# Anomaly detection in large graphs

*Christos Faloutsos*
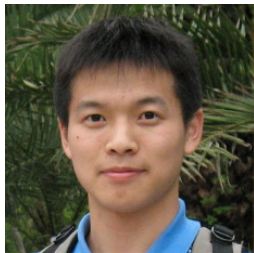
CMU
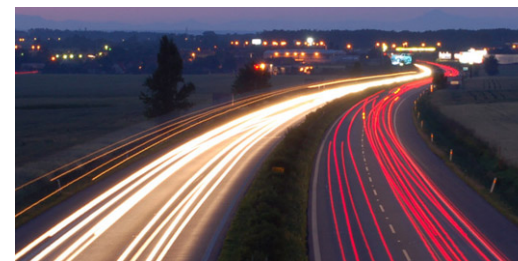
# **Thank you!**

- Annette Jiang (IEEE)

- Evan Butterfield (IEEE)

- Lei Li

# Roadmap

➡ • Introduction – Motivation

    – Why study (big) graphs?

• Part#1: Patterns in graphs
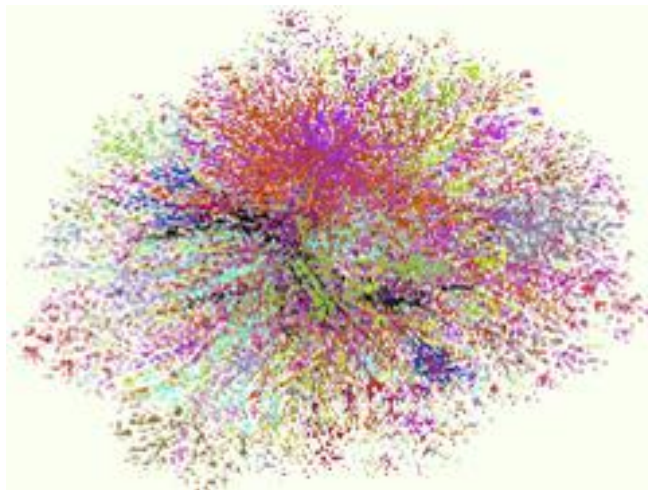
• Part#2: time-evolving graphs; tensors
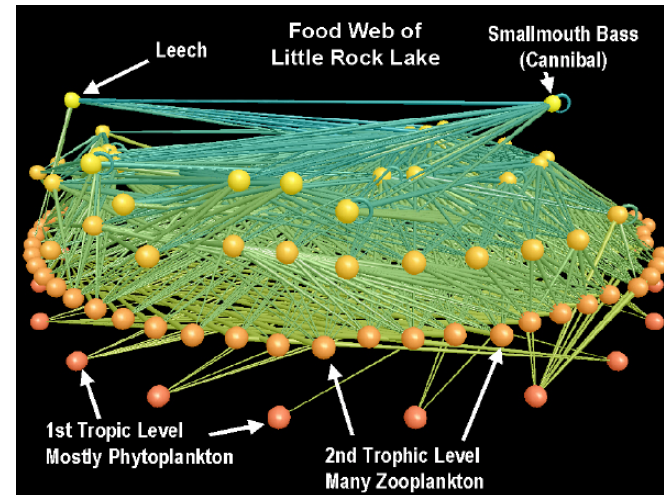
• Conclusions

# Graphs - why should we care?

>$10B; ~1B users

# Graphs – why should we care?
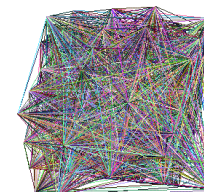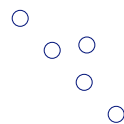


## Internet Map
## [lumeta.com]

## Food Web
## [Martinez '91]

# Graphs - why should we care?

- web-log ('blog') news propagation **YAHOO! BLOG**
- computer network security: email/IP traffic and anomaly detection
- Recommendation systems **NETFLIX**
- ....

- Many-to-many db relationship -> graph

# Motivating problems

- P1: patterns? Fraud detection?

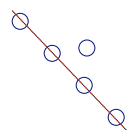- P2: patterns in time-evolving graphs / tensors

destination

source       time
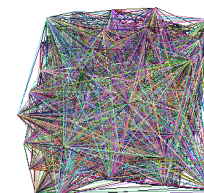
# Motivating problems

- P1: patterns? Fraud detection?

Patterns anomalies
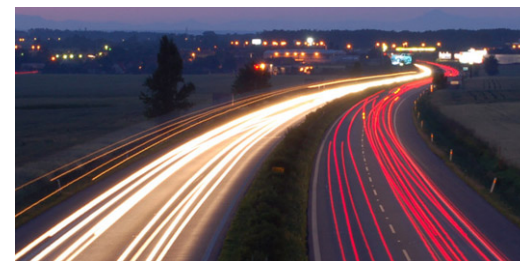
- P2: patterns in time-evolving graphs / tensors

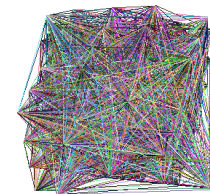destination
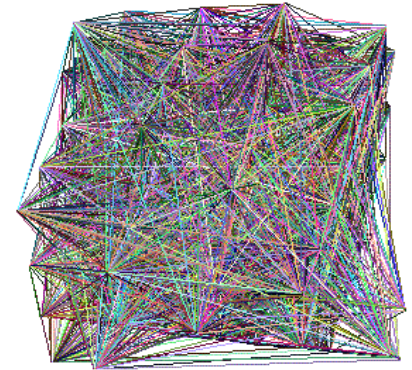
source    time

# Roadmap



- Introduction – Motivation
  - Why study (big) graphs?
- **→** Part#1: Patterns & fraud detection
- Part#2: time-evolving graphs; tensors
- Conclusions

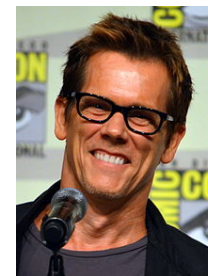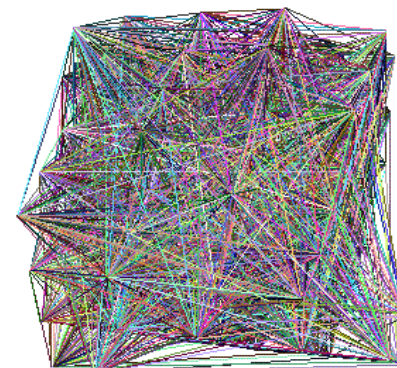# Part 1: Patterns, & fraud detection

# Laws and patterns

- Q1: Are real graphs random?
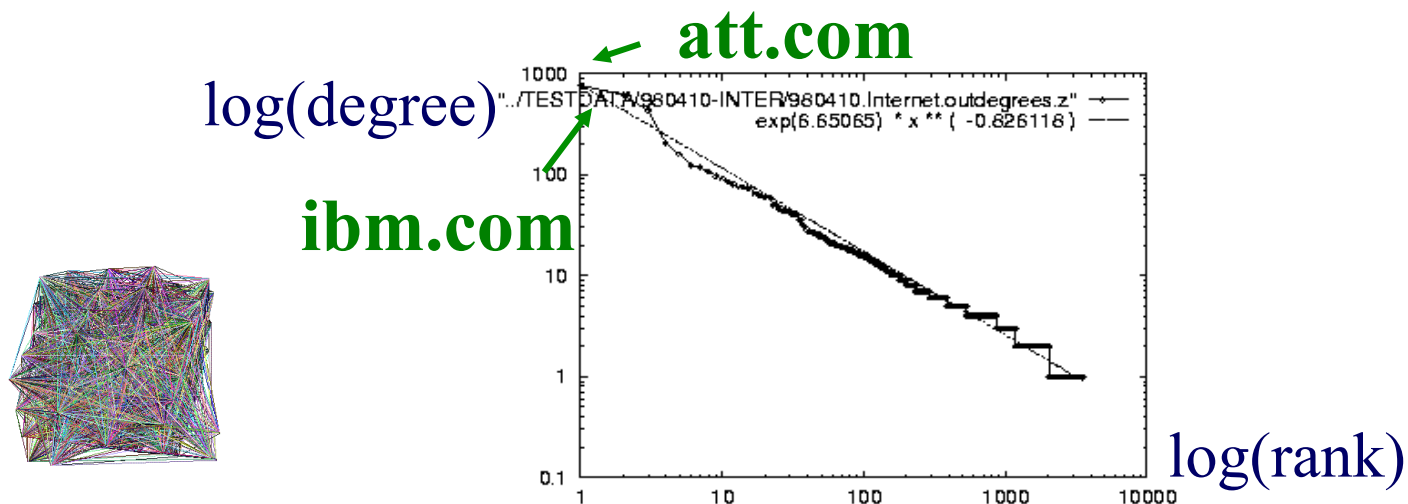
# Laws and patterns

- Q1: Are real graphs random?
- A1: NO!!
  - Diameter ('6 degrees'; 'Kevin Bacon')
  - in- and out- degree distributions
  - other (surprising) patterns
- So, let's look at the data

# Solution# S.1

- Power law in the degree distribution [Faloutsos x 3 SIGCOMM99]

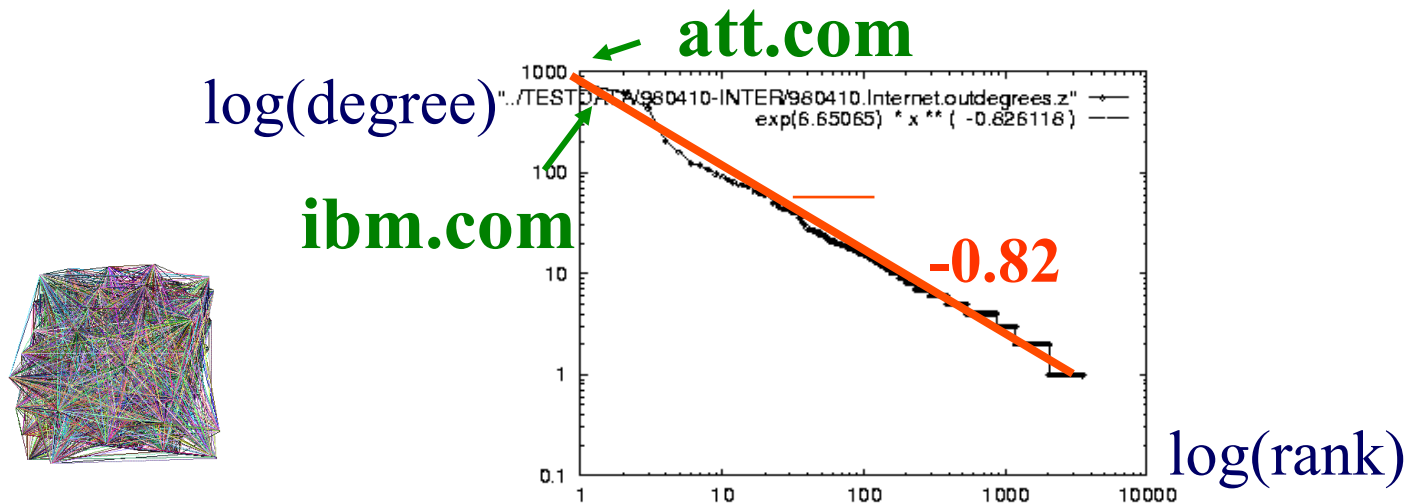**internet domains**

**att.com**

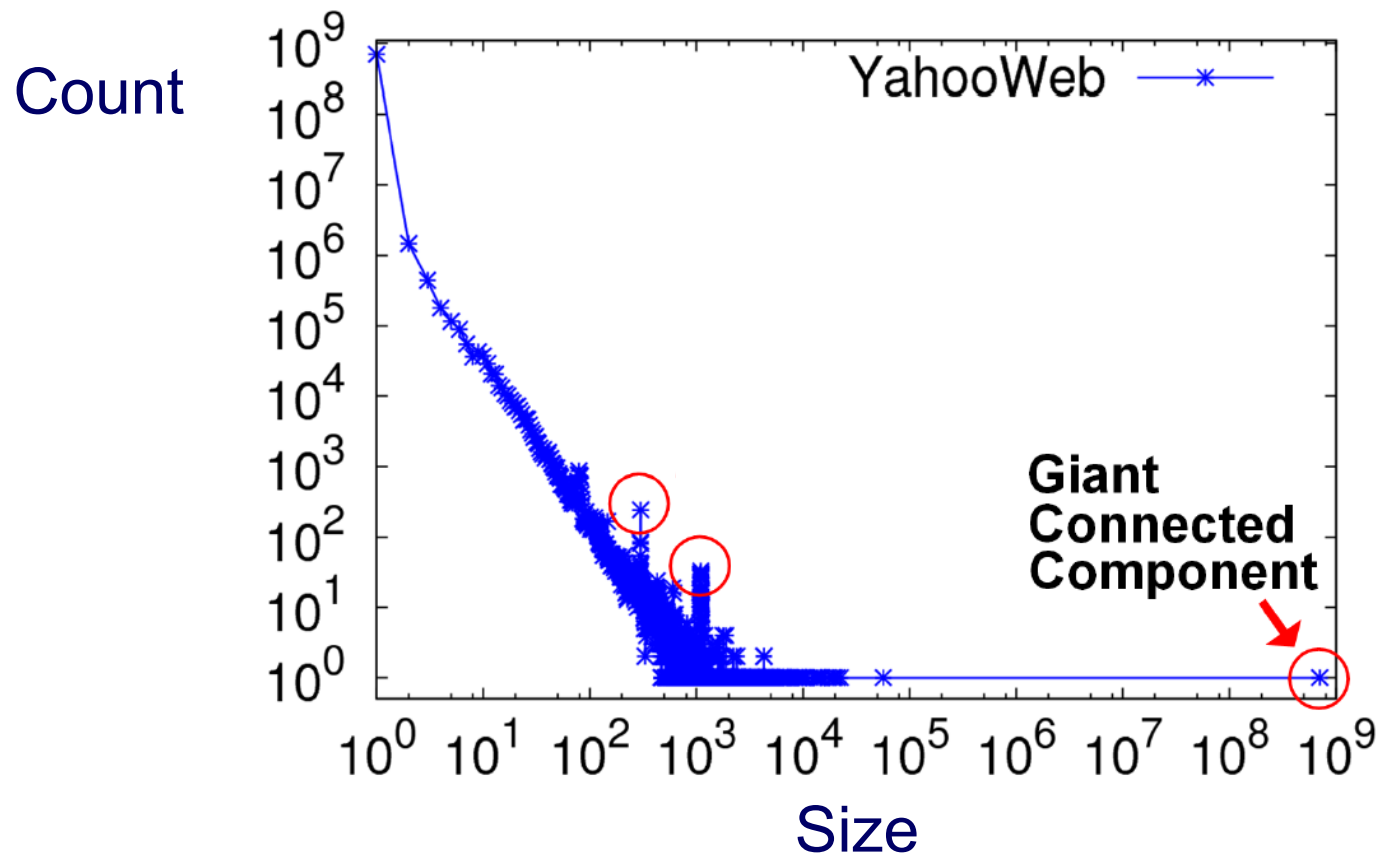log(degree)

