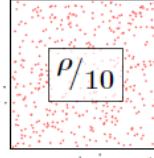
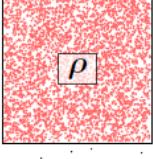
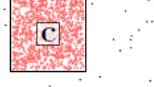


M4-2: Proposed a **principled** suspiciousness metric.

S4-2: Detected  $\sim 6M$  hijacking tweets of  $3$  hashtags  
by  $\sim 600$  users from  $\sim 300$  IP addresses in  $\sim 40$  days.

ICDM'15: cited by **42**; TKDE'16: cited by **8**.

# S4-2: Evaluating Suspiciousness across Dimensions

Density Axiom	>	Contrast Axiom	>	
	>		>	
Size Axiom	>	Concentration Axiom	>	
	>		>	

# Q5: Knowledge from Behavioral Content

## From Words, Topics, to Networks

“Modeling Complex  
Behaviors in Social  
Media”, July 2015. 



清华大学

Tsinghua University

# Q5: Knowledge from Behavioral Content

## From Words, Topics, to Networks

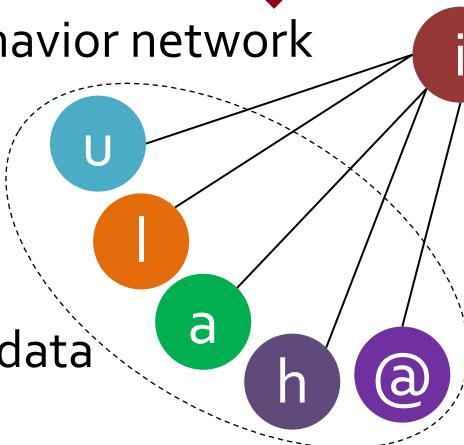
"Modeling Complex Behaviors in Social Media", July 2015. 



清华大学

Tsinghua University

Behavior network



Information network  
(entities, attributes...)

Integrating

Structuring

Rich unstructured text data



tweets, news...



product/restaurant review...



Publications:  
PubMed, DBLP...

Structured data



# Q5': Entity and Attribute Discovery

Given full **text** of all the Data Science publications

Q'5-1. Who has studied the biggest number of **datasets** of **large scale**?

Q'5-2. Who study **truly big data**, and who always claim their work is on **big data** but their datasets are not **big** at all?

Q'5-3. For a **dataset** and a **problem**, who are the **experts**? How can we organize a **team** to solve the problem?

# Q5': Entity and Attribute Discovery

Given full **text** of all the Data Science publications

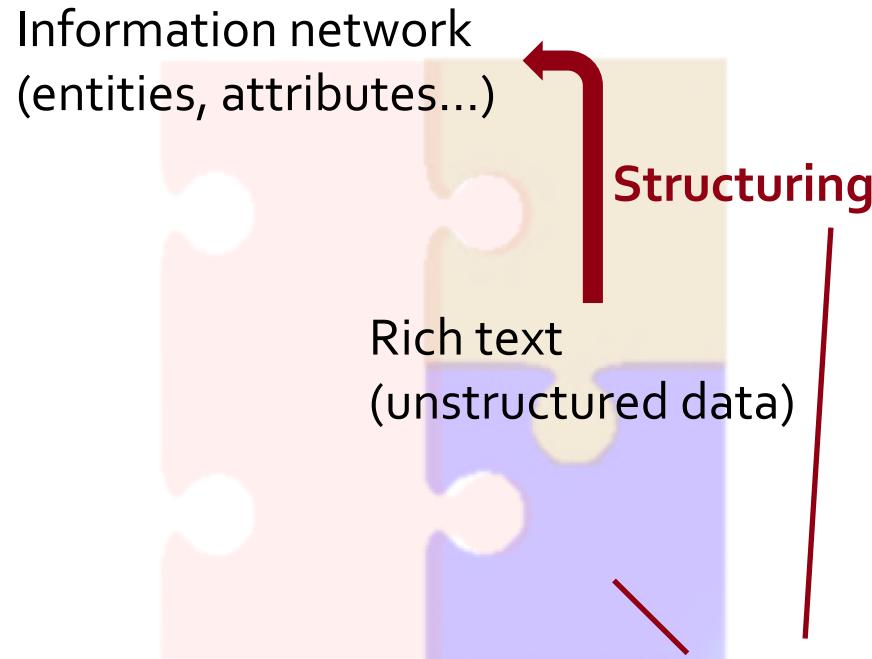
Q'5-1. Who has studied the biggest number of **datasets** of **large scale**?

Q'5-2. Who study **truly big data**, and who always claim their work is on **big data** but their datasets are not **big** at all?

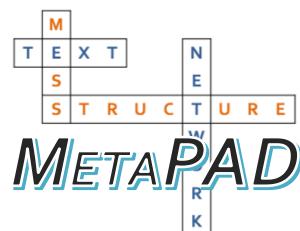
Q'5-3. For a **dataset** and a **problem**, who are the **experts**? How can we organize a **team** to solve the problem?

Sorry, I don't have answers now... But...

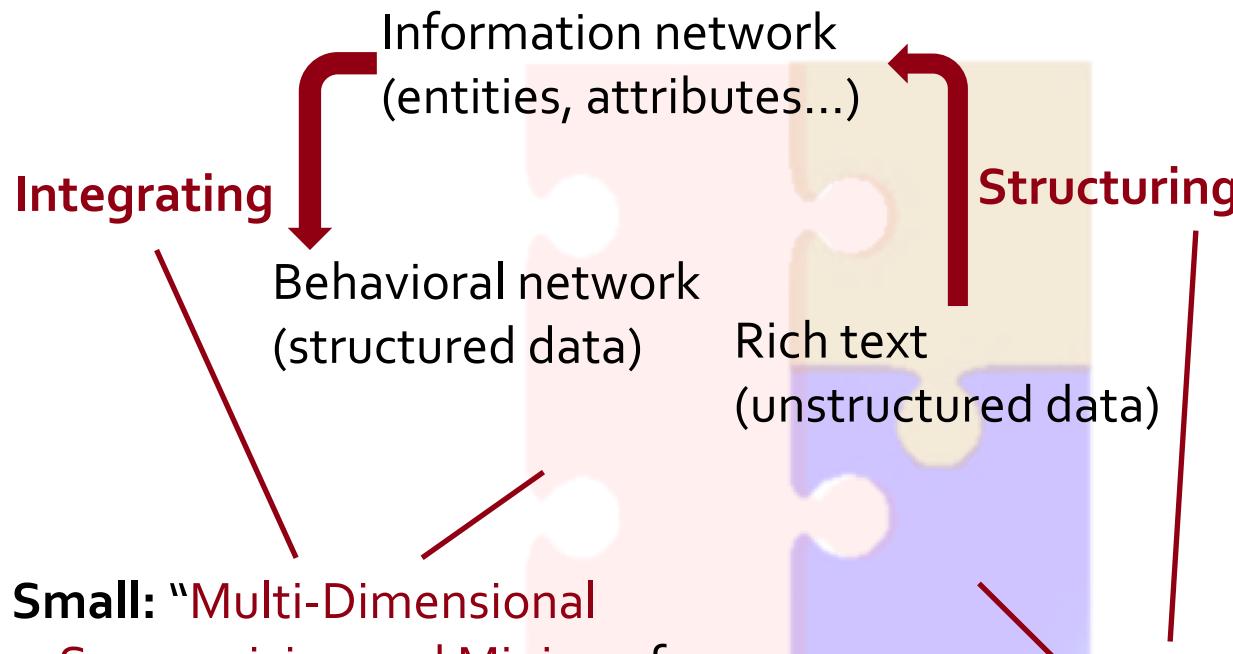
# S5: Multiple Proposal Writing and Papers



1. **NSF III: Medium:** Collaborative: “StructNet: **Constructing** and **Mining** Structure-Rich **Information Networks** for Scientific Research”. (Funded 2017)
2. **KDD’17:** “Meta Pattern-Driven **Attribute Discovery** from Massive Text Corpora”. (Accepted)



# S5: Multiple Proposal Writing and Papers



3. **NSF III: Small:** “**Multi-Dimensional Structuring, Summarizing and Mining of Social Media Data**”, NSF IIS 16-18481. Jiawei Han, PI. (Submitted Nov 15, funded Aug 16)  
4. **KDD’16:** “**CatchTartan: Representing and Summarizing Dynamic Multicontextual Behaviors**”. **Oral** (Acc. = 8.9%).

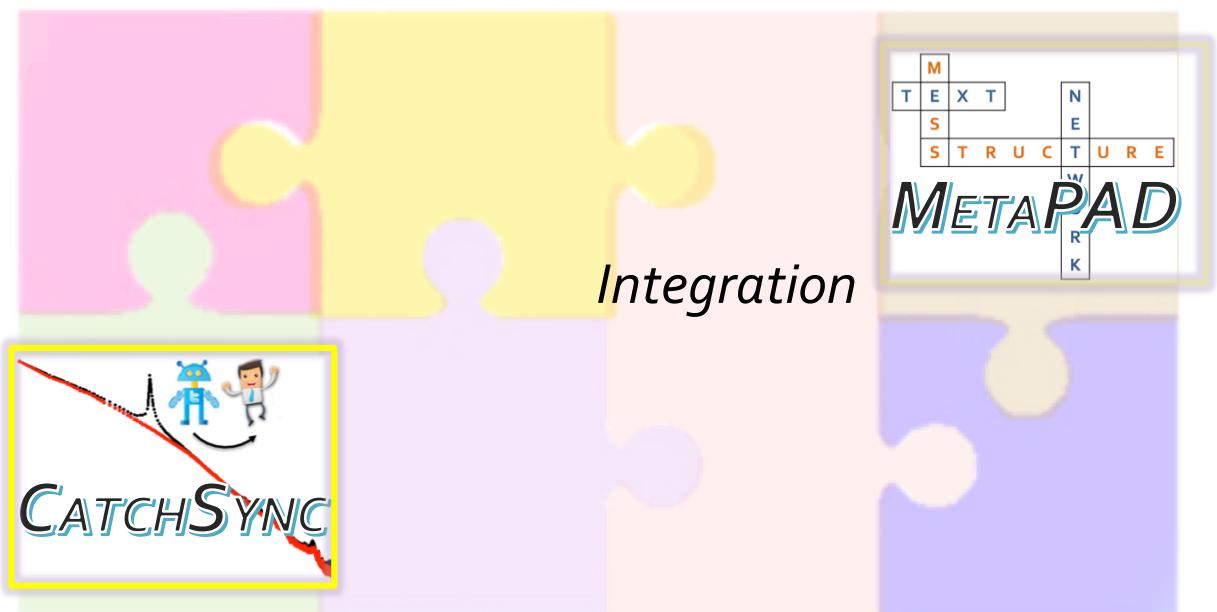
1. **NSF III: Medium:** Collaborative: “**StructNet: Constructing and Mining Structure-Rich Information Networks for Scientific Research**”. (Funded 2017)  
2. **KDD’17:** “**Meta Pattern-Driven Attribute Discovery** from Massive Text Corpora”. (Accepted)

# Outline

*Intelligence:*  
Behavior prediction  
and recommendation

*Trustworthiness:*  
Suspicious behavior  
detection

Social contexts      Spatiotemporal contexts      Behavioral content



# CatchSync: Catching Synchronized Behavior in Large Directed Graphs

Joint work with Peng Cui, Shiqiang Yang (Tsinghua),

Alex Beutel, and Christos Faloutsos (CMU)

**ACM SIGKDD 2014 Best Paper Finalist**

(among **151** accepted research papers, **1,036** submissions)



# Q3: Catching Zombie Followers



# Q3: Catching Zombie Followers



engineers



product managers

Knowledge  
from  
manual  
inspection:

#followees,  
#followers, #tweets,  
#hashtags, #urls...

Learning models (classifiers)

Fake account detection [Egele and Stringhini et al. NDSS'13; Yang and Wilson et al. TKDD'14; Viswanath and Bashir et al. USENIX Security Symposium'14]



Poor accuracy  
**(serious complaints** from users)

# Is this account a zombie follower???

Aisling Walsh  
@xAsherzka

Joined April 2009

[Tweet to Aisling Walsh](#)

Who to follow · Refresh · View all

- John Legere @JohnLe...  
[Follow](#)  
Promoted
- Dong Zhou @dongz9  
Followed by Peng Wang 王鵬 and others  
[Follow](#)
- Justin Zeus @askzy9  
Followed by Ruizhe, Li and others  
[Follow](#)

Find friends

Trends · Change

- #ThatsContinental  
Allowing curiosity to chart your course.  
[Promoted by LincolnMotorCompany](#)
- #2017in3words  
26.1K Tweets
- #nationalbaconday  
5,915 Tweets
- #NewYearsEveEve  
2,501 Tweets

FOLLOWING 20 FOLLOWERS 3 0 tweet

Rachel Maddow MSN... @maddow  
I see political people... (Retweets do not imply endorsement.)

Trent Reznor @trent\_reznor  
Nine Inch Nails, How To Destroy Angels and other things.

Guardian Tech @guardiantech

Jason Sweeney @sween  
limited edition, macaroni and glitter on construction paper.

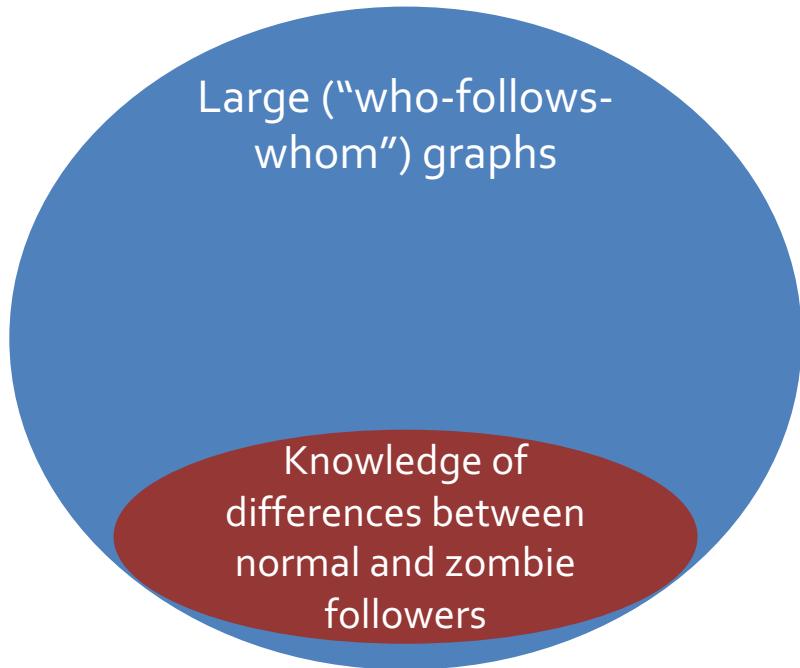
richard bacon @richardpbacon

woot.com @woot  
Check out who we're following for other Woot accounts, and follow us on Facebook for extra excitement: [facebook.com/woot](http://facebook.com/woot)

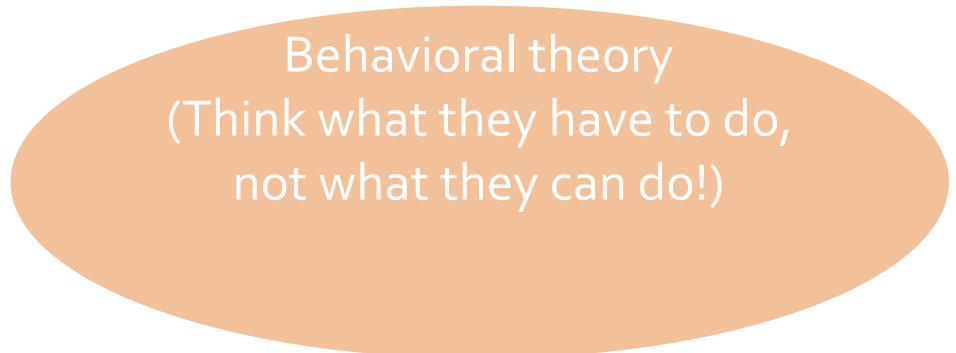
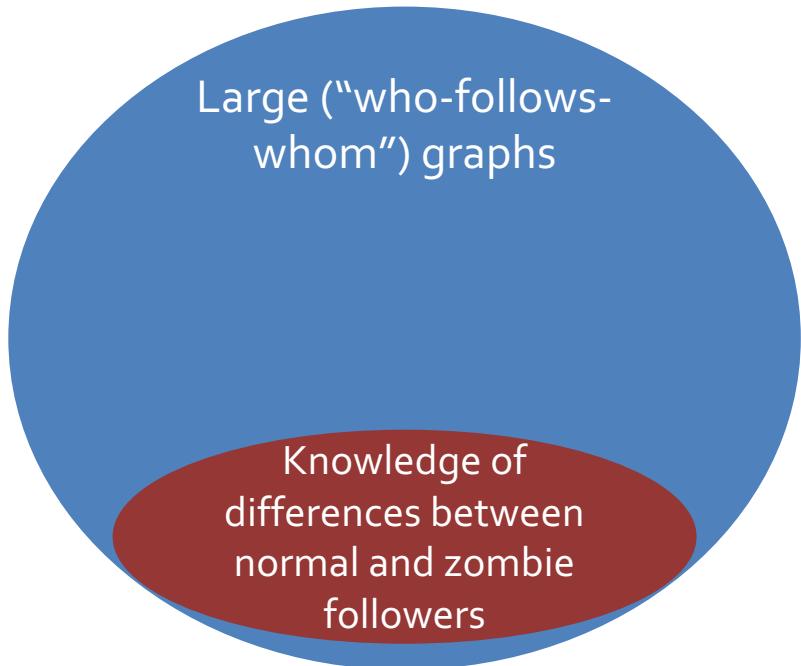
Hoppy New Year @markhoppus  
person

CBOE @CBOE

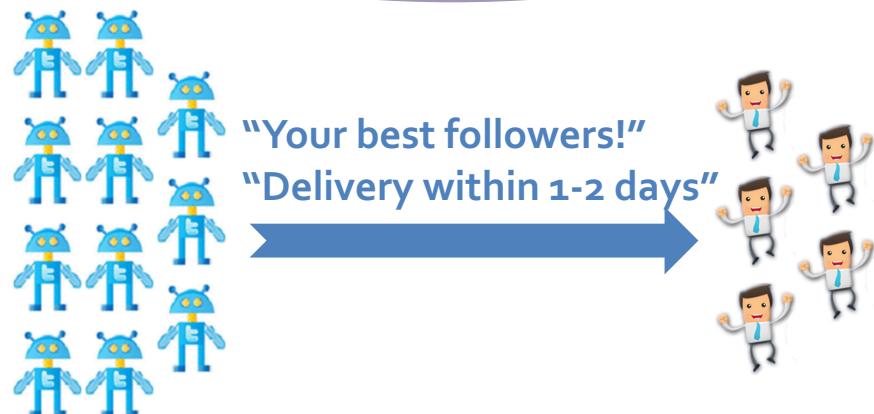
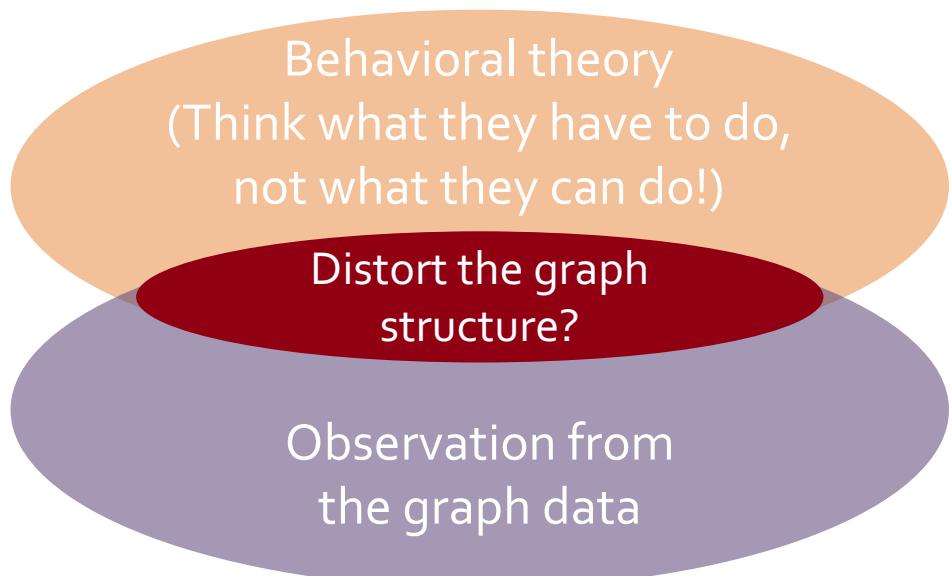
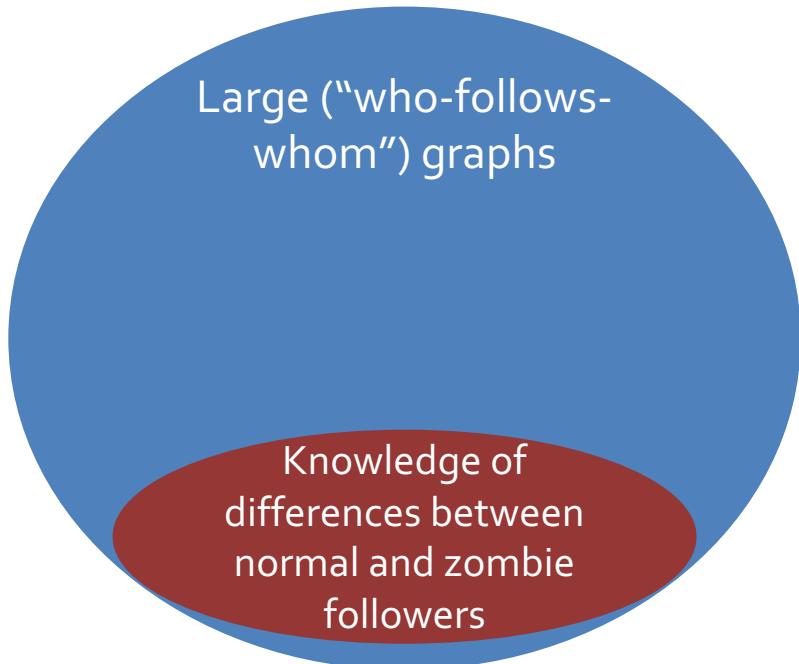
# Our Methodology (M1)



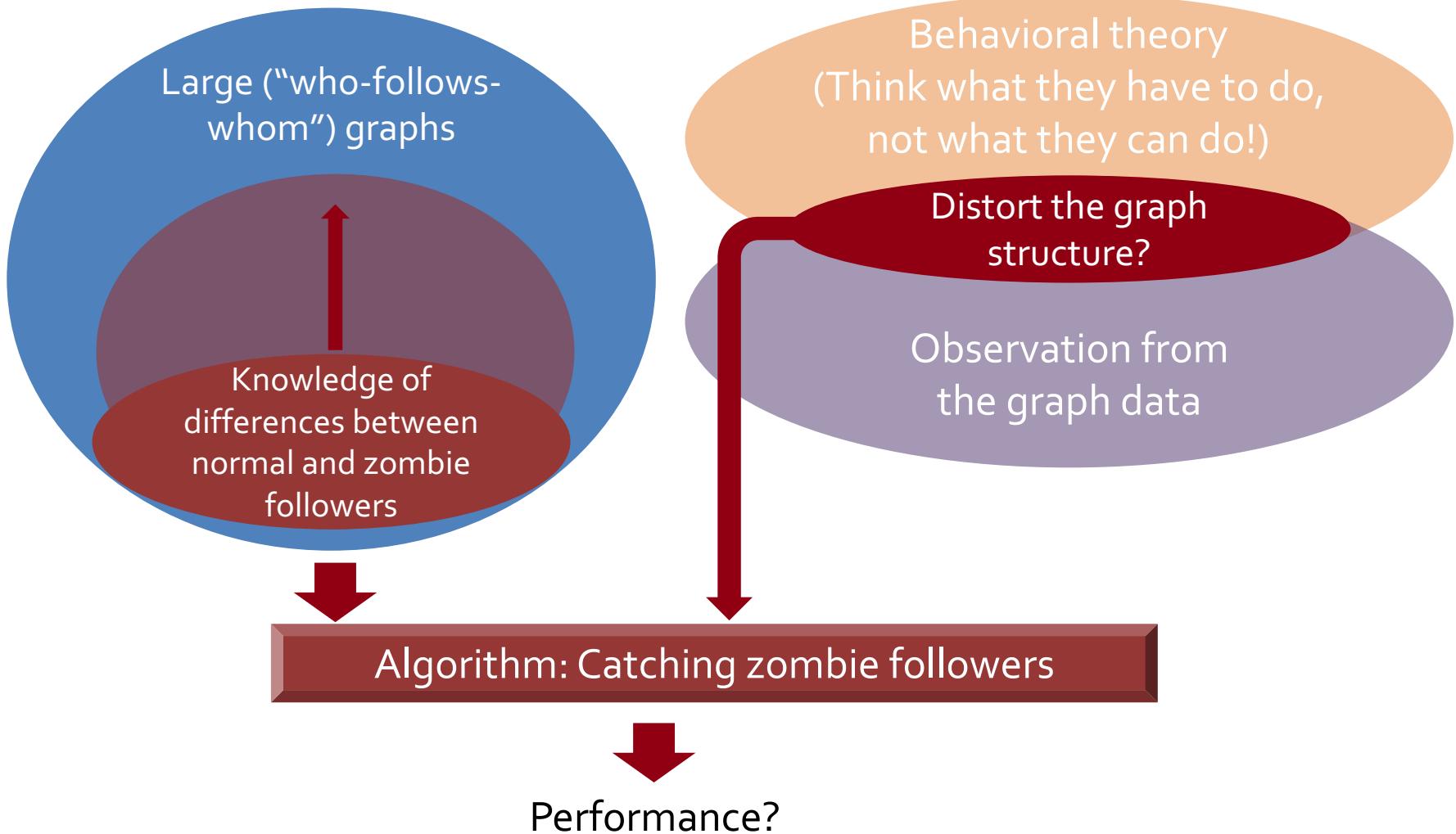
# Our Methodology (M1)