**ibm.com**



log(rank)

# Solution# S.1

- Power law in the degree distribution [Faloutsos x 3 SIGCOMM99]

**internet domains**

# S2: connected component sizes

- Connected Components – 4 observations:
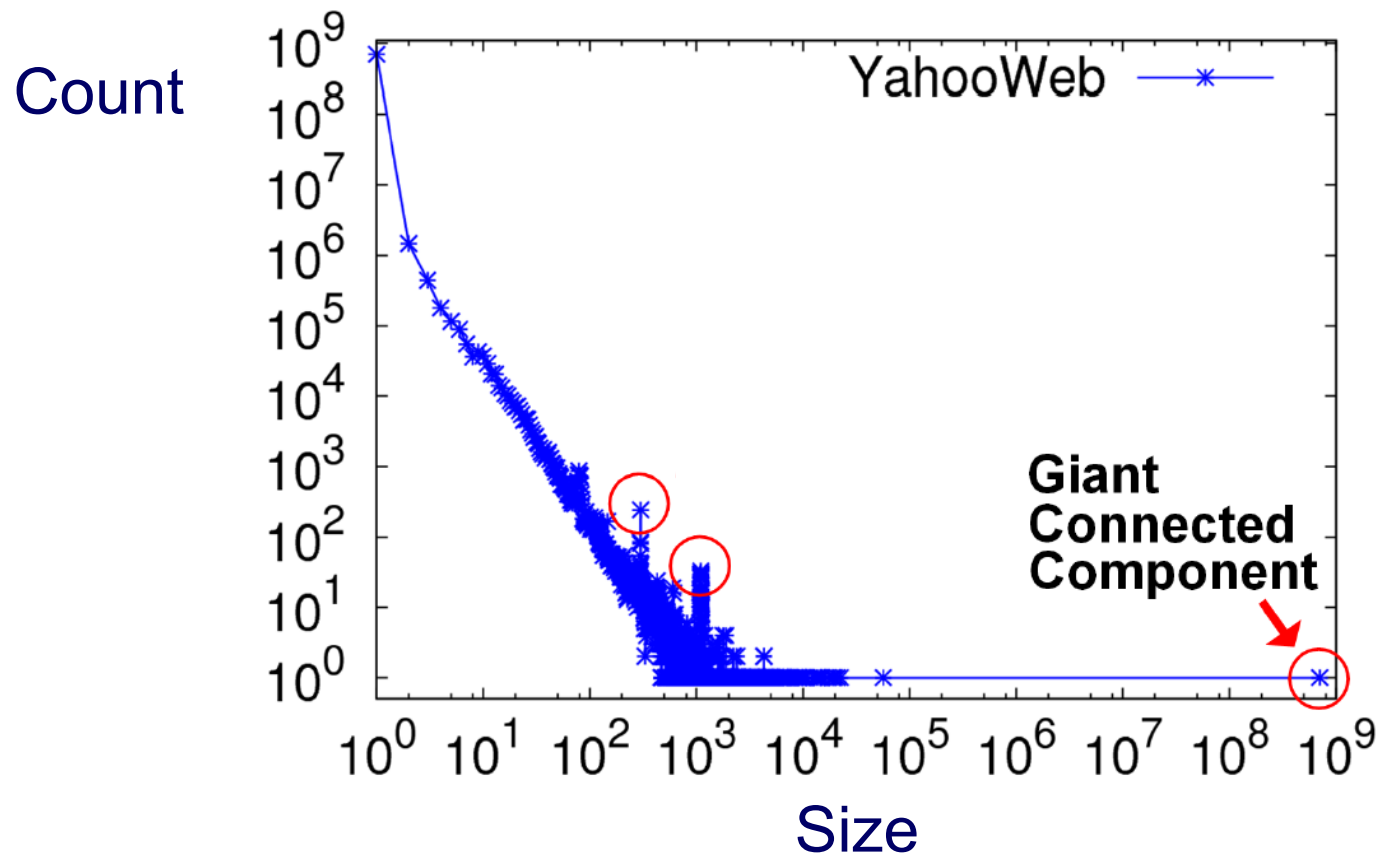


Count

1.4B nodes
6B edges

# S2: connected component sizes

- Connected Components



Count

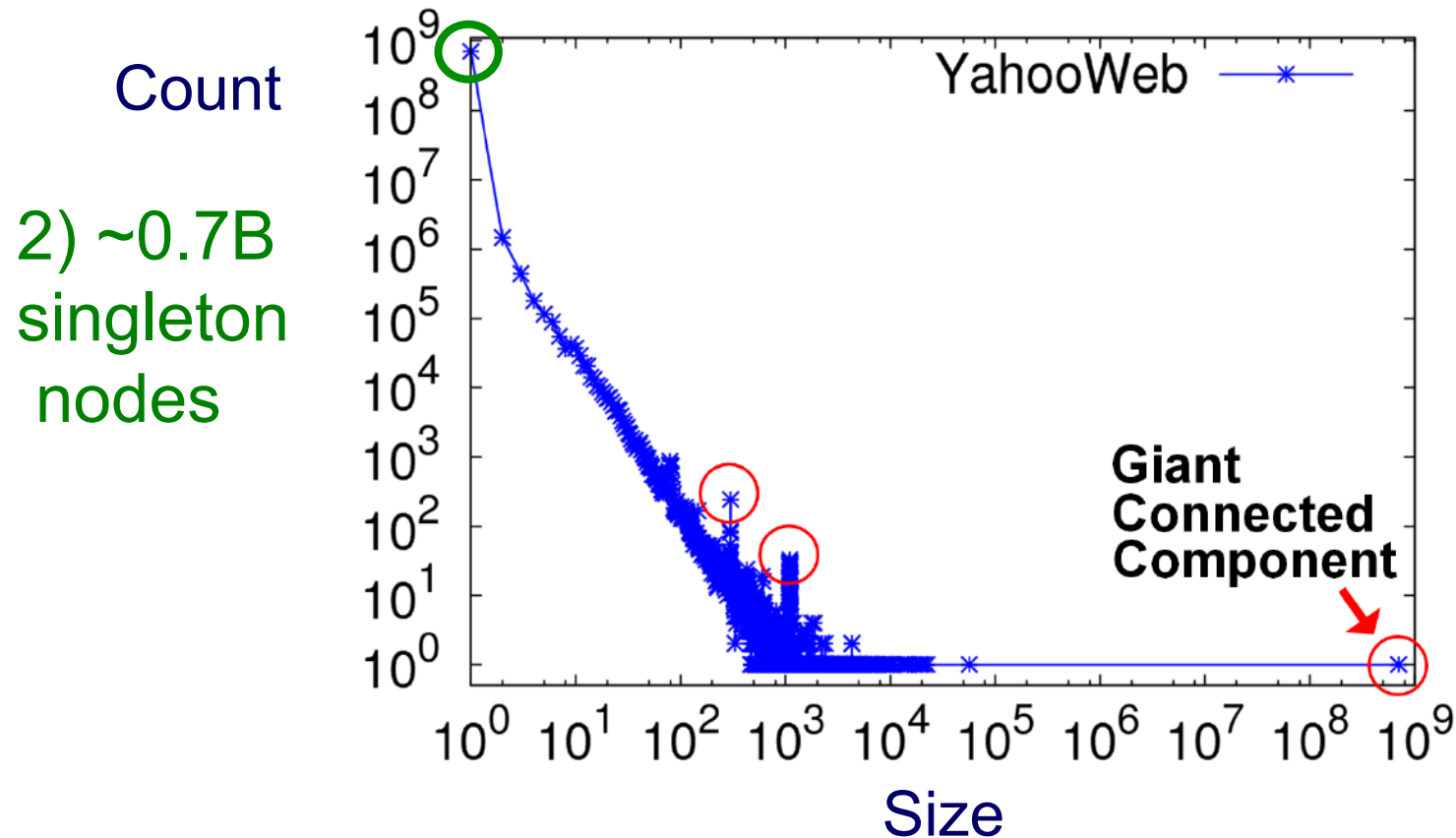1) 10K x larger than next

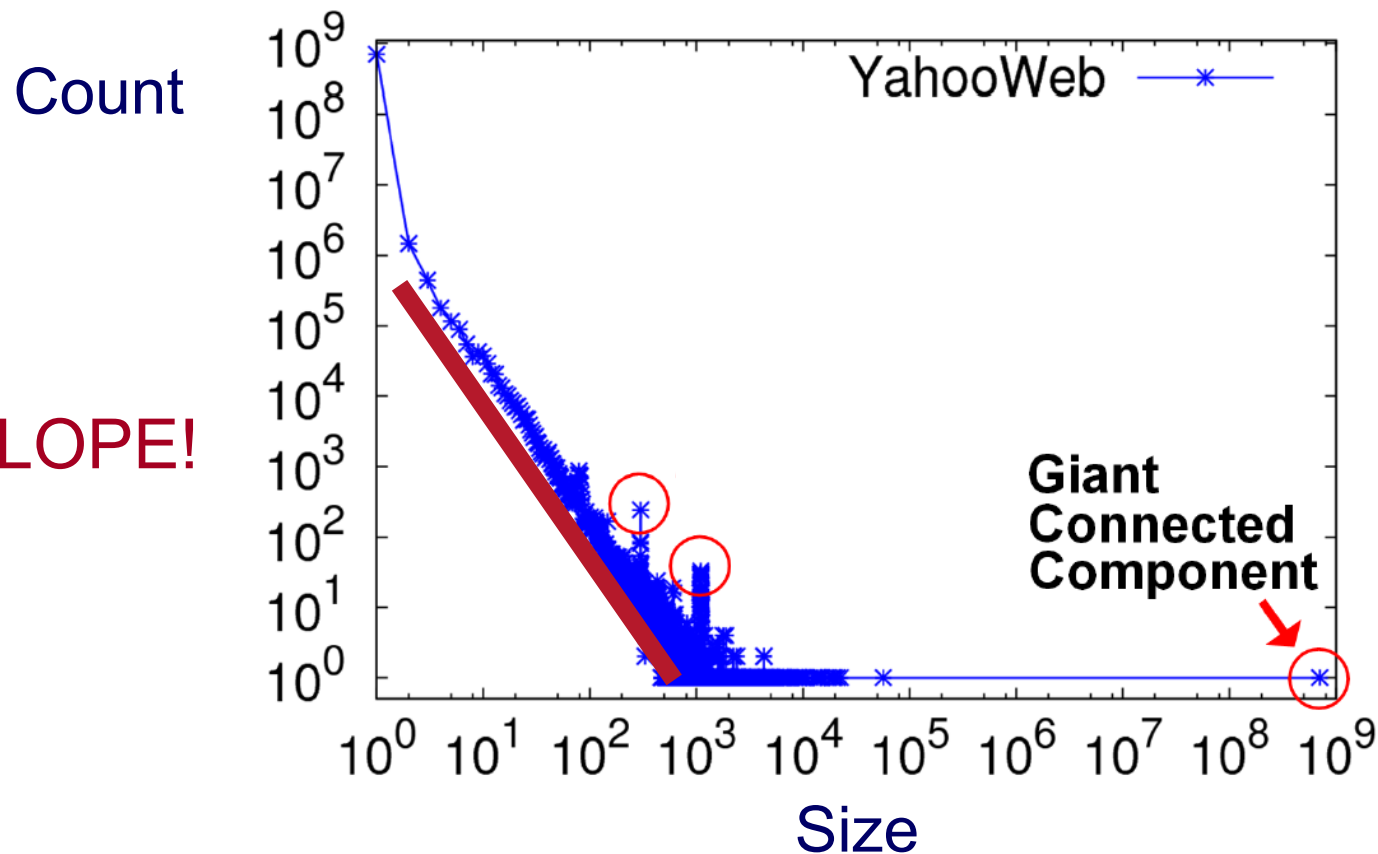# S2: connected component sizes

- Connected Components

**Count**

**2) ~0.7B singleton nodes**



**Size**

# S2: connected component sizes

- Connected Components



Count

3) SLOPE!

YahooWeb

Giant
Connected
Component

Size

# S2: connected component sizes

- Connected Components



Count

300-size cmpt X 500. Why?

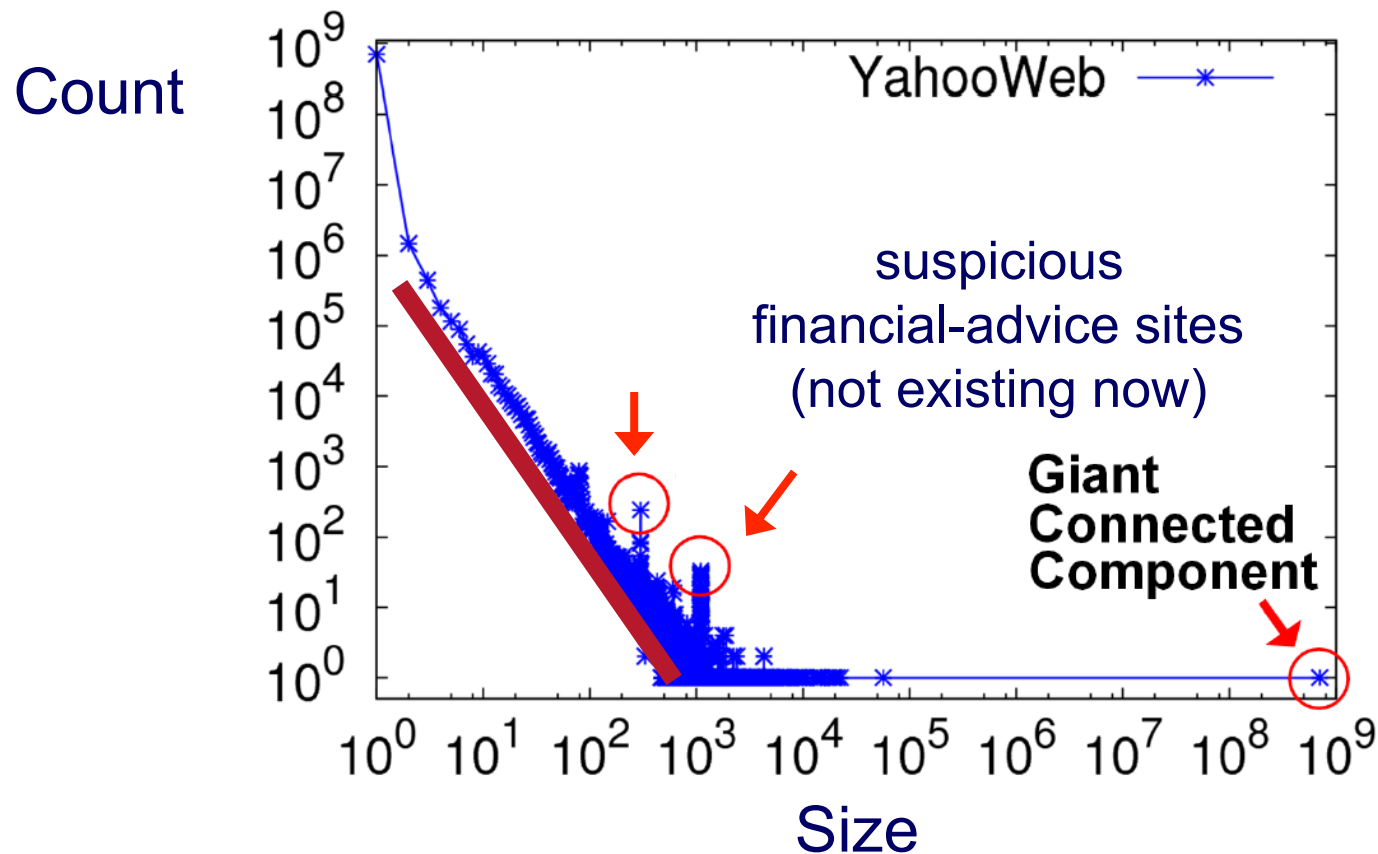1100-size cmpt X 65. Why?