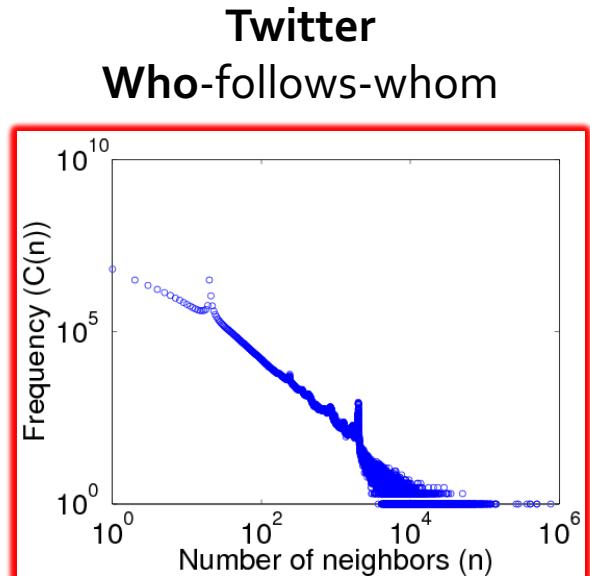
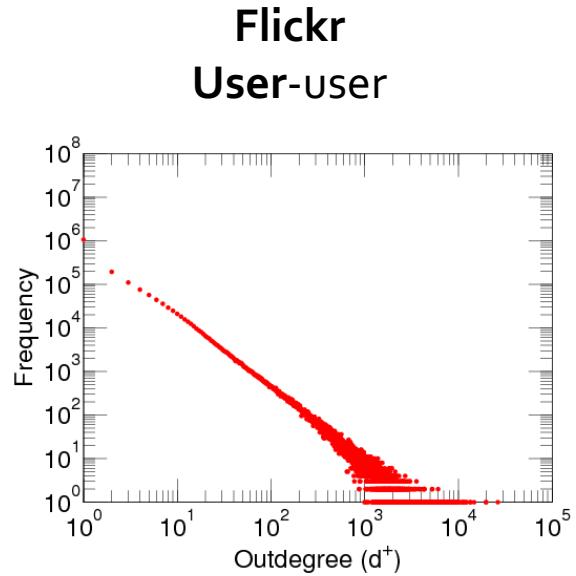
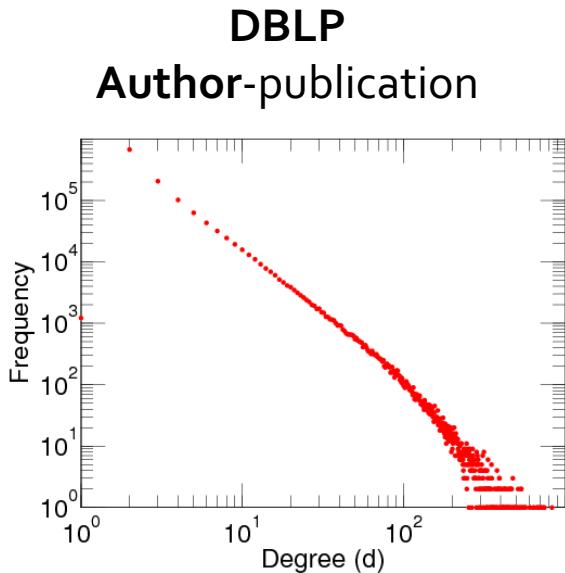
# Our Methodology (M1)



# Our Methodology (M1)



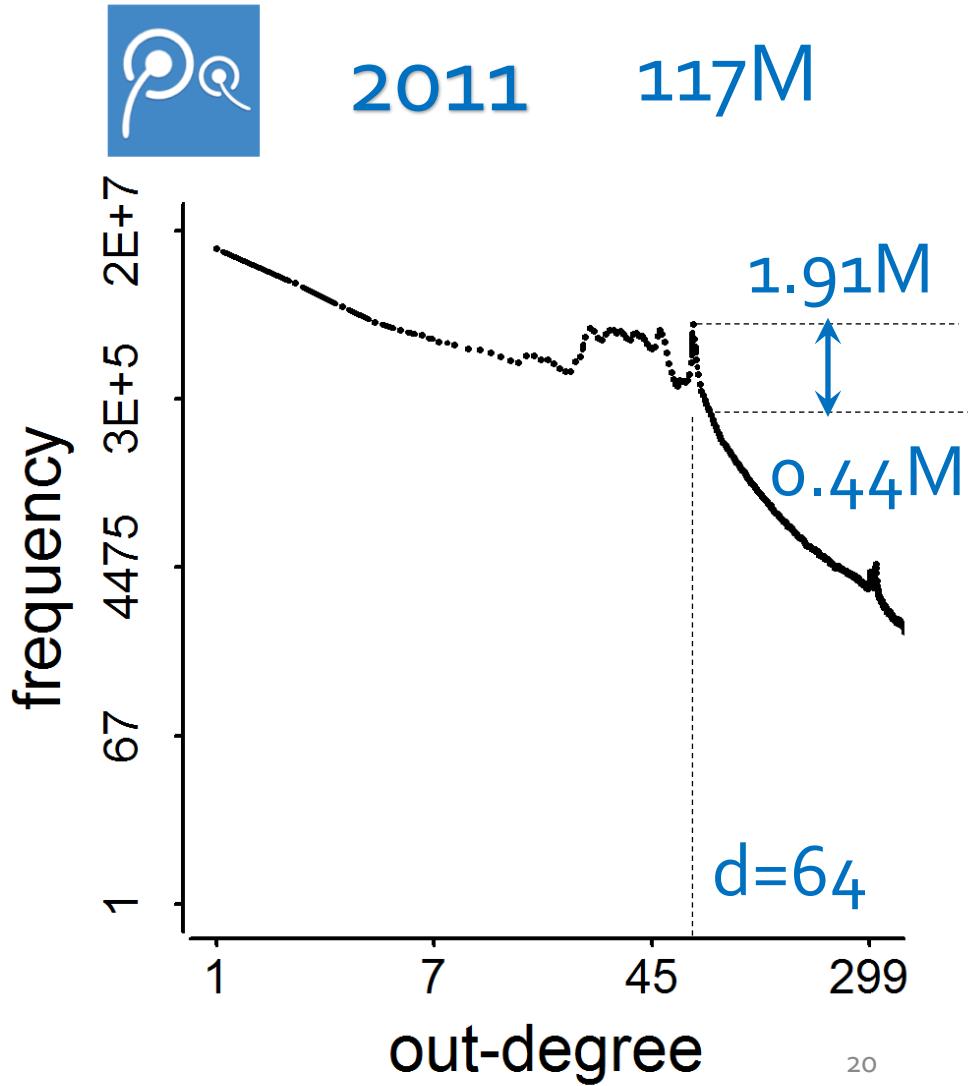
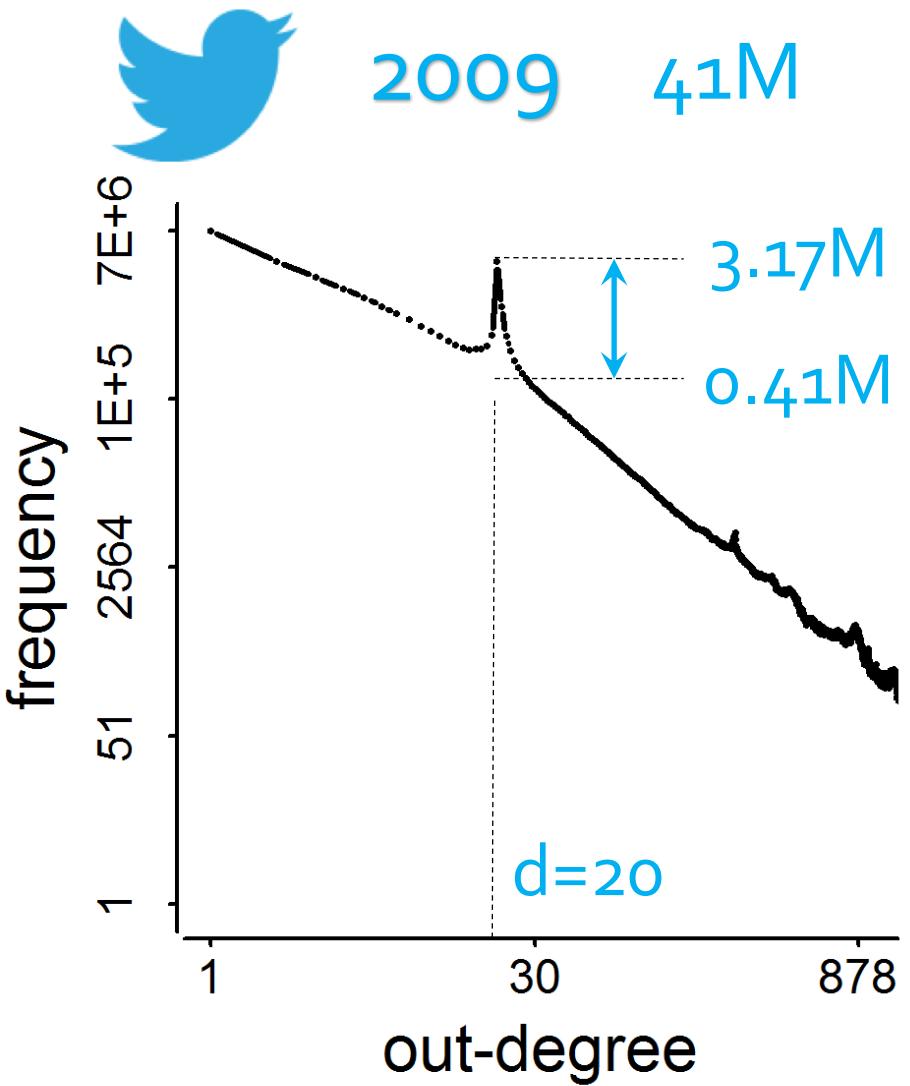
# Out-Degree Distributions: Power Law Expected



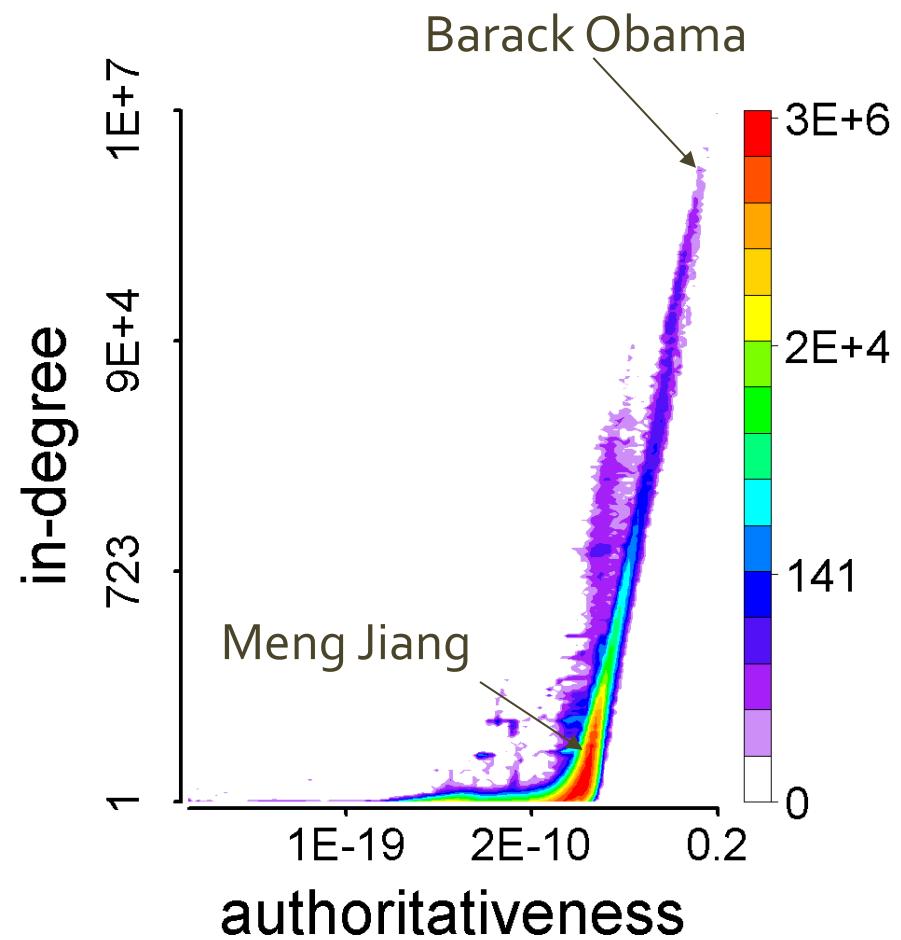
[[konect.uni-koblenz.de/networks/](http://konect.uni-koblenz.de/networks/)]

Power-law distributions in networks [Faloutsos et al.  
SIGCOMM'99; Chung et al. PNAS'02]

# Spikes!

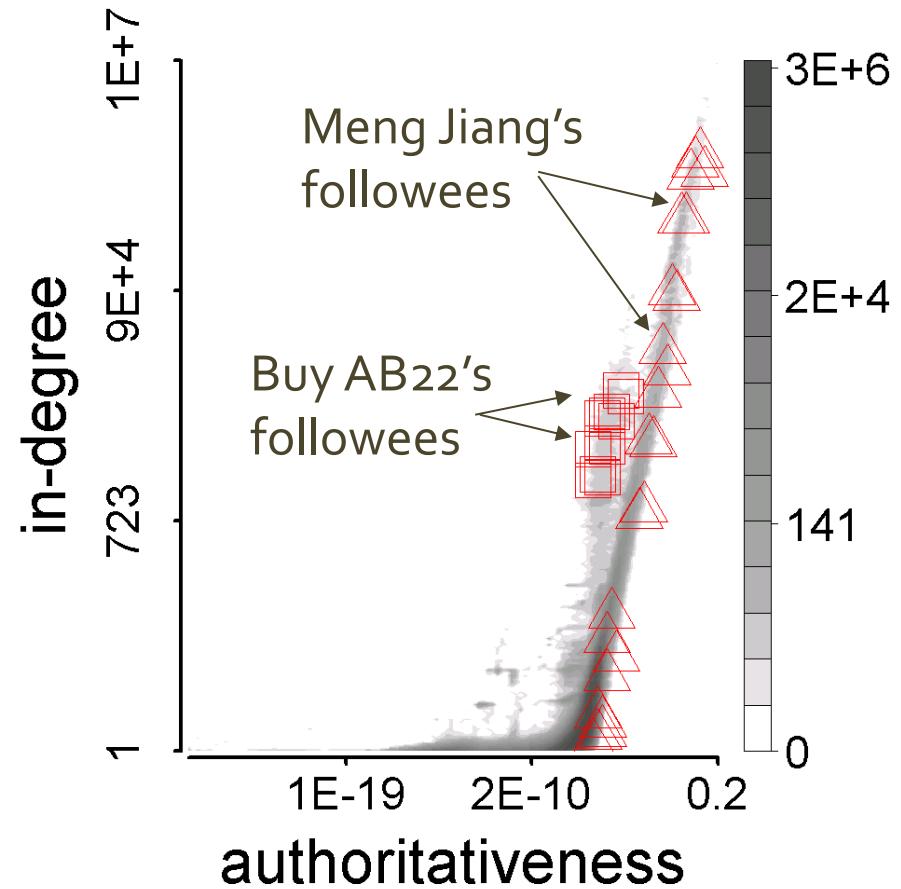


# How We/They Connect to Our/Their Followees

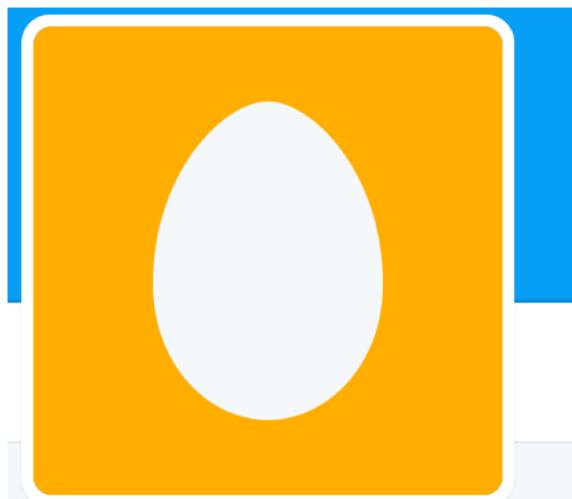


The HITS algorithm. Kleinberg. "Authoritative sources in a hyperlinked environment." JACM'99.

# How We/They Connect to Our/Their Followees



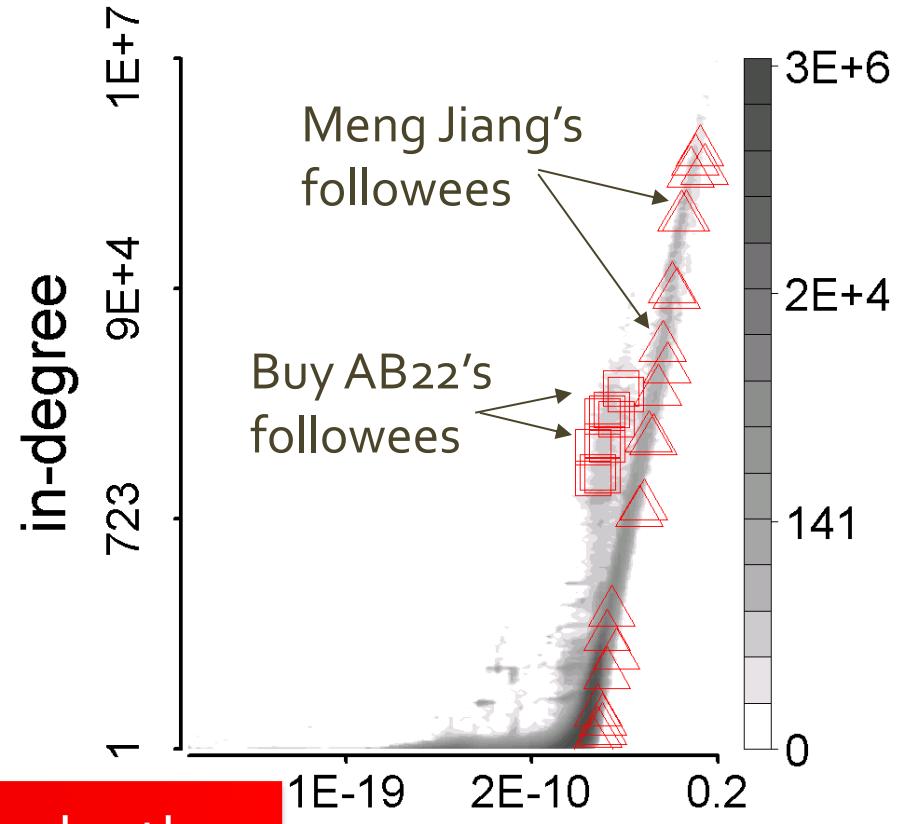
# How We/They Connect to Our/Their Followees



**Buy AB22 Propertwee**

@Buy\_AB22

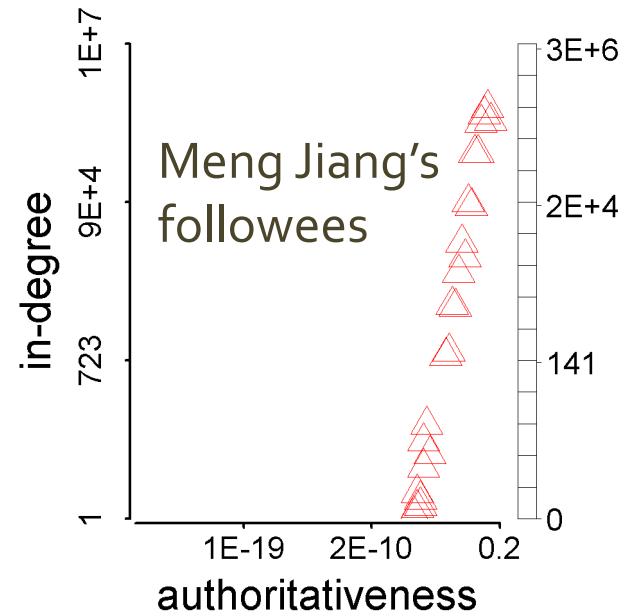
Joined May 2009



Synchronized: too similar with each other  
Abnormal: too different from the majority

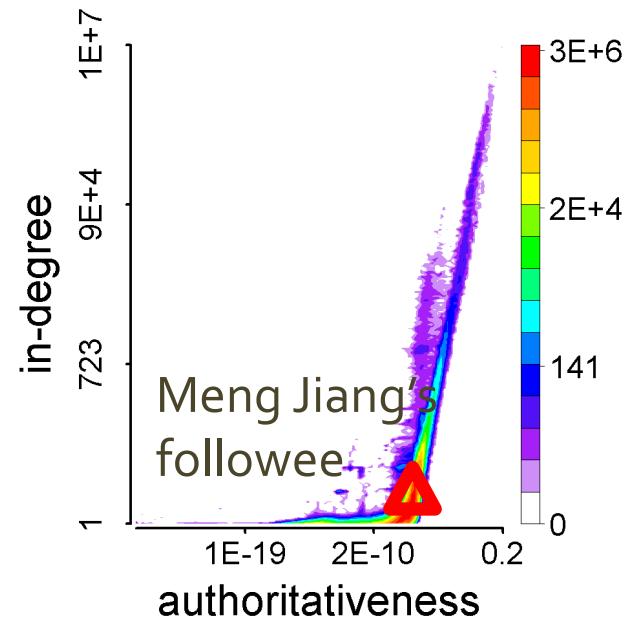
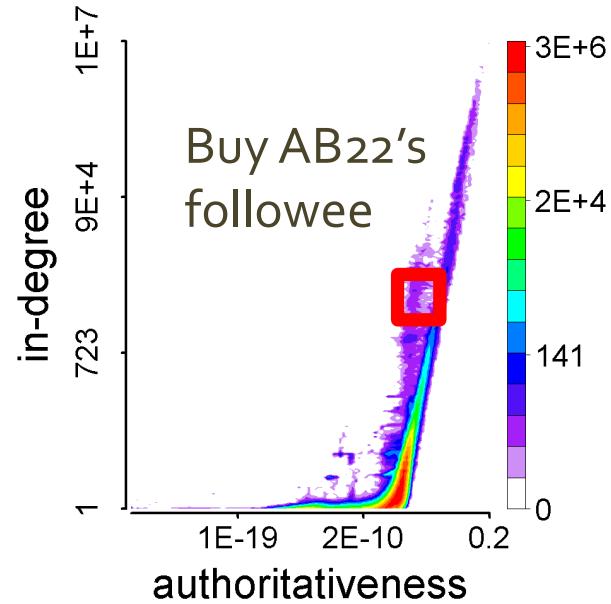
# Definition: Synchronicity

$$sync(u) = \frac{\sum_{(v,v') \in \mathcal{F}(u) \times \mathcal{F}(u)} \mathbf{p}(v) \cdot \mathbf{p}(v')}{d(u) \times d(u)}$$



# Definition: Normality

$$norm(u) = \frac{\sum_{(v,v') \in \mathcal{F}(u) \times \mathcal{U}} \mathbf{p}(v) \cdot \mathbf{p}(v')}{d(u) \times N}$$



# When is the Synchronicity Too High?

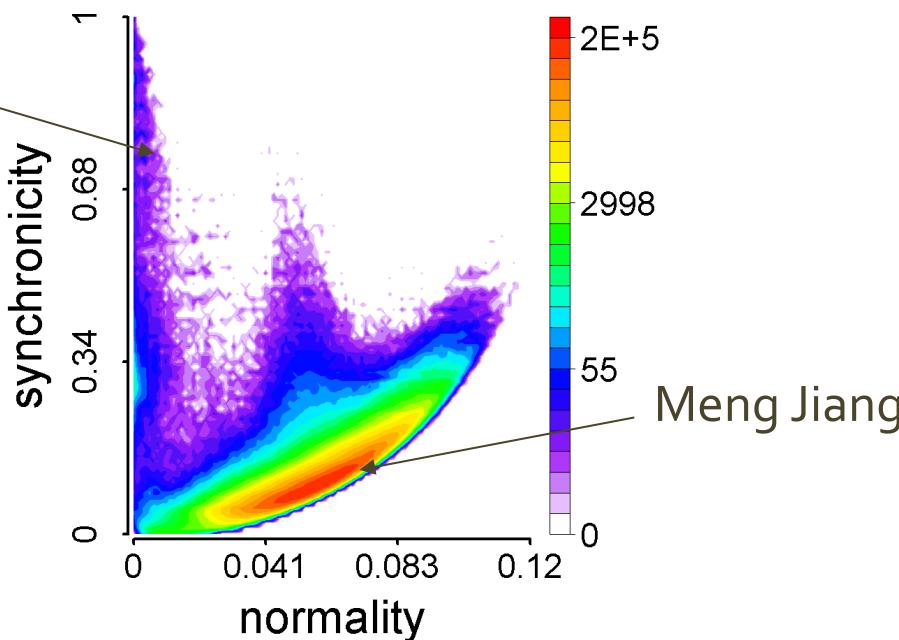
**Problem:** Given a normality value ( $n$ ) of a follower, find the minimal synchronicity value ( $s_{\min}$ ).