**Giant Connected Component**

4) Spikes!

YahooWeb

Size

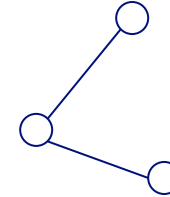# S2: connected component sizes

- Connected Components

# Roadmap



- Introduction – Motivation
- Part#1: Patterns in graphs
  - P1.1: Patterns: Degree; Triangles
  - P1.2: Anomaly/fraud detection
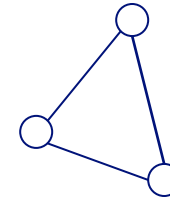- Part#2: time-evolving graphs; tensors
- Conclusions

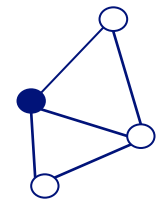# Solution# S.3: Triangle 'Laws'

- Real social networks have a lot of triangles
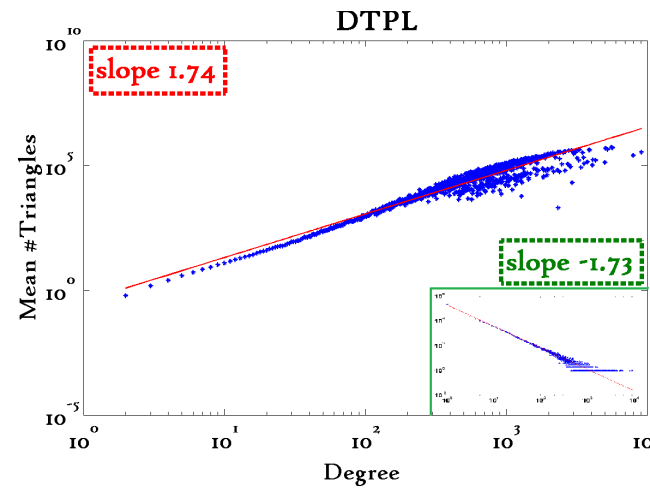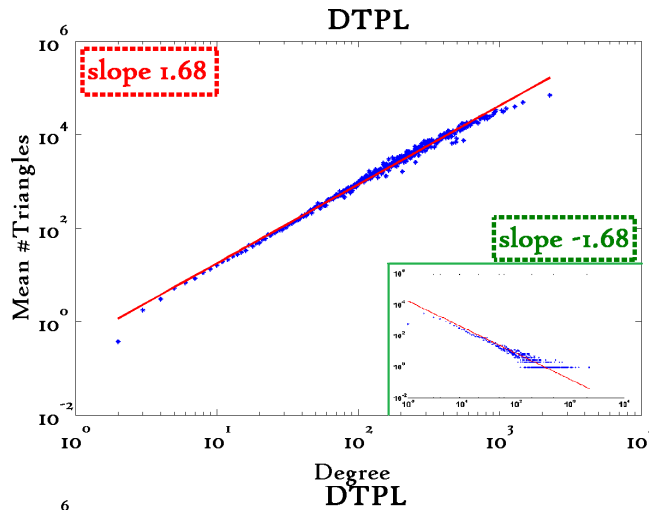
# Solution# S.3: Triangle 'Laws'

- ## Real social networks have a lot of triangles
  - Friends of friends are friends

- ## Any patterns?
  - 2x the friends, 2x the triangles ?

# Triangle Law: #S.3
## [Tsourakakis ICDM 2008]



Reuters

**DTPL**
slope 1.68
slope -1.68
Mean #Triangles
Degree

SN

**DTPL**
slope 1.74
slope -1.73
Mean #Triangles
Degree

Epinions

**DTPL**
slope 1.61
slope -1.59
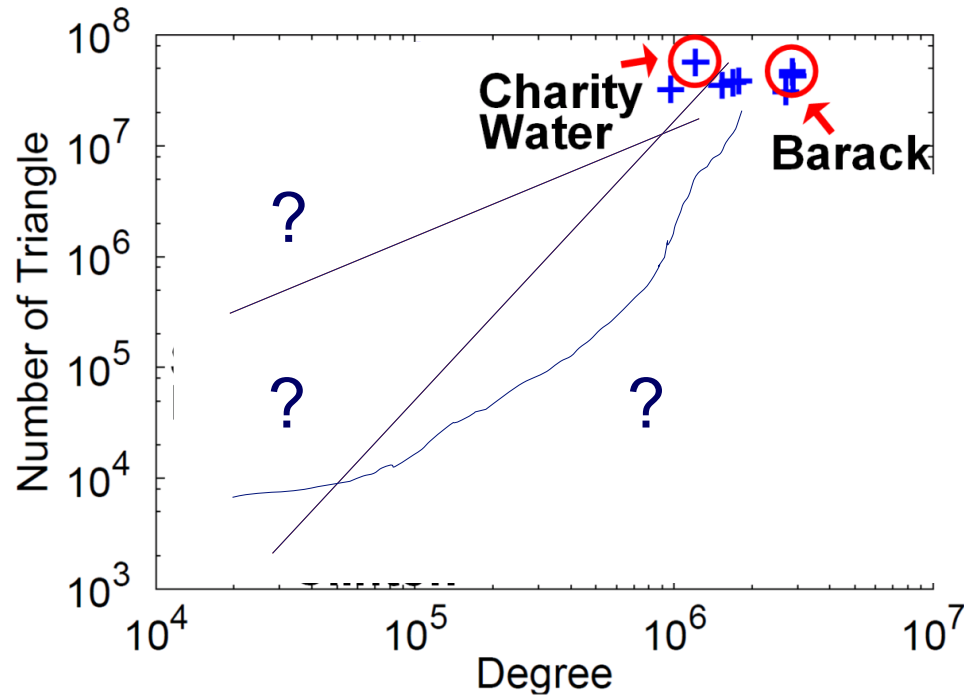Mean #Triangles
Degree

X-axis: degree
Y-axis: mean # triangles
$n$ friends -> ~$n^{1.6}$ triangles

# Triangle counting for large graphs?



Anomalous nodes in Twitter(~ 3 billion edges)
[U Kang, Brendan Meeder, +, PAKDD'11]

# Triangle counting for large graphs?



Anomalous nodes in Twitter(~ 3 billion edges)

[U Kang, Brendan Meeder, +, PAKDD'11]

# Triangle counting for large graphs?



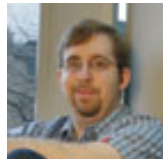Anomalous nodes in Twitter(~ 3 billion edges)
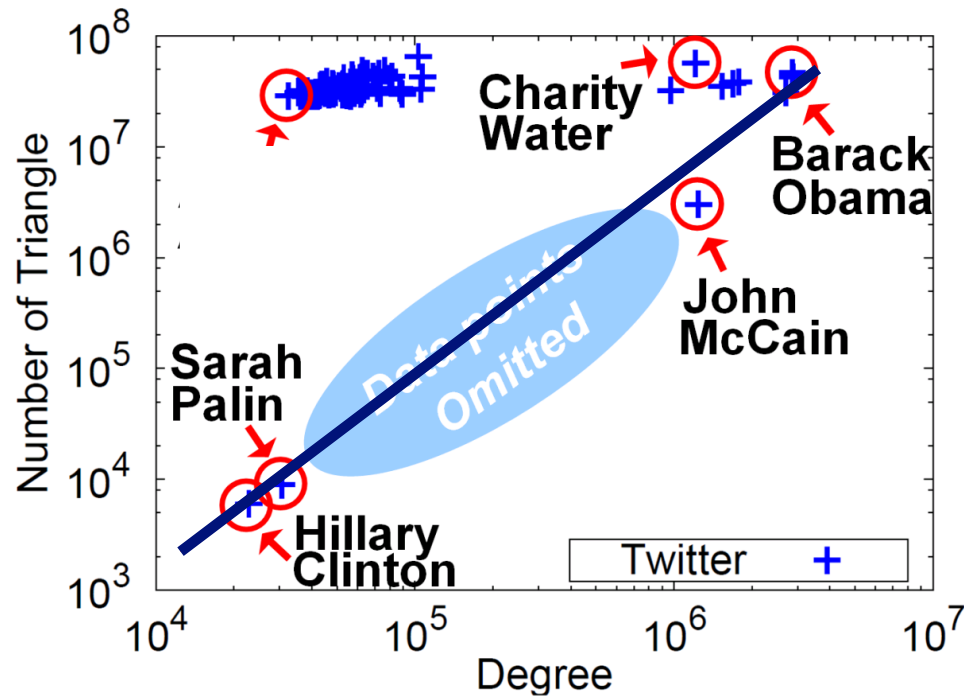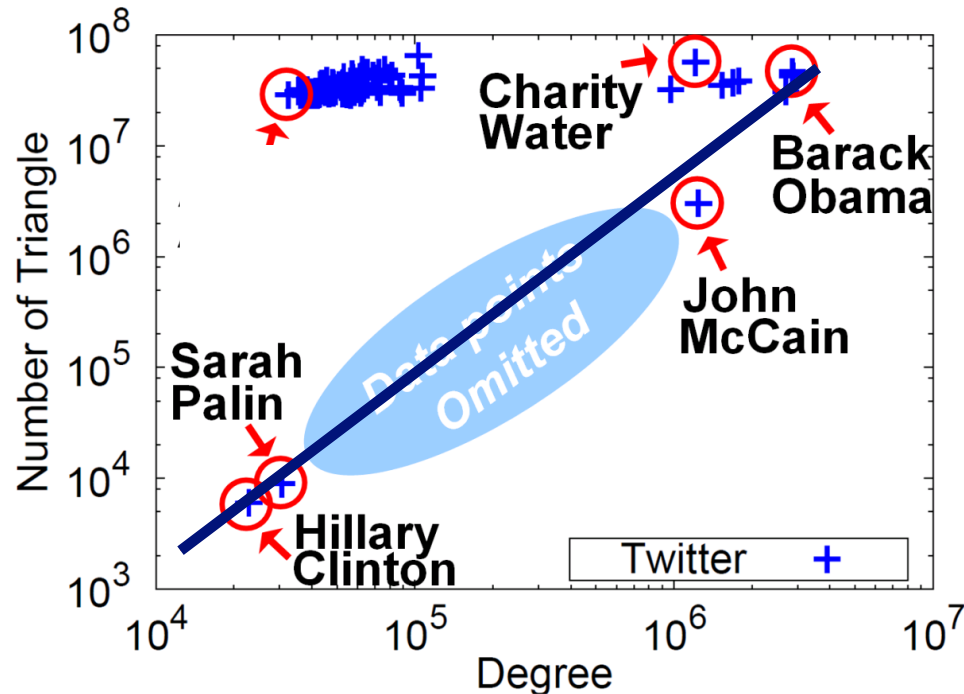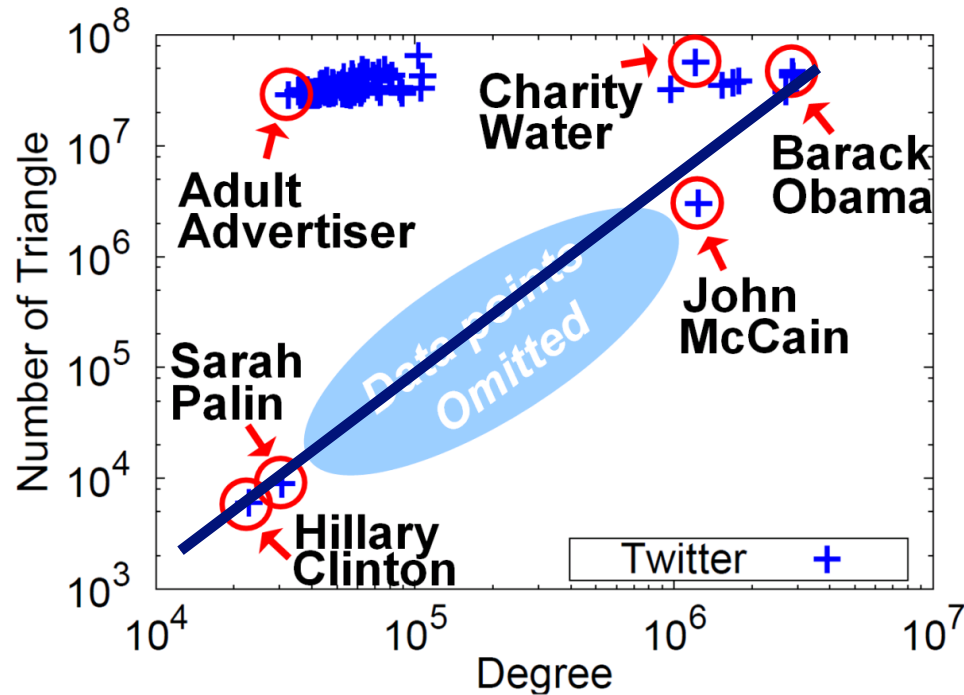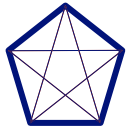[U Kang, Brendan Meeder, +, PAKDD'11]

# Triangle counting for large graphs?



Anomalous nodes in Twitter(~ 3 billion edges)

[U Kang, Brendan Meeder, +, PAKDD'11]

# Triangle counting for large graphs?



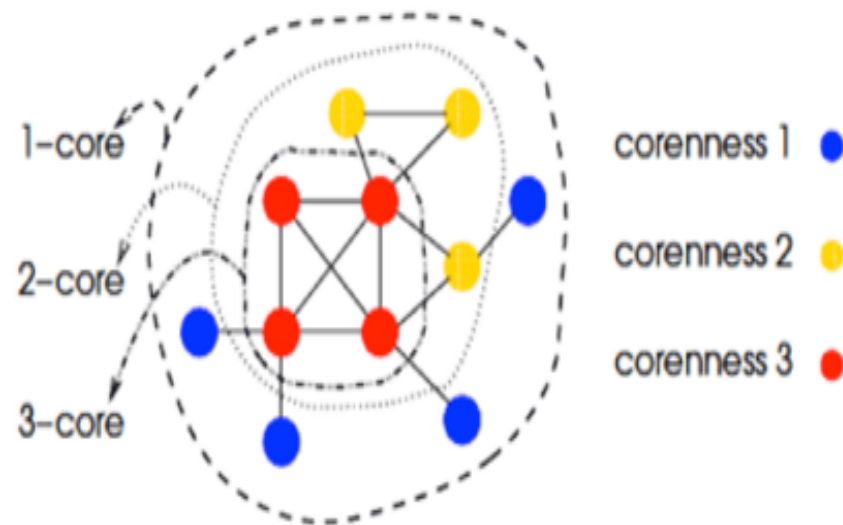Anomalous nodes in Twitter(~ 3 billion edges)

[U Kang, Brendan Meeder, +, PAKDD'11]

# S4: k-core patterns - dfn

- **k-core** (of a graph)
- **degeneracy** (of a graph)
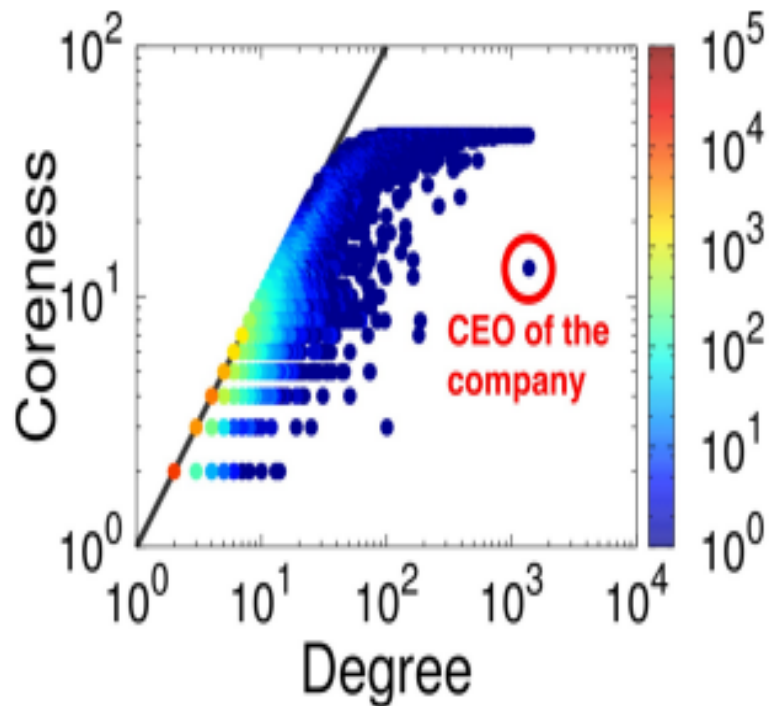- **coreness** (of a vertex)

# Mirror Pattern: Observation

- **coreness** (of a vertex): maximum k such that the vertex belongs to the k-core
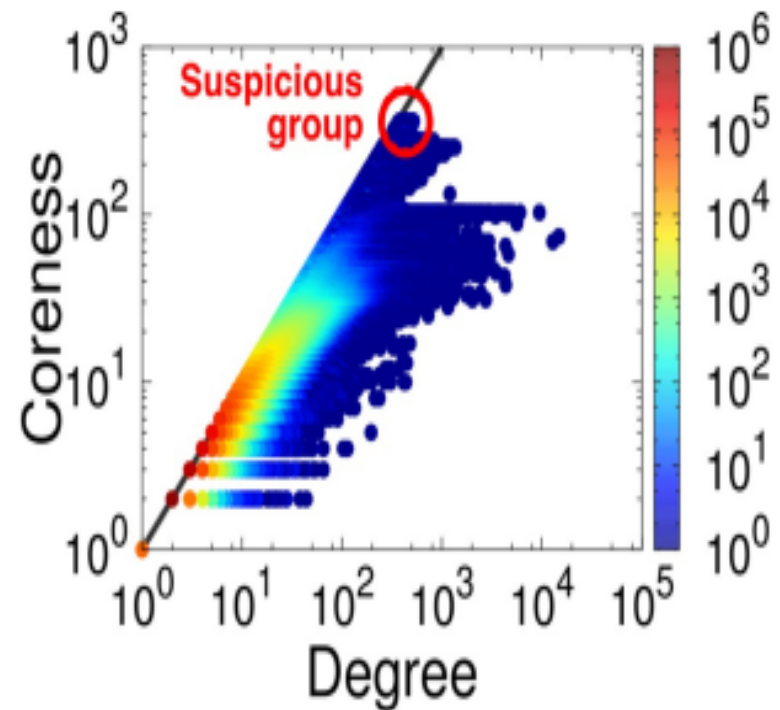
- Definition: **[Mirror Pattern]** *degree ~ coreness*



Rank correlation $\rho = 0.99$

# Mirror Pattern: Application

- Exceptions are 'strange'



Email ($\rho = 0.98$)

LiveJournal ($\rho = 0.99$)

# MORE Graph Patterns

| | Unweighted | Weighted |
|---|---|---|
| **Static** | ✔ L01. Power-law degree distribution [Faloutsos et al. `99, Kleinberg et al. `99, Chakrabarti et al. `04, Newman `04]<br>✔ L02. Triangle Power Law (TPL) [Tsourakakis `08]<br>✔ L03. Eigenvalue Power Law (EPL) [Siganos et al. `03]<br>L04. Community structure [Flake et al. `02, Girvan and Newman `02] | L10. Snapshot Power Law (SPL) [McGlohon et al. `08] |
| **Dynamic** | L05. Densification Power Law (DPL) [Leskovec et al. `05]<br>L06. Small and shrinking diameter [Albert and Barabási `99, Leskovec et al. `05]<br>L07. Constant size 2nd and 3rd connected components [McGlohon et al. `08]<br>L08. Principal Eigenvalue Power Law ($\lambda_1$PL) [Akoglu et al. `08]<br>L09. Bursty/self-similar edge/weight additions [Gomez and Santonja `98, Gribble et al. `98, Crovella and | L11. Weight Power Law (WPL) [McGlohon et al. `08] |

*RTG: A Recursive Realistic Graph Generator using Random Typing* Leman Akoglu and Christos Faloutsos. *PKDD*'09.
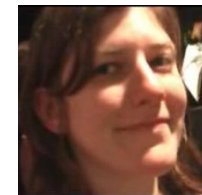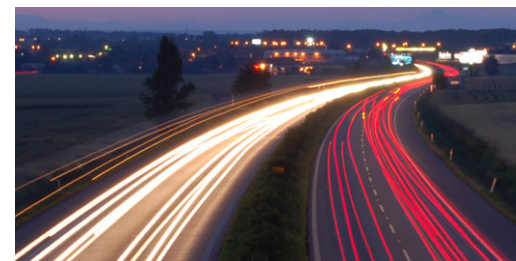
# MORE Graph Patterns

| | Unweighted | Weighted |
|---|---|---|
| Static | L01. Power-law degree distribution [Faloutsos et al. `99, Kleinberg et al. `99, Chakrabarti et al. `04, Newman `04] <br> L02. Triangle Power Law (TPL) [Tsourakakis `08] <br> L03. Eigenvalue Power Law (EPL) [Siganos et al. `03] <br> L04. Community structure [Flake et al. `02, Girvan and Newman `02] | L10. Snapshot Power Law (SPL) [McGlohon et al. `08] |
| Dynamic | L05. Densification Power Law (DPL) [Leskovec et al. `05] <br> L06. Small and shrinking diameter [Albert and Barabási `99, Leskovec et al. `05] <br> L07. Constant size 2nd and 3rd connected components [McGlohon et al. `08] <br> L08. Principal Eigenvalue Power Law ($\lambda_1$PL) [Akoglu et al. `08] <br> L09. Bursty/self-similar edge/weight additions [Gomez and Santonja `98, Gribble et al. `98, Crovella and Bestavros `99, McGlohon et al. `08] | L11. Weight Power Law (WPL) [McGlohon et al. `08] |

• Mary McGlohon, Leman Akoglu, Christos Faloutsos. *Statistical Properties of Social Networks.* in "Social Network Data Analytics" (Ed.: Charu Aggarwal)



• Deepayan Chakrabarti and Christos Faloutsos, *Graph Mining: Laws, Tools, and Case Studies* Oct. 2012, Morgan Claypool.

# Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
  - P1.1: Patterns
  - P1.2: Anomaly / fraud detection
    - No labels – spectral
    - With labels: Belief Propagation
  
  Patterns  anomalies
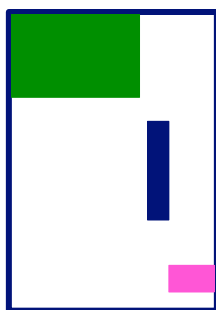- Part#2: time-evolving graphs; tensors
- Conclusions

# How to find 'suspicious' groups?

- 'blocks' are normal, right?

idols

fans

# Except that:



- 'blocks' are normal, igh?
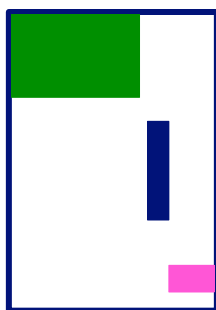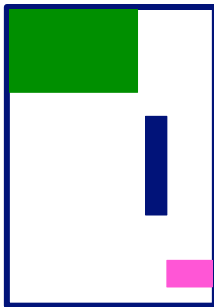- 'hyperbolic' communities are more realistic [Araujo+, PKDD'14]

# Except that:

- 'blocks' are usually suspicious
- 'hyperbolic' communities are more realistic [Araujo+, PKDD'14]

Q: Can we spot blocks, easily?

# Except that:



- 'blocks' are usually suspicious
- 'hyperbolic' communities are more realistic [Araujo+, PKDD'14]
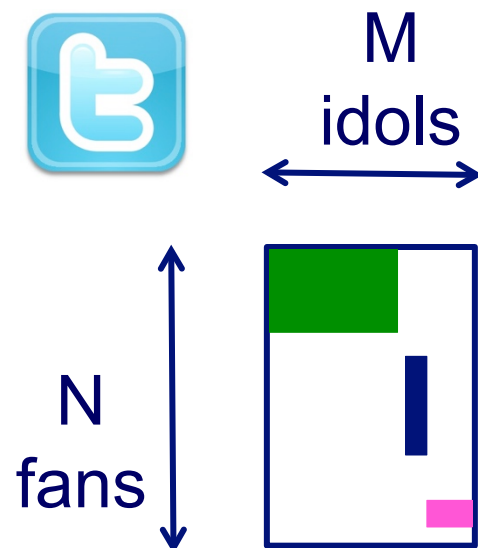
Q: Can we spot blocks, easily?
A: Silver bullet: SVD!

# Crush intro to SVD

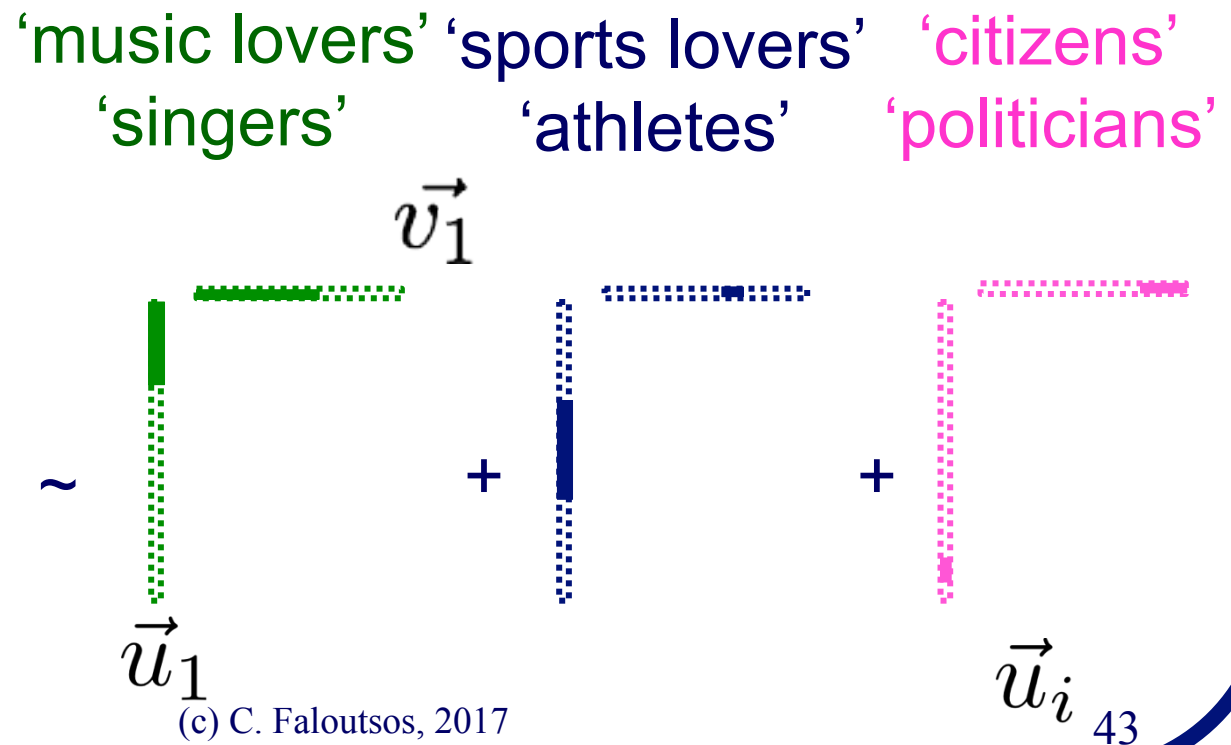- Recall: (SVD) matrix factorization: finds blocks



'music lovers' 'sports lovers' 'citizens'
'singers' 'athletes' 'politicians'

M idols

N fans

$\vec{v_1}$

$\sim$ $\quad$ + $\quad$ +

$\vec{u_1}$ $\qquad$ $\vec{u_i}$

Toutiao/Byte-Dance $\qquad$ (c) C. Faloutsos, 2017 $\qquad$ 42

# Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks

M idols

N fans

'music lovers' 'sports lovers' 'citizens'
'singers' 'athletes' 'politicians'

$$\vec{v}_1$$

$$\sim \qquad + \qquad +$$

$$\vec{u}_1 \qquad\qquad\qquad \vec{u}_i$$

Toutiao/Byte-Dance

(c) C. Faloutsos, 2017

43

# Inferring Strange Behavior from Connectivity Pattern in Social Networks
## PAKDD'14

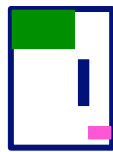Meng Jiang, Peng Cui, Shiqiang Yang (Tsinghua)
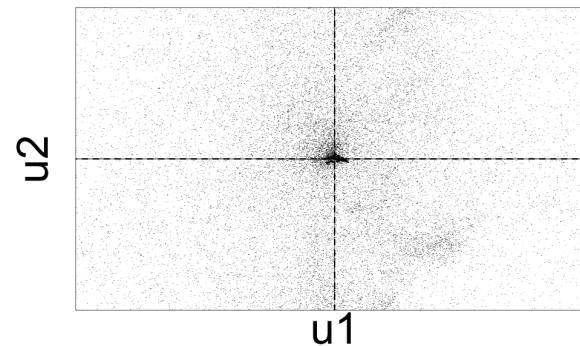Alex Beutel, Christos Faloutsos (CMU)