**Theorem:**

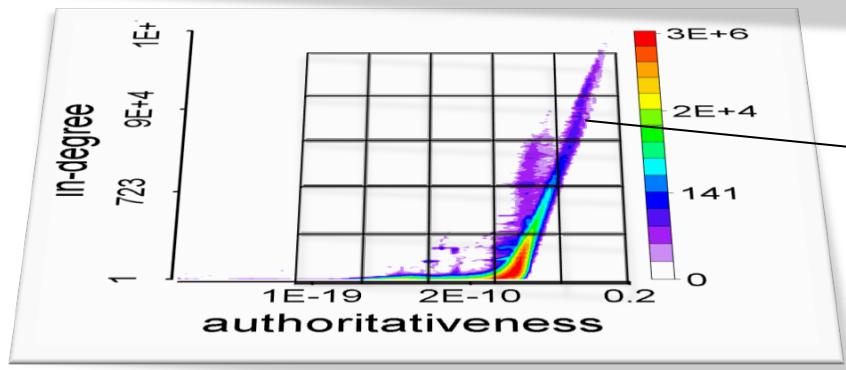
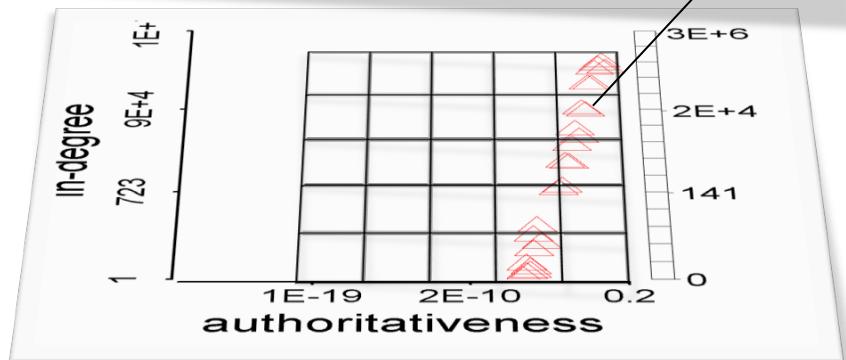
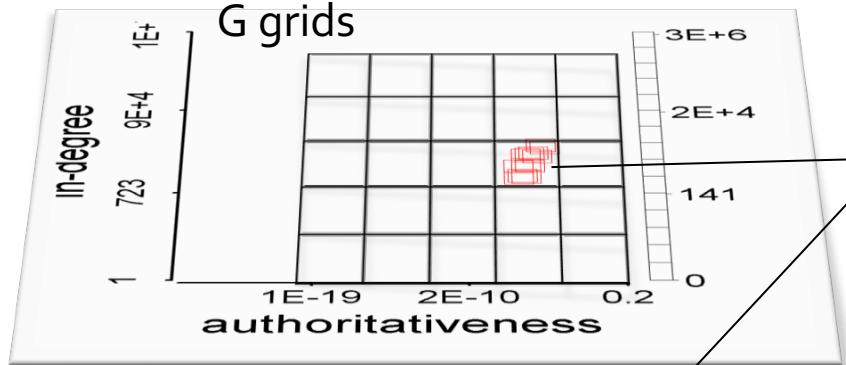
$$s_{\min} = \frac{-G n^2 + 2 n - s_b}{1 - G s_b} \quad (\text{parabolic lower limit})$$

**Our CatchSync:**

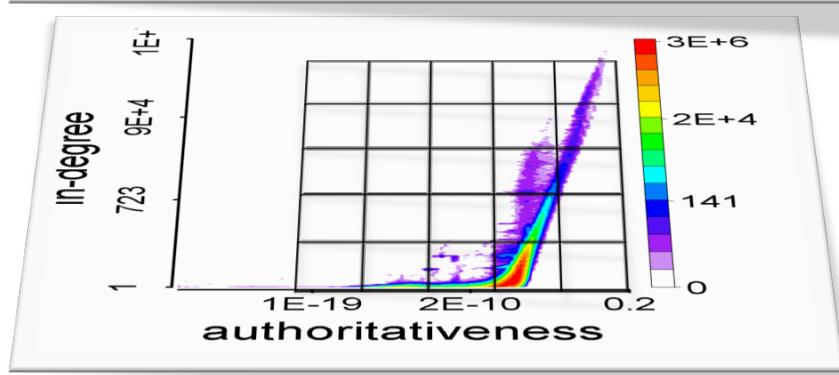
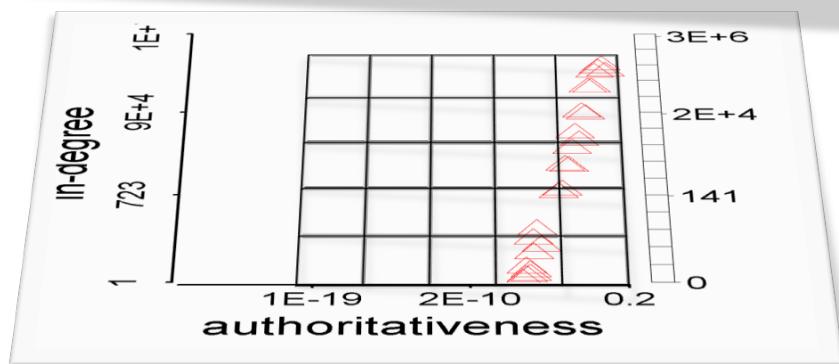
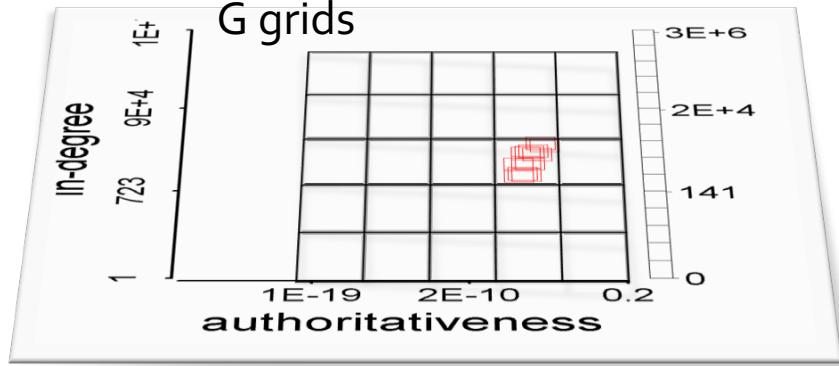
Buy AB22 &  
Aisling Walsh



# Proof



# Proof



**Lagrange multiplier:**

$$\text{minimize } s(f_g) = \sum f_g^2$$

$$\text{subject to } \sum f_g = 1, \sum f_g b_g = n$$

**Lagrange function:**

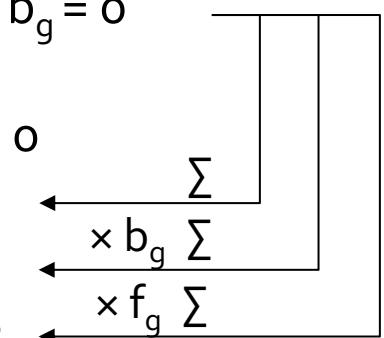
$$F(f_g, \lambda, \mu) = (\sum f_g^2) + \lambda (\sum f_g - 1) + \mu (\sum f_g b_g - n)$$

**Gradients:**

$$\left\{ \begin{array}{l} \nabla_{f_g} F = 2 f_g + \lambda + \mu b_g = 0 \\ \nabla_{\lambda} F = \sum f_g - 1 = 0 \\ \nabla_{\mu} F = \sum f_g b_g - n = 0 \end{array} \right.$$

$$\left\{ \begin{array}{l} 2 + \lambda G + \mu = 0 \\ 2 n + \lambda + \mu s_b = 0 \end{array} \right.$$

$$\left\{ \begin{array}{l} 2 s_{\min} + \lambda + \mu n = 0 \\ \sum b_g \sum \\ \sum f_g \sum \end{array} \right.$$

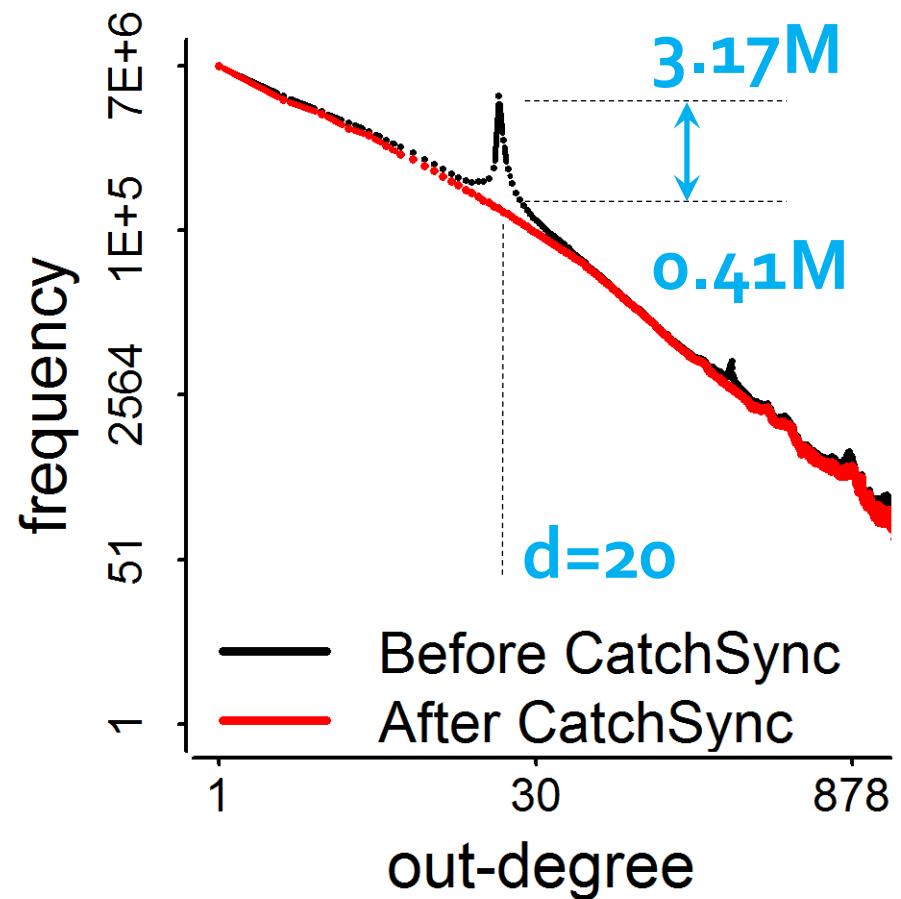
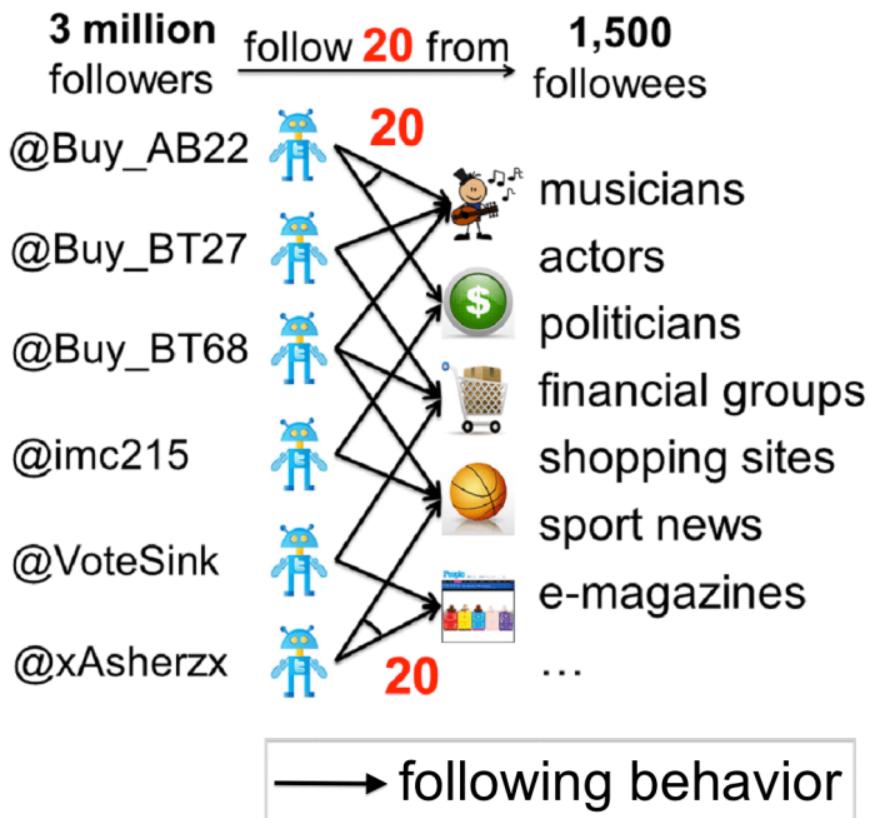


$$\text{where } s_b = \sum b_g^2.$$

Therefore,

$$s_{\min} = \frac{-G n^2 + 2 n - s_b}{1 - G s_b}$$

# The Distribution was Recovered!



# Impact

- Cited by **76**; WWW'14: cited by **21**. TKDD'16: cited by **21**.  
Synchronized behavior in cyber attacks.
  - ACM SIGSAC Conference on **Computer and Communications Security** (CCS), 2015 : two papers
  - **Network and Distributed System Security** Symposium (NDSS), 2016: one paper
  - The 31st Annual **Computer Security Applications** Conference (ACSAC), 2015: one paper
  - IEEE Transactions on **Information Forensics and Security** (TIFS), 2016: one paper
  - International Journal of **Digital Crime and Forensics** (IJDCF), 2016: one paper
  - ...

# Impact

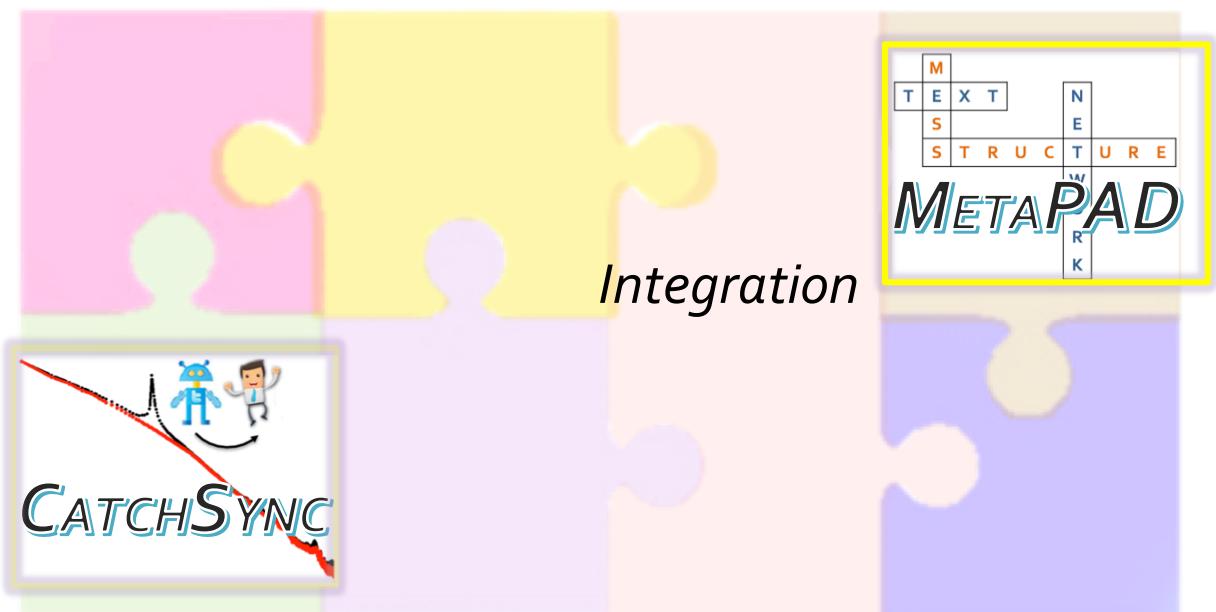
- Cited by **76**; WWW'14: cited by **21**. TKDD'16: cited by **21**. Synchronized behavior in cyber attacks.
- Taught in
  - CMU 15-826: [Multimedia Databases and Data Mining](#)
  - UMich EECS 598: [Graph Mining and Exploration at Scale](#)
  - ASONAM'16 Tutorial: “[Identifying Malicious Actors on Social Media](#)” by S. Kumar, F. Spezzano, V.S. Subrahmanian
- Endless games! First proposed **Camouflage** in PAKDD'14.
  - Cited by **47**.
  - Cited by *KDD'16 Best Research Paper*: the authors (B. Hooi *et al.*) provided theoretical bounds to prevent the camouflage.