# *Lockstep* and *Spectral Subspace Plot*

- Case #0: No lockstep behavior in random power law graph of 1M nodes, 3M edges

- Random $\longrightarrow$ "Scatter"
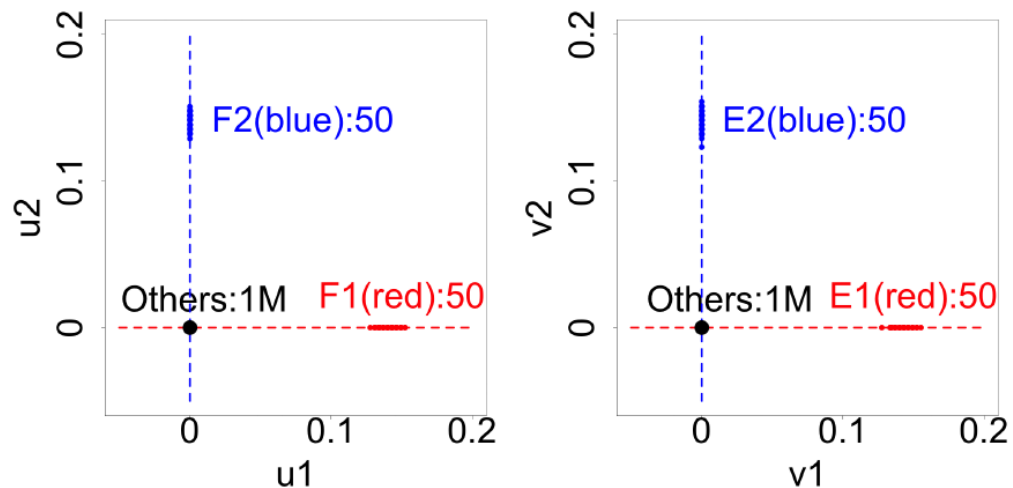
Adjacency Matrix

Spectral Subspace Plot

# *Lockstep* and *Spectral Subspace Plot*

- Case #1: non-overlapping lockstep
- "Blocks" $\longleftrightarrow$ "Rays"

Adjacency Matrix                    Spectral Subspace Plot



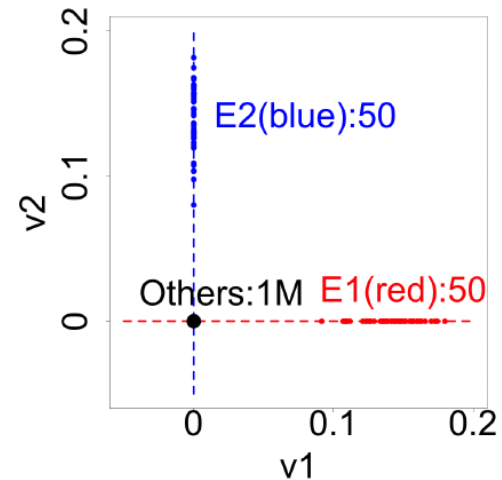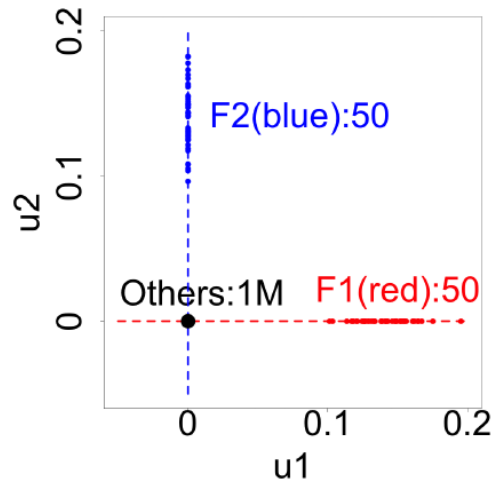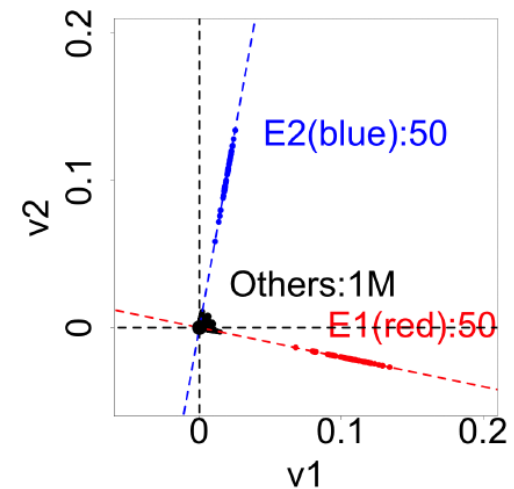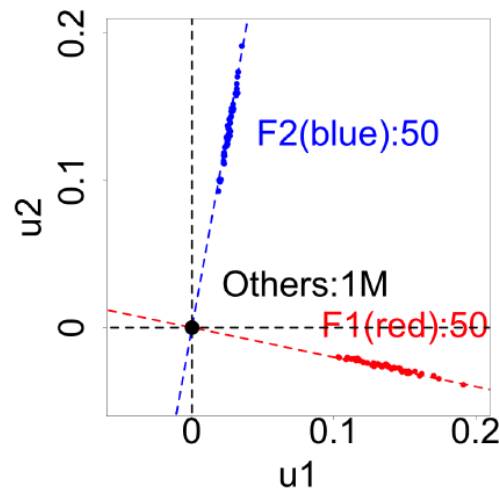Rule 1 (short "rays"): two blocks, high density (90%), no "camouflage", no "fame"

# *Lockstep* and *Spectral Subspace Plot*

- Case #2: non-overlapping lockstep
- "Blocks; low density" ⟵⟶ Elongation
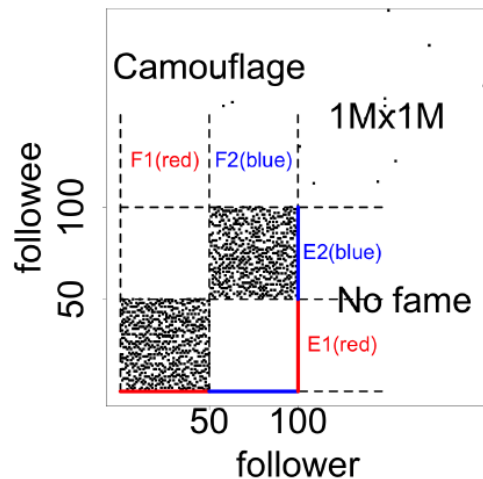
Adjacency Matrix          Spectral Subspace Plot



Rule 2 (long "rays"): two blocks, low density (50%), no "camouflage", no "fame"

# *Lockstep* and *Spectral Subspace Plot*

- Case #3: non-overlapping lockstep
- "**Camouflage**" (or "Fame") ⟷ Tilting "Rays"

### Adjacency Matrix

### Spectral Subspace Plot



Rule 3 (tilting "rays"): two blocks, with "camouflage", no "fame"

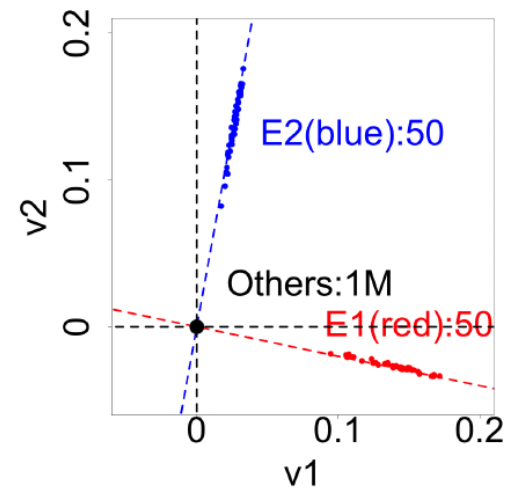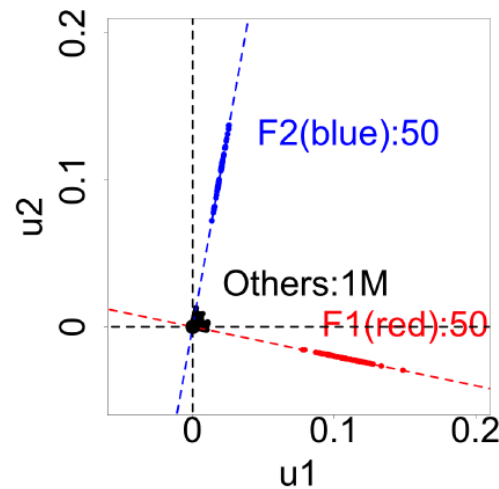# *Lockstep* and *Spectral Subspace Plot*

- Case #3: non-overlapping lockstep
- "Camouflage" (or "**Fame**") ⟷ Tilting "Rays"

Adjacency Matrix                    Spectral Subspace Plot



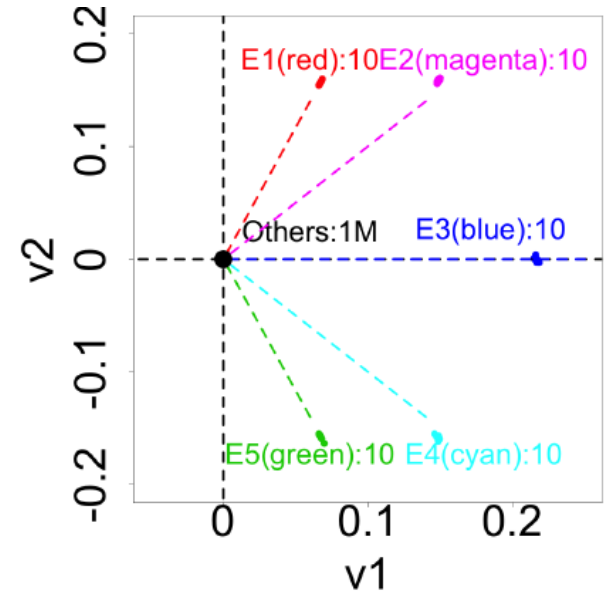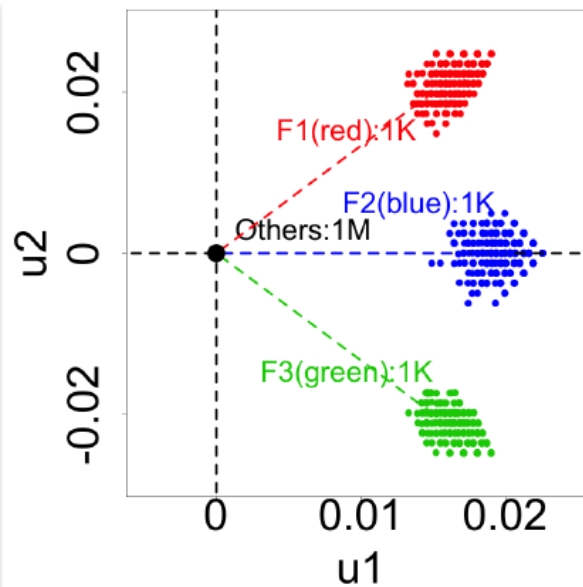Rule 3 (tilting "rays"): two blocks, no "camouflage", with "fame"

# *Lockstep* and *Spectral Subspace Plot*

- Case #4:      ?      lockstep
- "?"                              "Pearls"

⟷

**Adjacency Matrix**          **Spectral Subspace Plot**
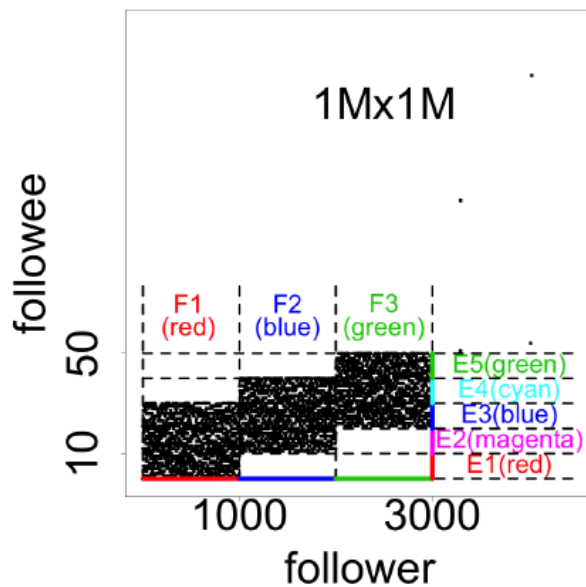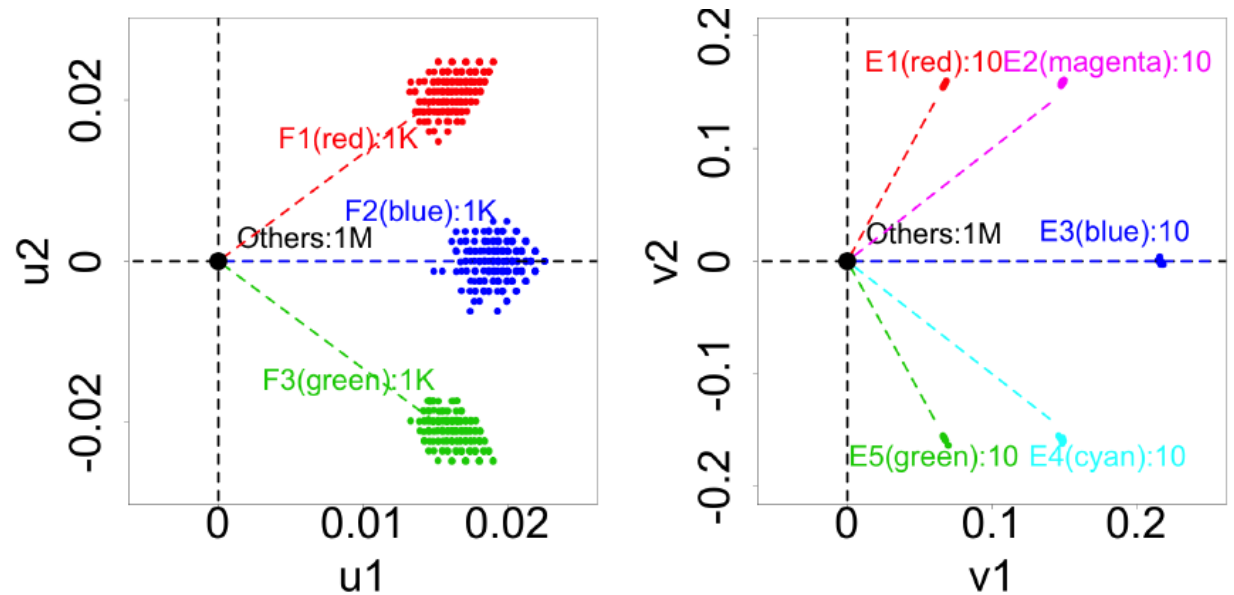
# *Lockstep* and *Spectral Subspace Plot*

- Case #4: **overlapping** lockstep
- "**Staircase**" ⟷ "Pearls"

Adjacency Matrix          Spectral Subspace Plot

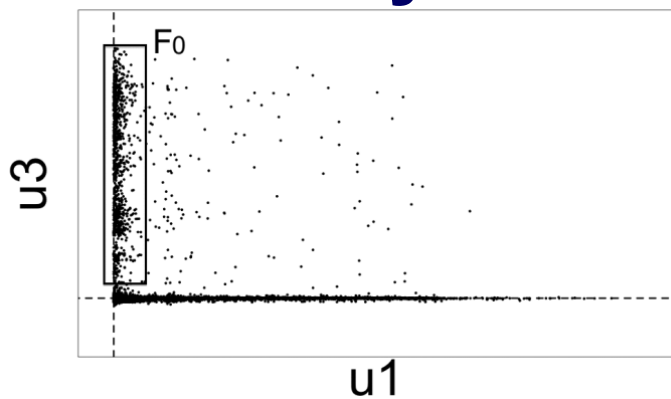Rule 4 ("pearls"): a "staircase" of three partially overlapping blocks.

# Dataset

- Tencent Weibo 

- 117 million nodes (with profile and UGC data)
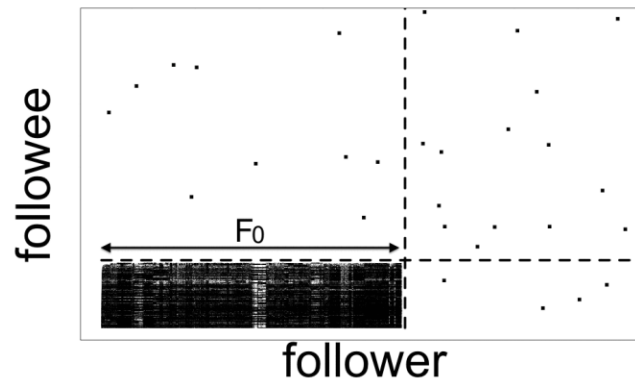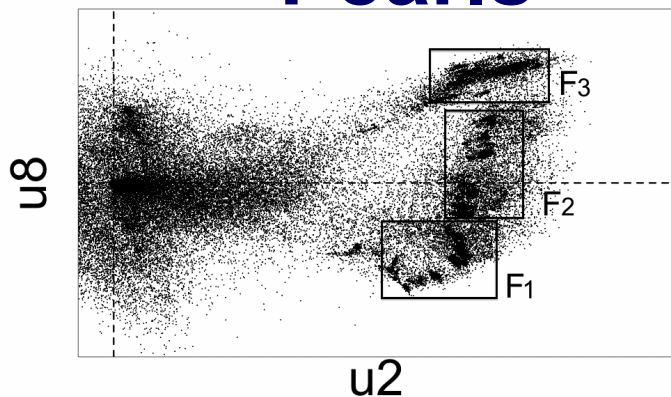
- 3.33 billion directed edges
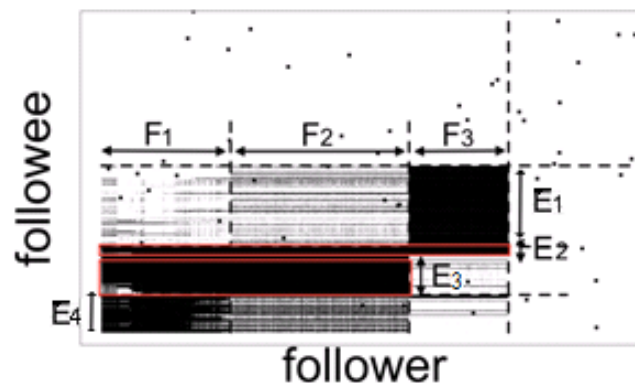
# Real Data

## "Rays"



## "Block"



## "Pearls"



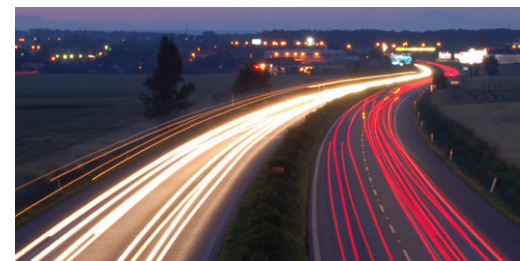## "Staircase"

# Real Data

- Spikes on the out-degree distribution
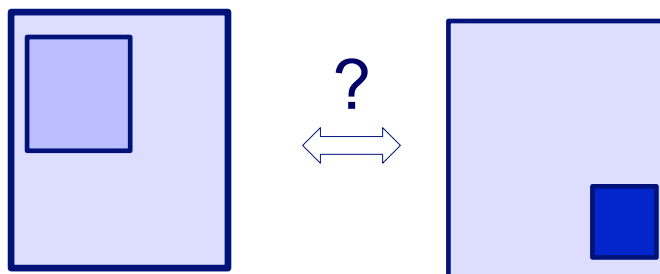
# **Roadmap**



- Introduction – Motivation
- Part#1: Patterns in graphs
  - P1.1: Patterns
  - P1.2: Anomaly / fraud detection
    - No labels – spectral methods
      - **Suspiciousness**
    - With labels: Belief Propagation
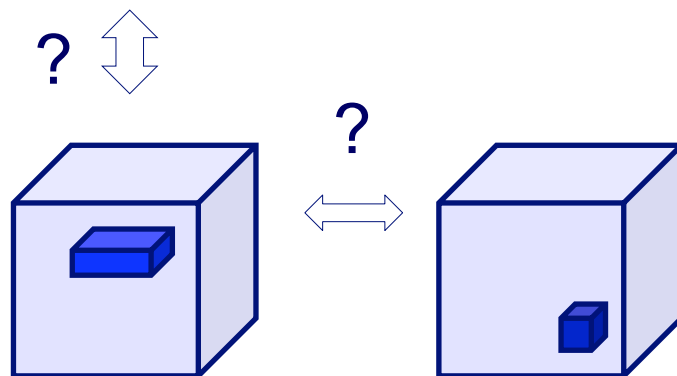- Part#2: time-evolving graphs; tensors
- Conclusions

# Suspicious Patterns in Event Data

2-modes

$n$-modes

?

?

?

**A General Suspiciousness Metric for Dense Blocks in Multimodal Data,** Meng Jiang, Alex Beutel, Peng Cui, Bryan Hooi, Shiqiang Yang, and Christos Faloutsos, *ICDM*, 2015.
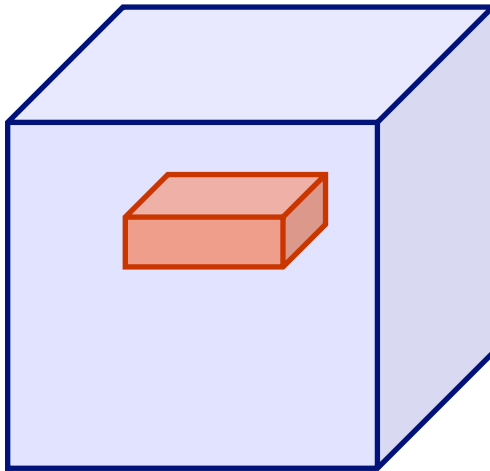
# Suspicious Patterns in Event Data

## Which is more suspicious?

20,000 Users ←→ 225 Users

Retweeting same 20 tweets ←→ Retweeting same 1 tweet

6 times each      vs.      15 times each

All in 10 hours ←→ All in 3 hours

All from 2 IP addresses

# Suspicious Patterns in Event Data

## Which is more suspicious?
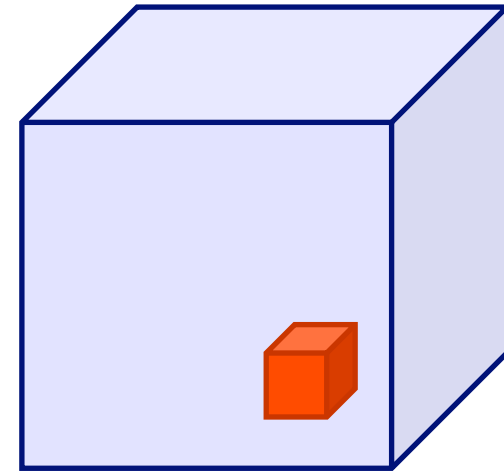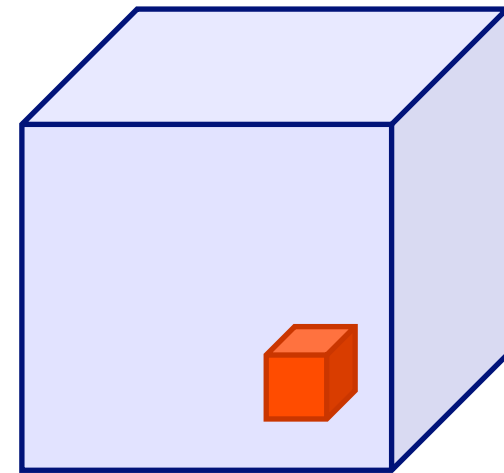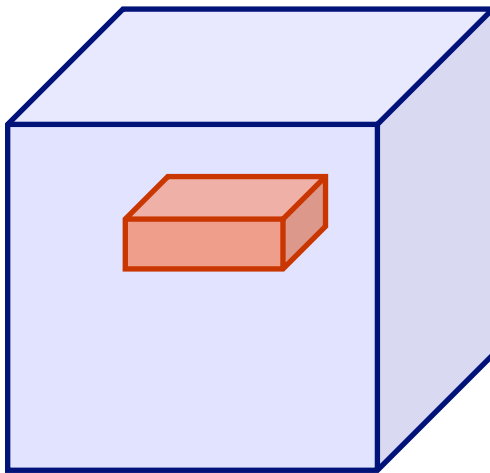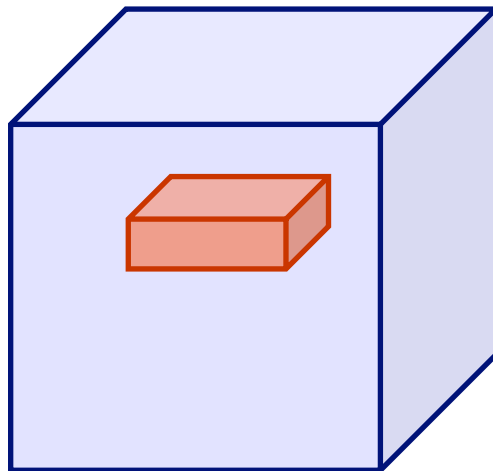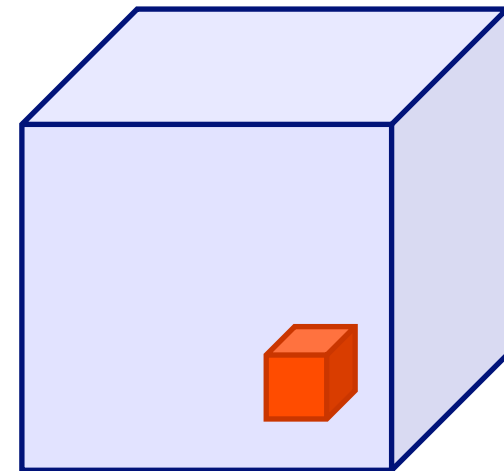
20,000 Users
Retweeting same 20 tweets
  6 times each
All in 10 hours

vs.

225 Users
Retweeting same 1 tweet
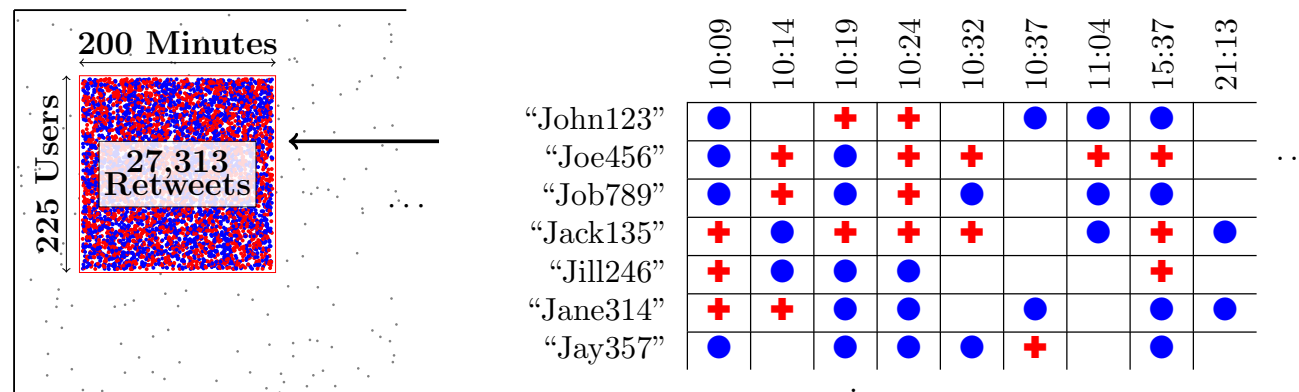  15 times each
All in 3 hours
All from 2 IP addresses

Answer: volume $* D_{KL}(p \| p_{background})$

58

# Suspicious Patterns in Event Data

## Which is more suspicious?

| 20,000 Users | | 225 Users |
| Retweeting same 20 tweets | | Retweeting same 1 tweet |
| 6 times each | vs. | 15 times each |
| All in 10 hours | | All in 3 hours |
| | | All from 2 IP addresses |

size        contrast

**Answer: volume * $D_{KL}(p \| p_{background})$**
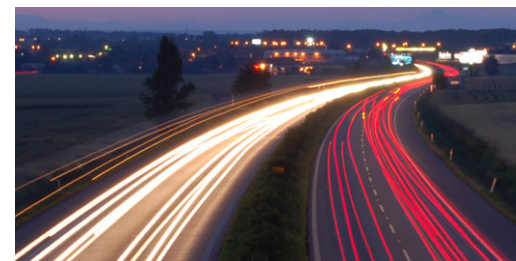
59

# Suspicious Patterns in Event Data



Retweeting: "Galaxy Note Dream Project: Happy Happy Life Traveling the World"
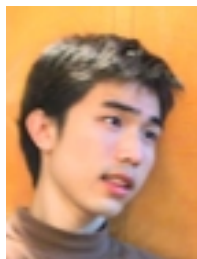
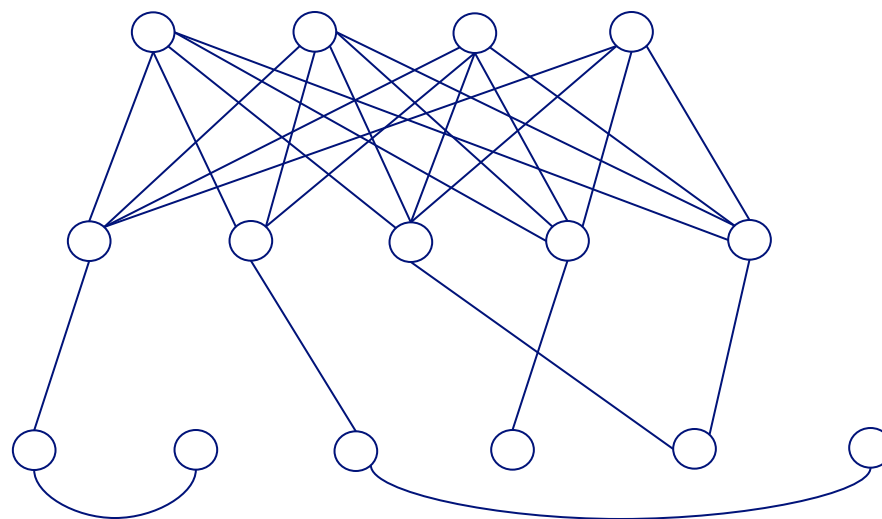|  | # | User $\times$ tweet $\times$ IP $\times$ minute | Mass $c$ | Suspiciousness |
|---|---|---|---|---|
| **CROSSSPOT** | 1 | $14 \times 1 \times 2 \times 1,114$ | 41,396 | 1,239,865 |
|  | 2 | $225 \times 1 \times 2 \times 200$ | 27,313 | 777,781 |
|  | 3 | $8 \times 2 \times 4 \times 1,872$ | 17,701 | 491,323 |
| HOSVD | 1 | $24 \times 6 \times 11 \times 439$ | 3,582 | 131,113 |
|  | 2 | $18 \times 4 \times 5 \times 223$ | 1,942 | 74,087 |
|  | 3 | $14 \times 2 \times 1 \times 265$ | 9,061 | 381,211 |