# Outline

*Intelligence:*  
Behavior prediction  
and recommendation

*Trustworthiness:*  
Suspicious behavior  
detection

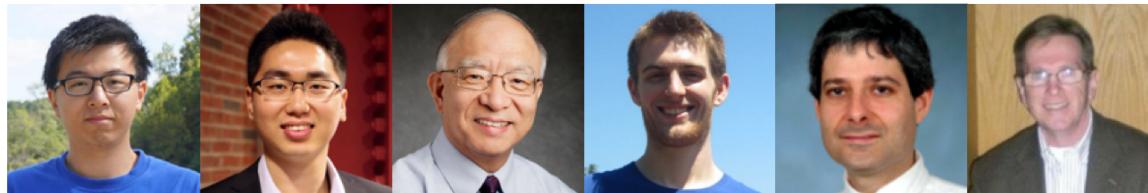
Social contexts      Spatiotemporal contexts      Behavioral content



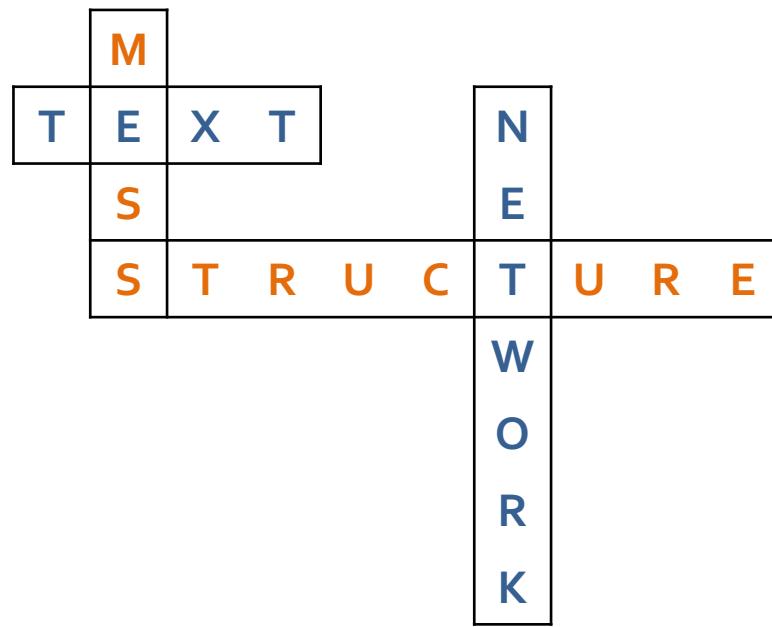
# MetaPAD: Meta Pattern-driven Attribute Discovery from Massive Text Corpora

Joint work with Jingbo Shang, Xiang Ren, Jiawei Han (UIUC),  
Taylor Cassidy, Lance Kaplan, and Timothy Hanratty (US Army Research Lab)

**ACM SIGKDD 2017**



# Motivation

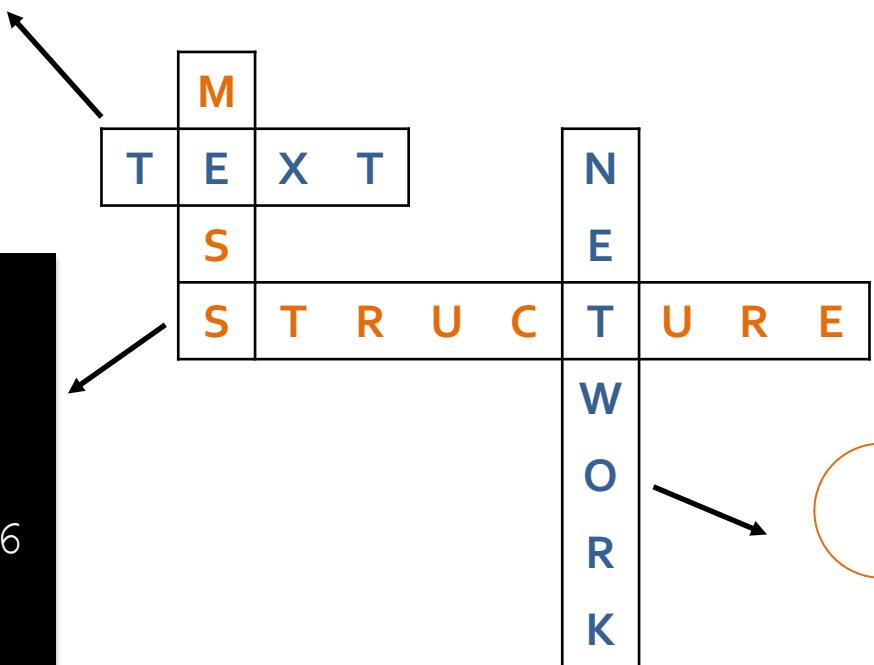


# Motivation

Given a sentence “President Blaise Compaoré’s government of **Burkina Faso** was founded...”, ...



was, 13897  
...  
president, 2769  
...  
government, 1886  
...  
blaise, 42  
...  
compaore, 15  
...



age: 65  
Blaise Compaoré:  
\$PERSON.POLITICIAN  
/ president  
Burkina Faso: \$COUNTRY  
population: 17 million

# Q'5: Attribute Discovery

Given a text corpus,

Task 1:  $\langle$ entity, attribute name, attribute value $\rangle$

$\langle$ Burkina Faso, president, Blaise Compaoré $\rangle$

$\langle$ Burkina Faso, population, 17 million $\rangle$

*Instance-level*

$\langle$ Blaise Compaoré, age, 65 $\rangle$

Task 2:  $\langle$ entity type, attribute name $\rangle$

$\langle$ \$COUNTRY, president $\rangle$

*Type-level*

$\langle$ \$COUNTRY, population $\rangle$

$\langle$ \$PERSON, age $\rangle$

# Literature: Task 1

Stanford OpenIE [ACL'15], AI2's Open IE-Ollie [EMNLP'12]:

Learn syntactic and lexical patterns of expressing relations

Ignore entity-typing information!

Input: "President Blaise Compaoré's government of Burkina Faso was founded..."

Output: {President Blaise Compaoré, **have**, government of Burkina Faso} 😞

# Literature: Task 2

Google's Biperpedia+ARI [VLDB'14, WWW'16], ReNoun [EMNLP'15]:

"president of united states"



"A of E", "E 's A", "E A", "A, E"

Query log: Highly constrained and unavailable

"Barack Obama, President of U.S.,"



"O, A of S,", "S A O"

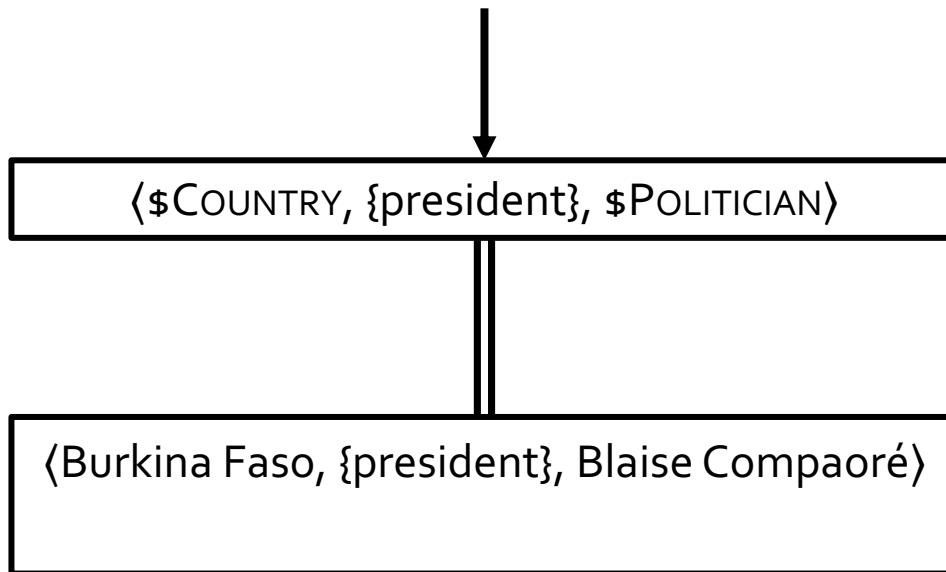
Annotated corpus: Domain limited and expensive

Input: "...Sunday night, Burkina Faso..." and the "A, E" pattern

Output: <\$COUNTRY, Sunday night> ☺

# Our Methodology (M5)

(#1) "President Blaise Compaoré's government of Burkina Faso was founded..."



# Our Methodology (M5)

- (#1) "President Blaise Compaoré's government of Burkina Faso was founded ..."
- (#2) "President Barack Obama's government of U.S. claimed that..."

Meta patterns:

[president \$PERSON.POLITICIAN 's government of \$LOCATION.COUNTRY] ...

Generate patterns with massive instances in the data

$\langle \$COUNTRY, \{president\}, \$POLITICIAN \rangle$

frequency↑

$\langle \text{Burkina Faso}, \{president\}, \text{Blaise Compaoré} \rangle$   
 $\langle \text{U.S.}, \{president\}, \text{Barack Obama} \rangle$

# Our Methodology (M5)

(#1) "President Blaise Compaoré's government of Burkina Faso was founded ..."

(#2) "President Barack Obama's government of U.S. claimed that..."

Meta patterns:

[president \$PERSON.POLITICIAN 's government of \$LOCATION.COUNTRY] ...

`($COUNTRY, {president}, $POLITICIAN)`

Generate massive triples by matching the meta patterns

`(Burkina Faso, {president}, Blaise Compaoré)`

`(U.S., {president}, Barack Obama)`

# Our Methodology (M5)

- (#1) "President Blaise Compaoré's government of Burkina Faso was founded ..."
- (#2) "President Barack Obama's government of U.S. claimed that..."
- (#3) "U.S. President Barack Obama visited ..."

Meta patterns:

〔president \$PERSON.POLITICIAN 's government of \$LOCATION.COUNTRY〕 was founded...  
〔\$LOCATION.COUNTRY president \$PERSON.POLITICIAN〕 ...

(\$COUNTRY, {president}, \$POLITICIAN)

frequency↑↑

Group synonymous patterns by massive triples

〈Burkina Faso, {president}, Blaise Compaoré〉  
〈U.S., {president}, Barack Obama〉

# Our Methodology (M5)

- (#1) "President Blaise Compaoré's government of Burkina Faso was founded ..."
- (#2) "President Barack Obama's government of U.S. claimed that..."
- (#3) "U.S. President Barack Obama visited ..."

Meta patterns:

[president \$PERSON.POLITICIAN 's government of \$LOCATION.COUNTRY was founded...]  
[\$LOCATION.COUNTRY president \$PERSON.POLITICIAN] ...

$\langle \$COUNTRY, \{president\}, \$POLITICIAN \rangle$

Adjust entity types in meta patterns  
for appropriate granularity with triples

$\langle Burkina Faso, \{president\}, Blaise Compaoré \rangle$   
 $\langle U.S., \{president\}, Barack Obama \rangle$

# Our Meta-Pattern Methodology

- (#1) "President Blaise Compaoré's government of Burkina Faso was founded ..."
- (#2) "President Barack Obama's government of U.S. claimed that..."
- (#3) "U.S. President Barack Obama visited ..."

Meta patterns:

*Meta pattern segmentation*

president \$PERSON.POLITICIAN 's government of \$LOCATION.COUNTRY was founded...  
\$LOCATION.COUNTRY president \$PERSON.POLITICIAN ...

$\langle \$COUNTRY, \{president\}, \$POLITICIAN \rangle$

*Joint  
extraction*

$\langle Burkina Faso, \{president\}, Blaise Compaoré \rangle$   
 $\langle U.S., \{president\}, Barack Obama \rangle$

*Adjust types for appropriate granularity*

*Group synonymous meta patterns*

# Our Meta-Pattern Methodology

- (#1) "President Blaise Compaoré's government of Burkina Faso was founded ..."
- (#2) "President Barack Obama's government of U.S. claimed that..."
- (#3) "U.S. President Barack Obama visited ..."

*Meta pattern segmentation*

Meta patterns:

president \$PERSON.POLITICIAN 's government of \$LOCATION.COUNTRY was founded...  
[\$LOCATION.COUNTRY president \$PERSON.POLITICIAN] ...

No heavy annotation required  
No domain knowledge required  
No query log required

*Adjust type appropriate granularity*

if we can recognize and type the entities in the same manner... *synonymous meta patterns*

`<Burkina Faso, {president}, Blaise Compaoré>`  
`<U.S., {president}, Barack Obama>`

# Han's Group Strength in Text Mining

“President Blaise Compaoré’s government of Burkina Faso was founded ...”

*Phrase mining (SegPhrase by Liu and Han et al. SIGMOD’15)*

“president blaise\_compaoré ’s government of burkina\_faso was founded ...”

*Entity recognition and typing with Distant Supervision  
(ClusType by Ren and Han et al. KDD’15)*

“president \$PERSON ’s government of \$LOCATION was founded ...”

*Fine-grained typing (PLE by Ren and Han et al. KDD’16)*

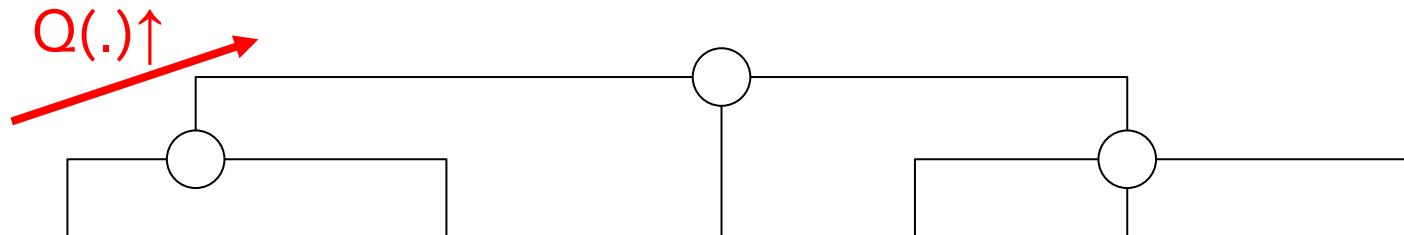
“president \$PERSON.POLITICIAN ’s government of \$LOCATION.COUNTRY was founded ...”

# Meta-Pattern Quality Assessment and Segmentation

A rich set of features:

- ✓ Frequency
- ✓ Concordance: "\$PERSON 's wife"
- ✓ Completeness: "\$COUNTRY president" vs "\$COUNTRY president \$POLITICIAN"
- ✓ Informativeness: "\$PERSON and \$PERSON" vs "\$PERSON 's wife, \$PERSON"

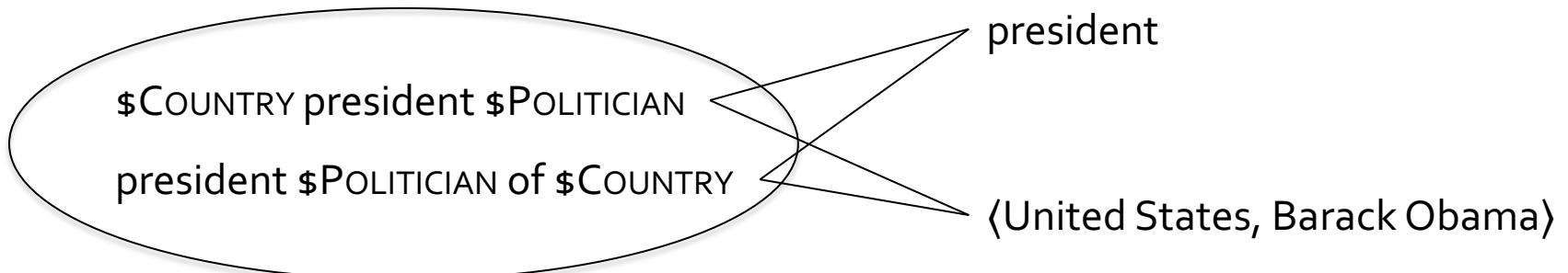
Regression Q(.): random forest with only 300 labels



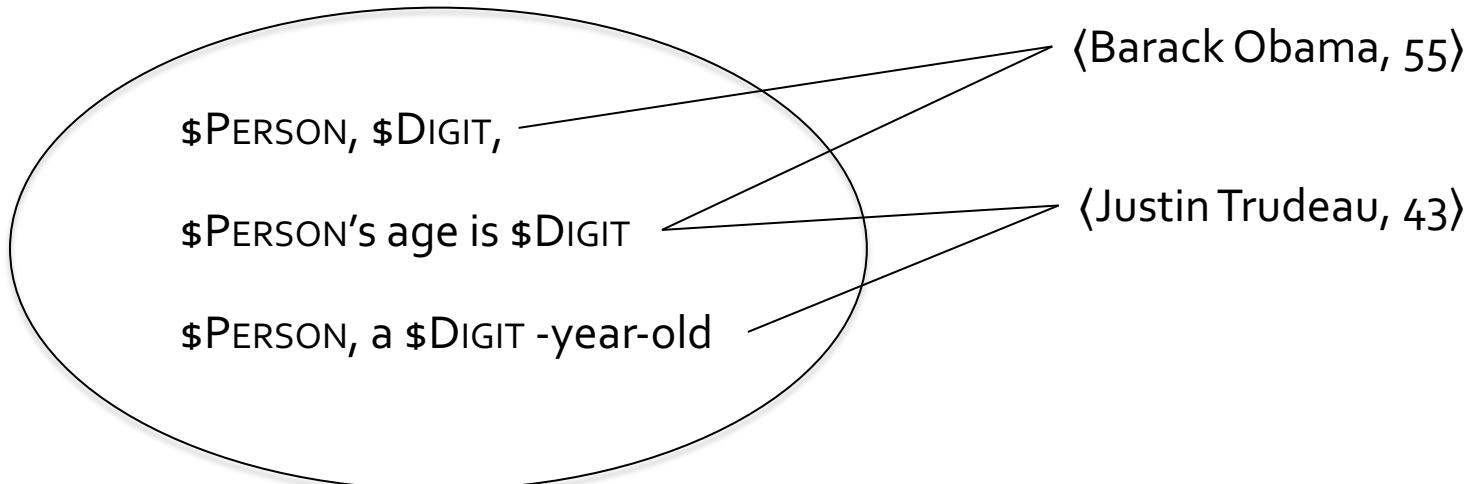
president \$PERSON.POLITICIAN 's government of \$LOCATION.COUNTRY

# Grouping Synonymous Patterns

$\langle \$COUNTRY, \text{president}, \$POLITICIAN \rangle$



$\langle \$PERSON, \{\text{age}, \text{-year-old}\}, \$DIGIT \rangle$



# Adjusting Types in Meta Patterns for Appropriate Granularity

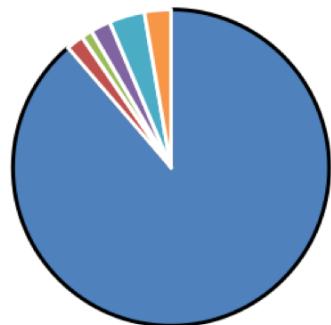
\$PERSON, \$DIGIT,

\$PERSON's age is \$DIGIT

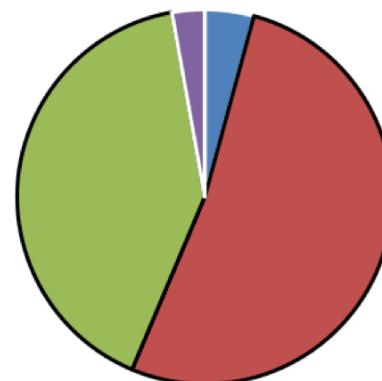
\$PERSON, a \$DIGIT -year-old

\$COUNTRY president \$POLITICIAN

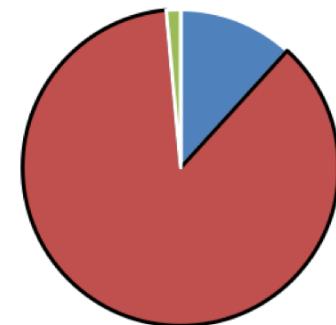
president \$POLITICIAN of \$COUNTRY



- \$PERSON
- \$ATTACKER
- \$ARTIST
- \$ATHLETE
- \$LOCATION
- \$COUNTRY
- \$ETHNICITY
- \$VICTIM
- \$CITY



- \$LOCATION
- \$COUNTRY
- \$ETHNICITY
- \$VICTIM
- \$CITY



- \$PERSON
- \$COUNTRY
- \$ARTIST

# Results in General Domain

Meta patterns	Entity	Attribute value
\$COUNTRY President \$POLITICIAN	United States	Barack Obama
\$COUNTRY's president \$POLITICIAN	Russia	Vladimir Putin
President \$POLITICIAN of \$COUNTRY	France	Francois Hollande
...	...	...
\$POLITICIAN's government of \$COUNTRY	Burkina Faso	Blaise Compaoré

Meta patterns	Entity	Attribute value
\$COMPANY CEO \$PERSON	Apple	Tim Cook
\$COMPANY chief executive \$PERSON	Facebook	Mark Zuckerberg
\$PERSON, the \$COMPANY CEO,	Hewlett-Packard	Carly Fiorina
...	...	...
\$COMPANY former CEO \$PERSON	Infor	Charles Phillips
\$PERSON, the \$COMPANY former CEO,	Afghan Citadel	Roya Mahboob

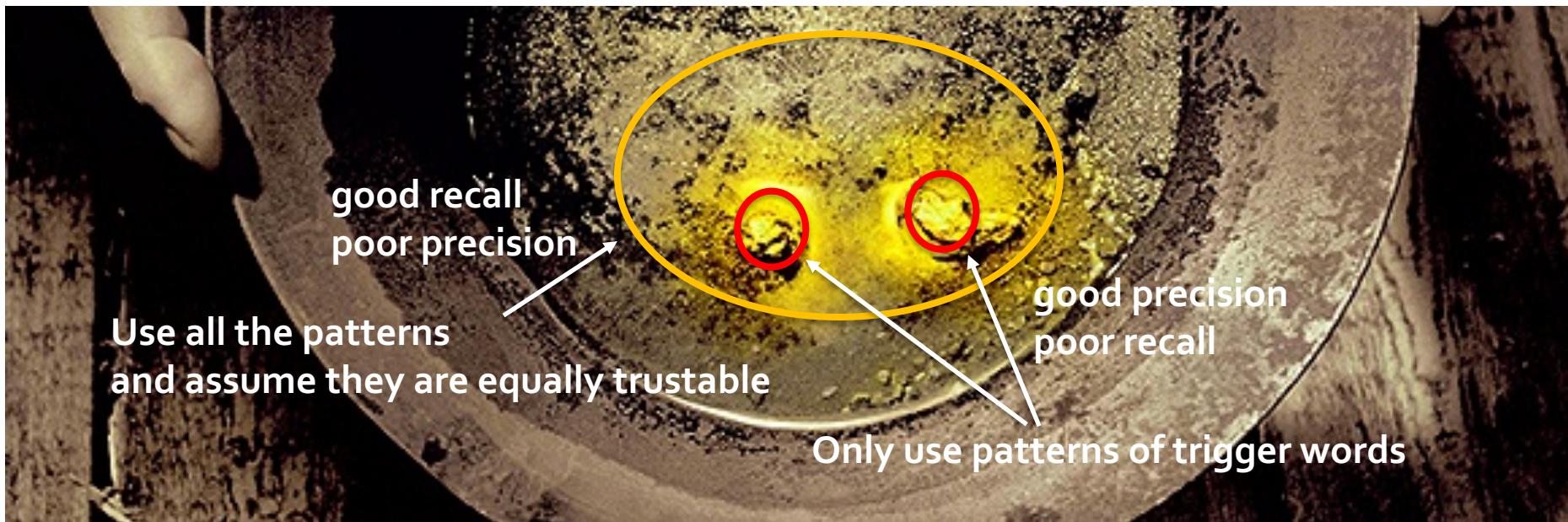
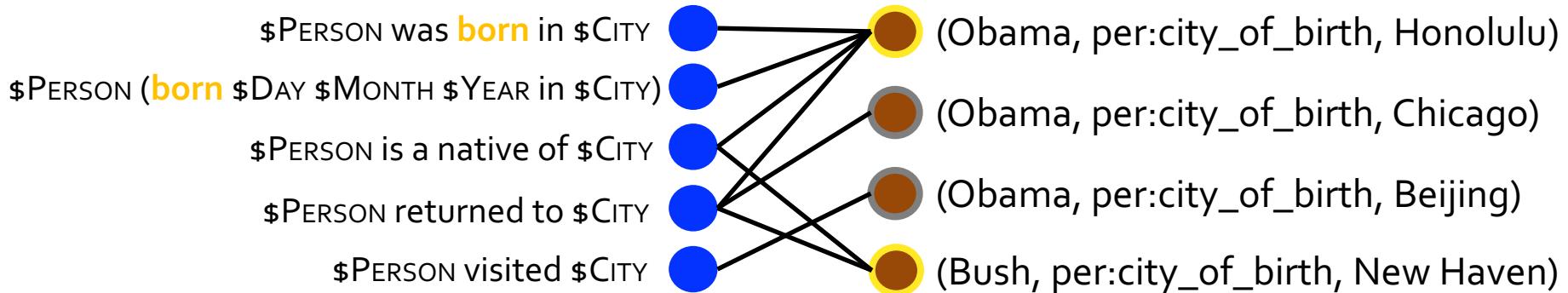
# Results in Biomedical Domain

Meta patterns	Entity	Attribute value
\$TREATMENT was used to treat \$DISEASE \$DISEASE using the \$TREATMENT \$TREATMENT has been used to treat \$DISEASE \$TREATMENT of patients with \$DISEASE ...	zoledronic acid therapy	Paget's disease of bone
	bisphosphonates	osteoporosis
	calcitonin	Paget's disease of bone
	calcitonin	osteoporosis
	...	...
	...	...
Meta patterns	Entity	Attribute value
\$BACTERIA was resistant to \$ANTIBIOTICS \$BACTERIA are resistant to \$ANTIBIOTICS \$BACTERIA is the most resistant to \$ANTIBIOTICS ...	corynebacterium striatum BM4687	gentamicin
	corynebacterium striatum BM4687	tobramycin
	methicillin-susceptible S aureus	vancomycin
	multidrug-resistant enterobacteriaceae	gentamicin
	...	...
	...	...