# Roadmap



- Introduction – Motivation
- Part#1: Patterns in graphs
  - P1.1: Patterns
  - P1.2: Anomaly / fraud detection
    - No labels – spectral methods
    - With labels: Belief Propagation
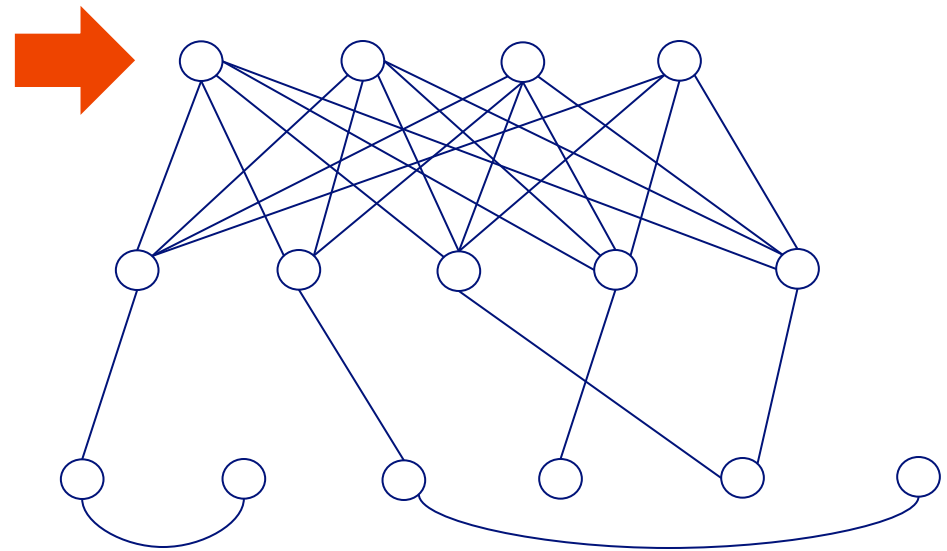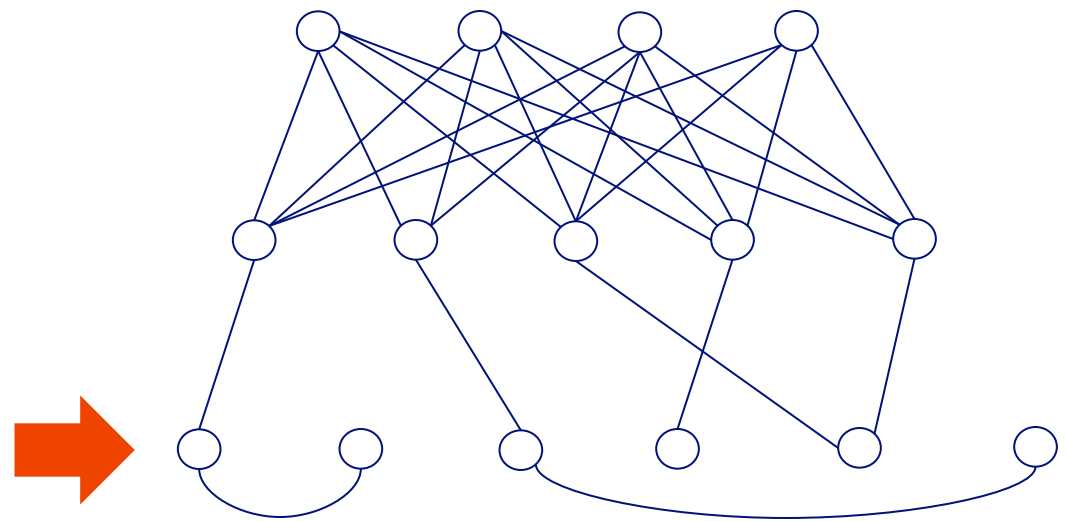- Part#2: time-evolving graphs; tensors
- Conclusions

# E-bay Fraud detection

w/ Polo Chau &
Shashank Pandit, CMU
[www'07]
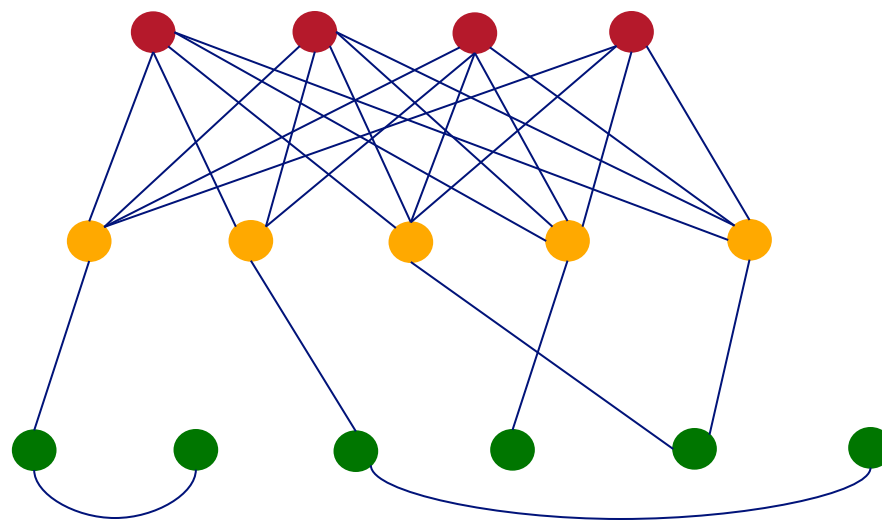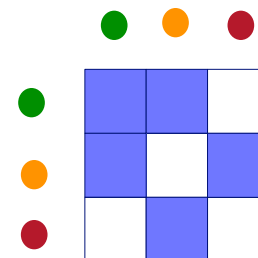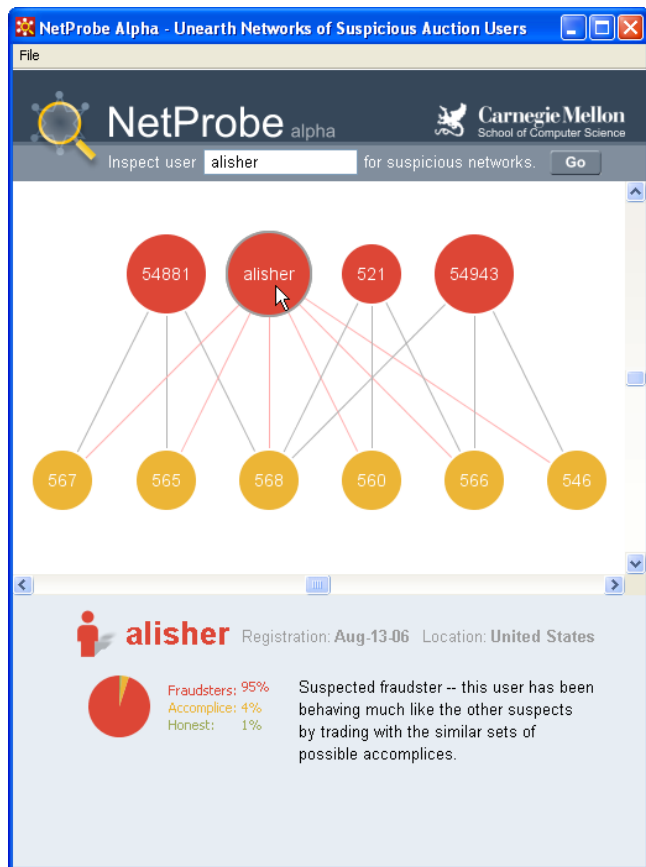
# E-bay Fraud detection

# E-bay Fraud detection

# E-bay Fraud detection - NetProbe

# Popular press

And less desirable attention:
- E-mail from 'Belgium police' ('copy of your code?')

# **Roadmap**



- Introduction – Motivation
- Part#1: Patterns in graphs
  - Patterns
  - Anomaly / fraud detection
    - No labels - Spectral methods
    - w/ labels: Belief Propagation – closed formulas
- Part#2: time-evolving graphs; tensors
- Conclusions

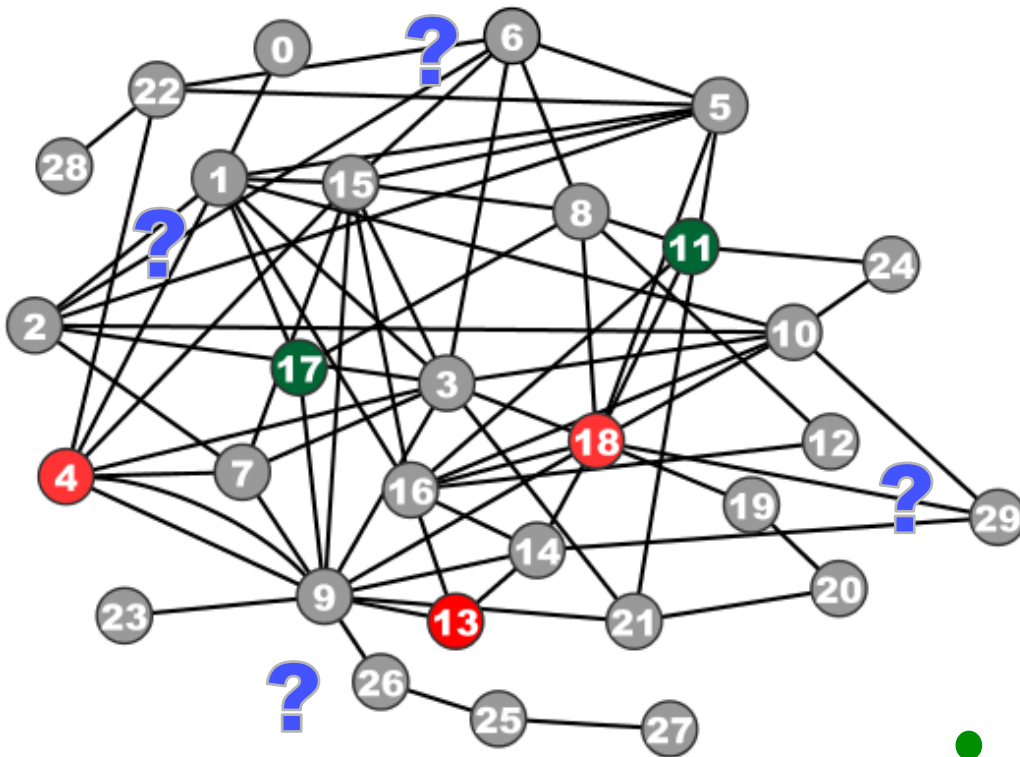# Unifying Guilt-by-Association Approaches: Theorems and Fast Algorithms
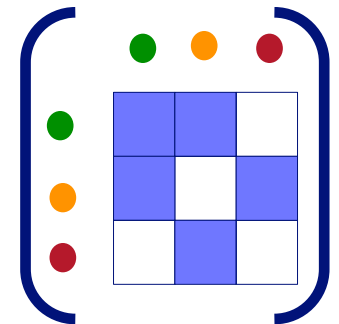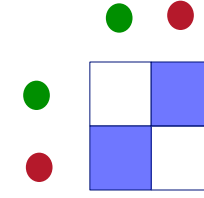
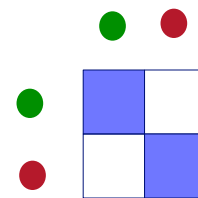**Danai Koutra**
U Kang
Hsing-Kuo Kenneth Pao

Tai-You Ke
Duen Horng (Polo) Chau
Christos Faloutsos

# Problem Definition: GBA techniques



**Given**: Graph; &
  few labeled nodes
**Find**: labels of rest
(assuming network
effects)

# Are they related?

- RWR (Random Walk with Restarts)
  - google's pageRank ('*if my friends are important, I'm important, too*')  Google
- SSL (Semi-supervised learning)
  - minimize the differences among neighbors
- BP (Belief propagation)
  - send messages to neighbors, on what you believe about them
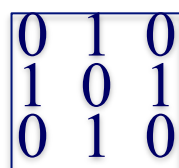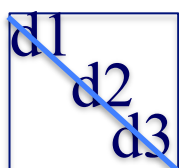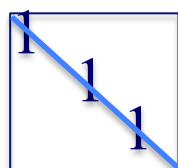
# Are they related? YES!