# Experimental Results

F1 score	<code>(entity type, attribute name)</code>	<code>(entity, attribute name, attribute value)</code>
Baselines		Stanford's OpenIE: 0.035
		AI2's Ollie: 0.131
	Biperpedia: 0.324	Google's ReNoun: 0.309
+Segmentation	+40.0%	+19.4%
+Type Adjustment	+6.5%	+15.0%
+Synonymous	+2.6%	
All	0.495 <b>relatively +52.9%</b>	0.424 <b>relatively +37.3%</b>

# On-going: Reliable Attribute Discovery

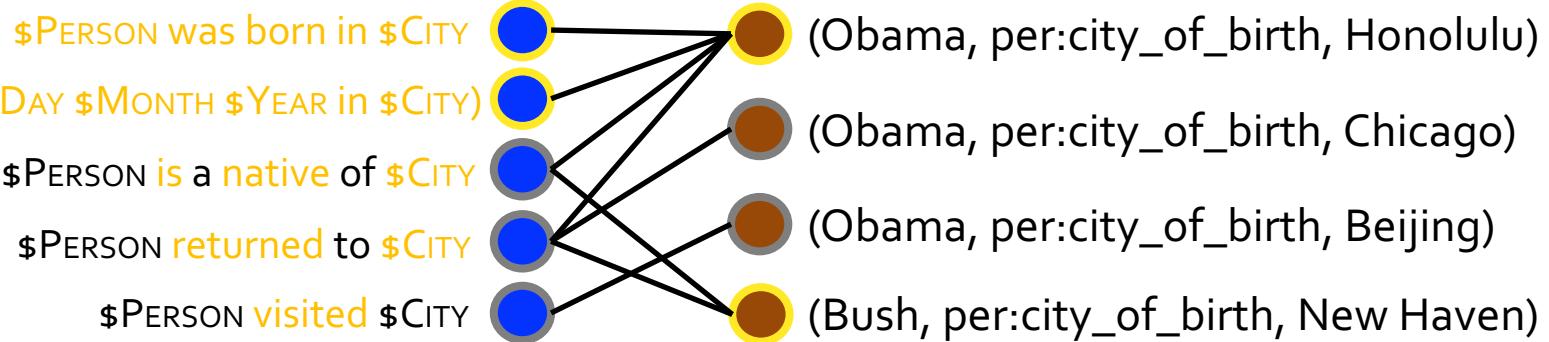


# On-going: Reliable Attribute Discovery



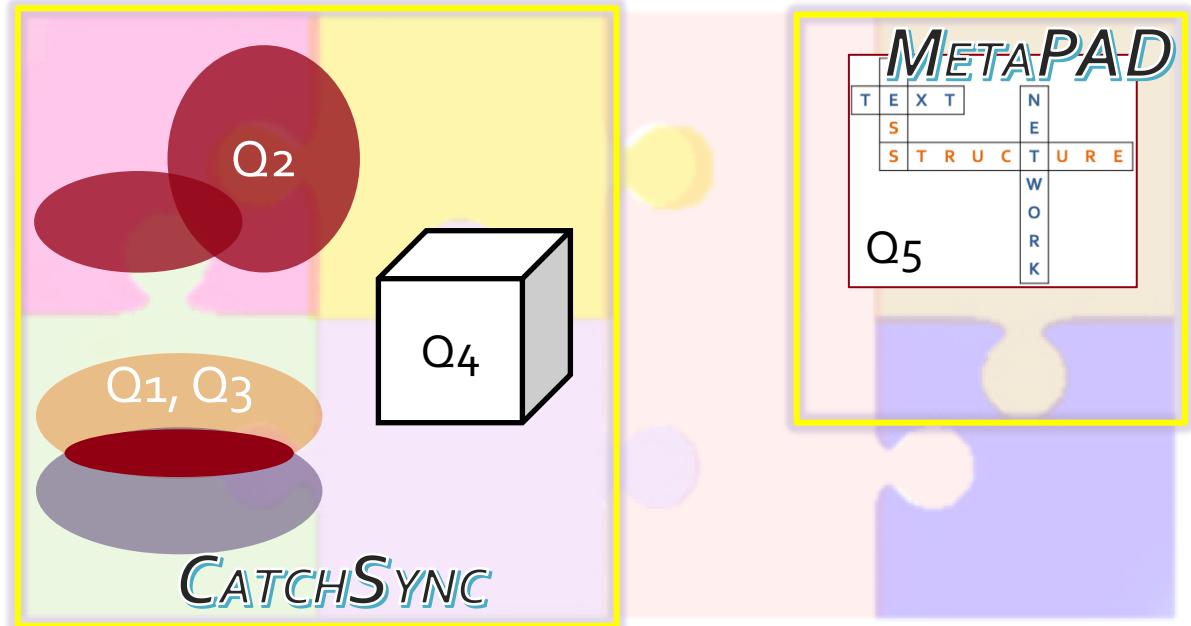
Modeling the **reliability** of extractors:

- Not every extraction is 100% correct.
- The meta patterns are not equally trustable.



# Summary

Social contexts      Spatiotemporal contexts      *Integration*      Behavioral content

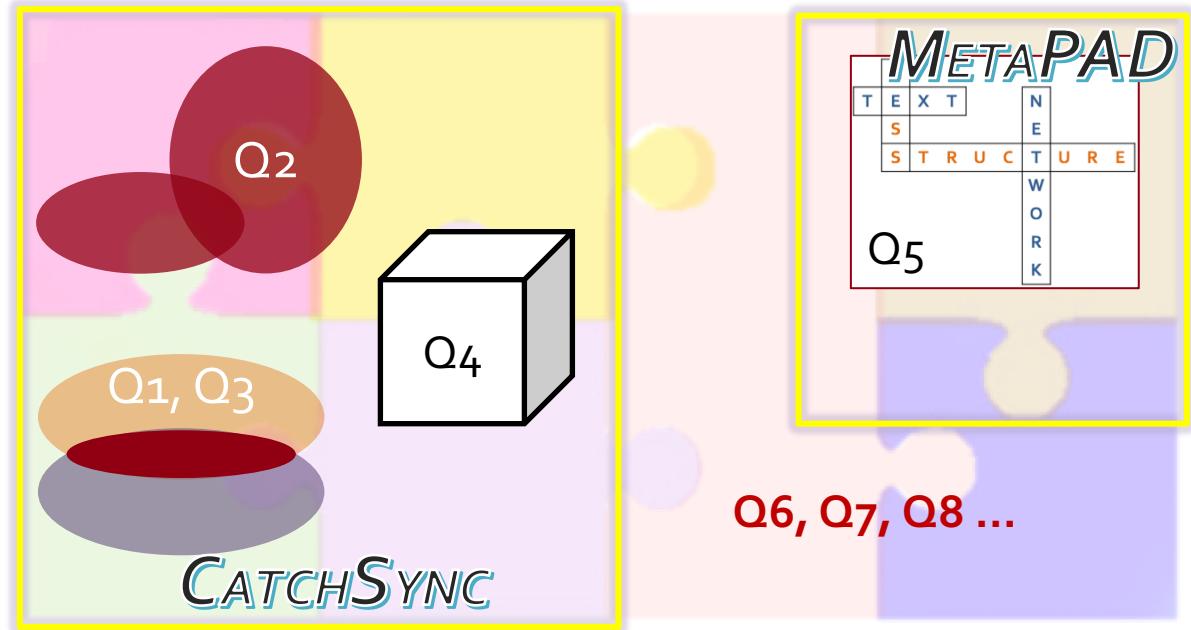


*Intelligence:*  
Behavior prediction  
and recommendation

*Trustworthiness:*  
Suspicious behavior  
detection

# Summary

Social contexts      Spatiotemporal contexts      *Integration*      Behavioral content



## Expedition!!!

- Data sciences + ML, AI, NLP, CyberSecurity ...
- Interdisciplinary research: Sociology, Psychology ...
- Transformative technologies: Change the games ☺

# References

- Breese, Heckerman, and Kadie. Empirical analysis of predictive algorithms for collaborative filtering. In the 4<sup>th</sup> conference on Uncertainty in artificial intelligence (UAI), 1998.
- Getoor and Sahami. Using probabilistic relational models for collaborative filtering. In Workshop on Web Usage Analysis and User Profiling (WEBKDD), 1999.
- Herlocker, Konstan, and Riedl. Explaining collaborative filtering recommendations. ACM Conference on Computer Supported Cooperative Work (CSCW), 2000.
- Herlocker, Konstan, Terveen, and Riedl. Evaluating collaborative filtering recommender systems. ACM Transactions on Information Systems (TOIS) 22, no.1, 2004.
- Yu and Ji. Unsupervised Person Slot Filling based on Graph Mining. ACL, 2016.
- Hu, Koren, and Volinsky. Collaborative filtering for implicit feedback datasets. In International Conference on Data Mining (ICDM), pages 263–272, 2008.
- Koren. Factorization meets the neighborhood: A multifaceted collaborative filtering model. ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD), 2008.
- Koren, Bell, and Volinsky. Matrix factorization techniques for recommender systems. Computer, 2009.
- Liu and Yang. Eigenrank: A ranking-oriented approach to collaborative filtering. International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR), 2008.
- Liu, Zhao, and Yang. Probabilistic latent preference analysis for collaborative filtering. ACM International Conference on Information and Knowledge Management (CIKM), 2009.
- Salganik, Dodds, and Watts. Experimental Study of Inequality and Unpredictability in an Artificial Cultural Market. Science, 2006.

# References

- Sarwar, Karypis, Konstan, and Riedl. Item-based collaborative filtering recommendation algorithms. International conference on World Wide Web (WWW), 2001.
- Burke. Hybrid recommender systems: Survey and experiments. User modeling and user-adapted interaction, 2002.
- Ma, Zhou, Lyu, and King. Improving recommender systems by incorporating social contextual information. ACM Transactions on Information Systems (TOIS), 2011.
- Ma, King, and Lyu. Learning to recommend with explicit and implicit social relations. ACM Transactions on Intelligent Systems and Technology (TIST), 2011.
- Tang, Hu, and Liu. Social recommendation: a review. Social Network Analysis and Mining, 2013.
- Han, Zhang, Ghalwash, Vucetic, and Obradovic. Joint Learning of Representation and Structure for Sparse Regression on Graphs. SIAM International Conference on Data Mining (SDM), 2016.
- Egele, Stringhini, Kruegel, and Vigna. COMPA: Detecting Compromised Accounts on Social Networks. The Network and Distributed System Security Symposium (NDSS), 2013.
- Yang, Wilson, Wang, Gao, Zhao, and Dai. Uncovering social network sybils in the wild. ACM Transactions on Knowledge Discovery from Data (TKDD), 2014.
- Viswanath, Bashir, Crovella, Guha, Gummadi, Krishnamurthy, and Mislove. Towards detecting anomalous user behavior in online social networks. USENIX Security Symposium (USENIX Security), 2014.
- Faloutsos, Faloutsos, and Faloutsos. On power-law relationships of the internet topology. ACM SIGCOMM computer communication review, 1999.
- Chung, and Lu. The average distances in random graphs with given expected degrees. Proceedings of the National Academy of Sciences, 2002.

# References

- Kleinberg. Authoritative sources in a hyperlinked environment. *Journal of ACM*, 1999.
- Hooi, Song, Beutel, Shah, Shin, and Faloutsos. Fraudar: Bounding graph fraud in the face of camouflage. *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2016.
- Angeli, Premkumar, and Manning. Leveraging linguistic structure for open domain information extraction. *Linguistics*, 2015.
- Schmitz, Bart, Soderland, and Etzioni. Open language learning for information extraction. *Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP)*, 2012.
- Gupta, Halevy, Wang, Whang, and Wu. Biperpedia: An ontology for search applications. *Proceedings of the VLDB Endowment (VLDB)*, 2014.
- Halevy, Noy, Sarawagi, Whang, and Yu. Discovering Structure in the Universe of Attribute Names. *International Conference on World Wide Web (WWW)*, 2016.
- Yahya, Whang, Gupta, and Halevy. ReNoun: Fact Extraction for Nominal Attributes. *Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP)*, 2014.
- Liu, Shang, Wang, Ren, and Han. Mining quality phrases from massive text corpora. *ACM SIGMOD International Conference on Management of Data (SIGMOD)*, 2015.
- Ren, El-Kishky, Wang, Tao, Voss, and Han. Clustype: Effective entity recognition and typing by relation phrase-based clustering. *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2015.

# Acknowledgement