- RWR (Random Walk with Restarts)
  - google's pageRank ('*if my friends are important, I'm important, too*')
- SSL (Semi-supervised learning)
  - minimize the differences among neighbors
- BP (Belief propagation)
  - send messages to neighbors, on what you believe about them

# Correspondence of Methods

| Method | Matrix | | Unknown | | known |
|---|---|---|---|---|---|
| **RWR** | $[\mathbf{I} - c\ \underline{\mathbf{A}}\mathbf{D}^{-1}]$ | $\times$ | $\mathbf{x}$ | $=$ | $(1-c)\mathbf{y}$ |
| **SSL** | $[\mathbf{I} + a(\mathbf{D} - \underline{\mathbf{A}})]$ | $\times$ | $\mathbf{x}$ | $=$ | $\mathbf{y}$ |
| **FABP** | $[\mathbf{I} + a\,\mathbf{D} - c'\underline{\mathbf{A}}]$ | $\times$ | $\mathbf{b_h}$ | $=$ | $\boldsymbol{\phi_h}$ |

$$
\begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix}
\quad
\begin{bmatrix} d1 & & \\ & d2 & \\ & & d3 \end{bmatrix}
\quad
\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}
\qquad
\begin{bmatrix} ? \end{bmatrix}
\qquad
\begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}
$$

adjacency matrix     final labels/ beliefs     prior labels/ beliefs
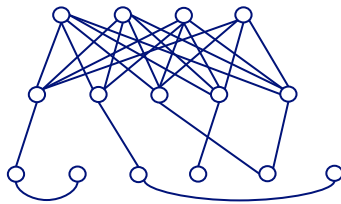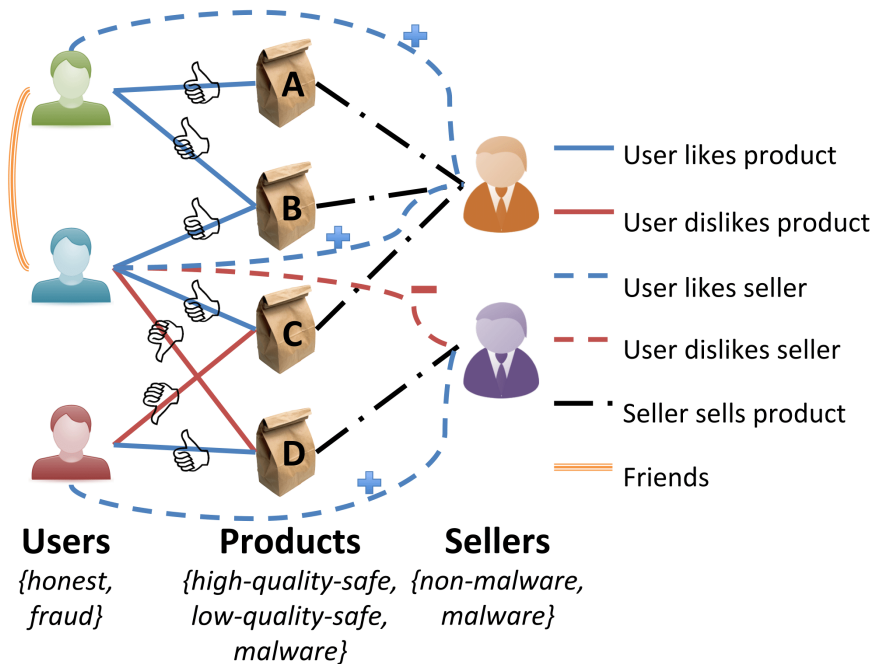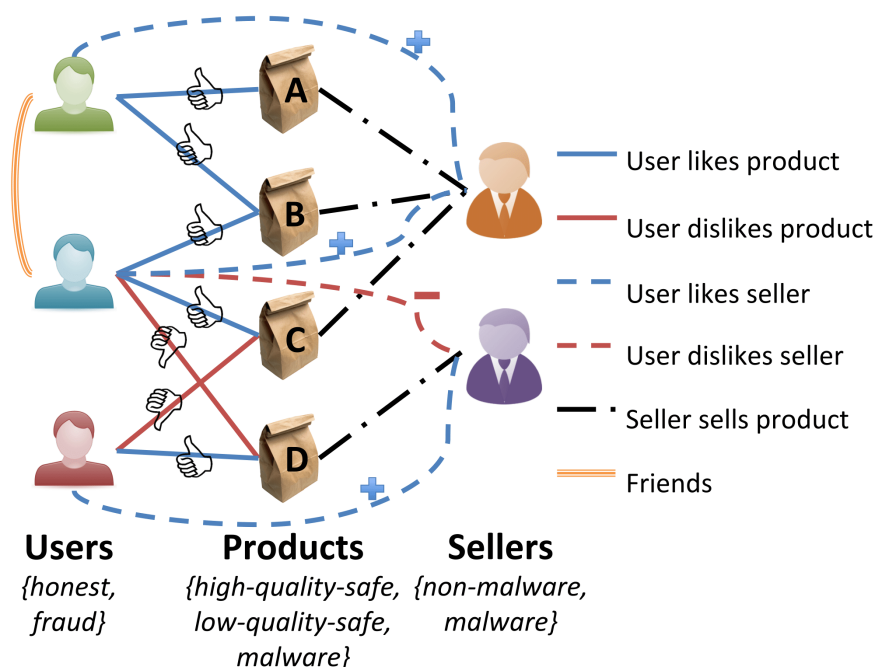
# Results: Scalability



FABP is **linear** on the number of edges.

# Problem: e-commerce ratings fraud



| | |
|---|---|
| ——— | User likes product |
| ——— | User dislikes product |
| – – – | User likes seller |
| – – – | User dislikes seller |
| – · – · | Seller sells product |
| ═══ | Friends |

**Users**
*{honest, fraud}*

**Products**
*{high-quality-safe, low-quality-safe, malware}*

**Sellers**
*{non-malware, malware}*

- **Given** a **heterogeneous** graph on users, products, sellers and positive/negative ratings with "seed labels"

- **Find** the top *k* most fraudulent users, products and sellers

# Problem: e-commerce ratings fraud



User likes product
User dislikes product
User likes seller
User dislikes seller
Seller sells product
Friends

**Users**
*{honest, fraud}*

**Products**
*{high-quality-safe, low-quality-safe, malware}*

**Sellers**
*{non-malware, malware}*

- **Given** a heterogeneous graph on users, products, sellers and positive/negative ratings with "seed labels"

- **Find** the top $k$ most fraudulent users, products and sellers

# Problem: e-commerce ratings fraud



**Theorem 1** (ZooBP). *If* $\mathbf{b}, \mathbf{e}, \mathbf{P}, \mathbf{Q}$ *are constructed as described above, the linear equation system approximating the final node beliefs given by BP is:*

$$\mathbf{b} = \mathbf{e} + (\mathbf{P} - \mathbf{Q})\mathbf{b} \qquad (\text{ZooBP}) \qquad (10)$$

# ZooBP: features

Fast; convergence guarantees.



**Near-perfect accuracy**　　　**linear in graph size**

# ZooBP in the real world



- Near 100% precision on top 300 users (Flipkart)
- Flagged users: suspicious
  - 400 ratings in 1 sec
  - 5000 good ratings and no bad ratings

# Summary of Part#1

- *many* patterns in real graphs
  - Power-laws everywhere
  - Long (and growing) list of tools for anomaly/ fraud detection



Patterns  anomalies

# Roadmap



- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs
  - P2.1: tools/tensors
  - P2.2: other patterns
- Conclusions

# Part 2: Time evolving graphs; tensors

# Graphs over time -> tensors!

- Problem #2.1:
  - Given who calls whom, and when
  - Find patterns / anomalies

*johnson*

smith

# Graphs over time -> tensors!

- Problem #2.1:
  - Given who calls whom, and when
  - Find patterns / anomalies

# Graphs over time -> tensors!

- Problem #2.1:
  - Given who calls whom, and when
  - Find patterns / anomalies

Tue

Mon

# Graphs over time -> tensors!

- Problem #2.1:
  - Given who calls whom, and when
  - Find patterns / anomalies



time

caller

callee

# Graphs over time -> tensors!

- Problem #2.1':
  - Given author-keyword-date
  - Find patterns / anomalies

date

author

keyword

MANY more settings, with >2 'modes'

# **Graphs over time -> tensors!**

- Problem #2.1'':
  – Given subject – verb – object facts
  – Find patterns / anomalies

verb

subject

object

MANY more settings,
with >2 'modes'

# Graphs over time -> tensors!

- Problem #2.1''':
  - Given <triplets>
  - Find patterns / anomalies

mode3

mode1

mode2

MANY more settings, with >2 'modes' (and 4, 5, etc modes)

# Answer : tensor factorization

- Recall: (SVD) matrix factorization: finds blocks



M products

N users

'meat-eaters' 'steaks'  'vegetarians' 'plants'  'kids' 'cookies'

$\vec{v}_1$

$\sim$ $\vec{u}_1$ $+$ $+$ $\vec{u}_i$

Toutiao/Byte-Dance

(c) C. Faloutsos, 2017

89

# Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks

'music lovers' 'sports lovers' 'citizens'
'singers' 'athletes' 'politicians'

M idols

N fans

$\vec{v_1}$

$\vec{u_1}$ $\vec{u_i}$

~ + +

Toutiao/Byte-Dance

(c) C. Faloutsos, 2017

90

# Answer: tensor factorization

- PARAFAC decomposition

# Answer: tensor factorization

- PARAFAC decomposition
- Results for who-calls-whom-when
  - 4M x 15 days

# Anomaly detection in time-evolving graphs

- Anomalous communities in phone call data:
  – European country, 4M clients, data over 2 weeks

| 1 caller | 5 receivers | 4 days of activity |
|---|---|---|



~200 calls to EACH receiver on EACH day!

# Anomaly detection in time-evolving graphs

- Anomalous communities in phone call data:
  - European country, 4M clients, data over 2 weeks



1 caller      5 receivers      4 days of activity

~200 calls to EACH receiver on EACH day!

# Anomaly detection in time-evolving graphs

- Anomalous communities in phone call data:
  – European country, 4M clients, data over 2 weeks

**Miguel Araujo**, Spiros Papadimitriou, Stephan Günnemann, Christos Faloutsos, Prithwish Basu, Ananthram Swami, Evangelos Papalexakis, Danai Koutra. *Com2: Fast Automatic Discovery of Temporal (Comet) Communities*. PAKDD 2014, Tainan, Taiwan.

# Roadmap



- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs
  - P2.1: tools/tensors
  - → P2.2: other patterns – inter-arrival time
- Conclusions

KDD 2015 – Sydney, Australia

# RSC: Mining and Modeling Temporal Activity in Social Media

Alceu F. Costa*    Yuto Yamaguchi    Agma J. M. Traina

Caetano Traina Jr.    Christos Faloutsos

*alceufc@icmc.usp.br

# Pattern Mining: Datasets

Reddit Dataset

Time-stamp from comments
21,198 users
20 Million time-stamps

Twitter Dataset

Time-stamp from tweets
6,790 users
16 Million time-stamps

For each user we have:

Sequence of postings time-stamps: $T = (t_1, t_2, t_3, ...)$
Inter-arrival times (IAT) of postings: $(\Delta_1, \Delta_2, \Delta_3, ...)$



Toutiao/Byte-Dance          (c) C. Faloutsos, 2017          98

# Pattern Mining

**Pattern 1:** Distribution of IAT is heavy-tailed

Users can be inactive for long periods of time before making new postings

IAT Complementary Cumulative Distribution Function (CCDF)
(log-log axis)

Toutiao/Byte-Dance

(c) C. Faloutsos, 2017

Reddit Users

Twitter Users

# Pattern Mining

**Pattern 1:** Distribution of IAT is heavy-tailed

Users can be inactive for long periods of time before making new postings

IAT Com~~plementary Cumulative Distribution~~ (CCDF)

**No surprises – Should we give up?**



Reddit Users

Twitter Users

# Human? Robots?

linear ↑

log →

# Human? Robots?

# Experiments: Can RSC-Spotter Detect Bots?

## Precision vs. Sensitivity Curves

Good performance: curve close to the top

Twitter

Precision > 94%
Sensitivity > 70%

With strongly imbalanced datasets
# humans >> # bots



Legend: RSC–Spotter — Entropy [6] — All Fe... — IAT Hist. — Weekday Hist.

# Experiments: Can RSC-Spotter Detect Bots?

## Precision vs. Sensitivity Curves

Good performance: curve close to the top

Reddit

Precision > 96%
Sensitivity > 47%
With strongly imbalanced datasets
# humans >> # bots

Precision

Sensitivity (Recall)

RSC–Spotter     Entropy [6]     − − − All Features

IAT Hist.     Weekday Hist.

# Roadmap



- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs
  - P2.1: tools/tensors
  - P2.2: other patterns
    - inter-arrival time
    - Network growth
- Conclusions

# PROBLEM: n(t) and e(t), over time?

- n(t): the number of nodes.
- e(t): the number of edges.
- E.g.:
  - How many members will 🐦 have next month?
  - How many friendship links will 💬 have next year?

- Linear?
- Exponential?
- Sigmoid?



Innovators 2.5 %  Early Adopters 13.5 %  Early Majority 34 %  Late Majority 34 %  Laggards 16 %

C
C/2
0

# Datasets

- **WeChat 2011/1-2013/1   300M nodes, 4.75B links**

- ArXiv      1992/3-2002/3      17k nodes,     2.4M links

- Enron     1998/1-2002/7      86K nodes,     600K links

- Weibo    2006                    165K nodes,   331K links

# A: Power Law Growth



Cumulative growth（Log-Log scale）

# Proposed: NetTide Model

- ## Nodes n(t)
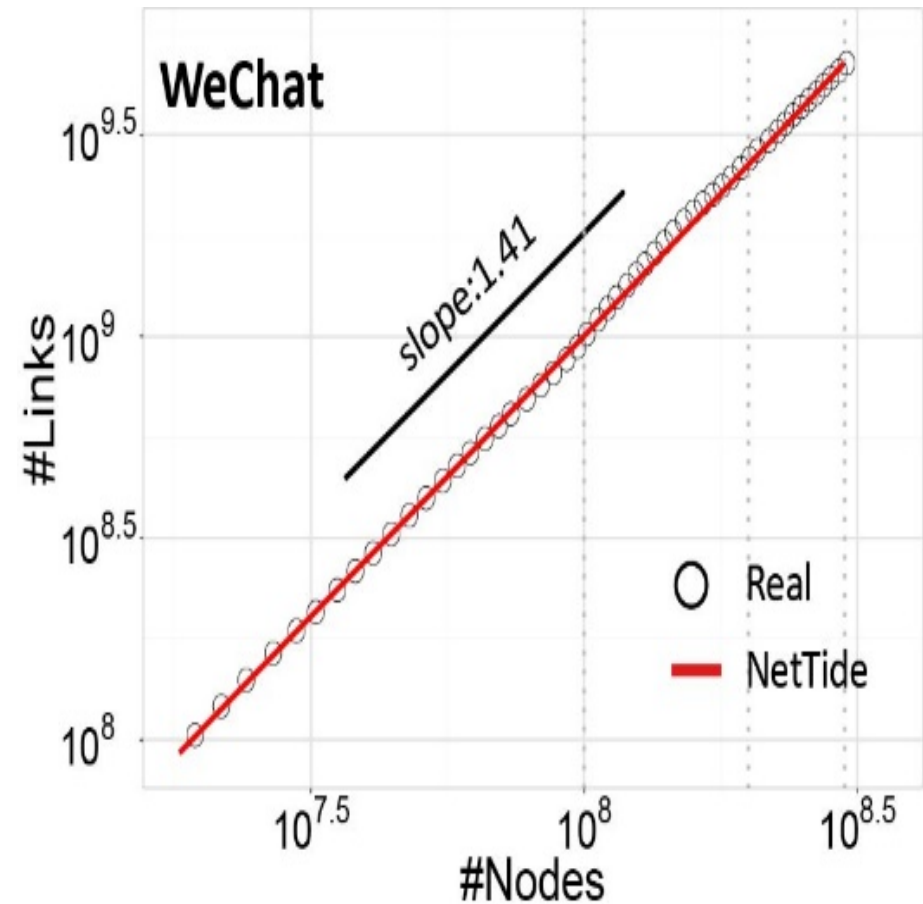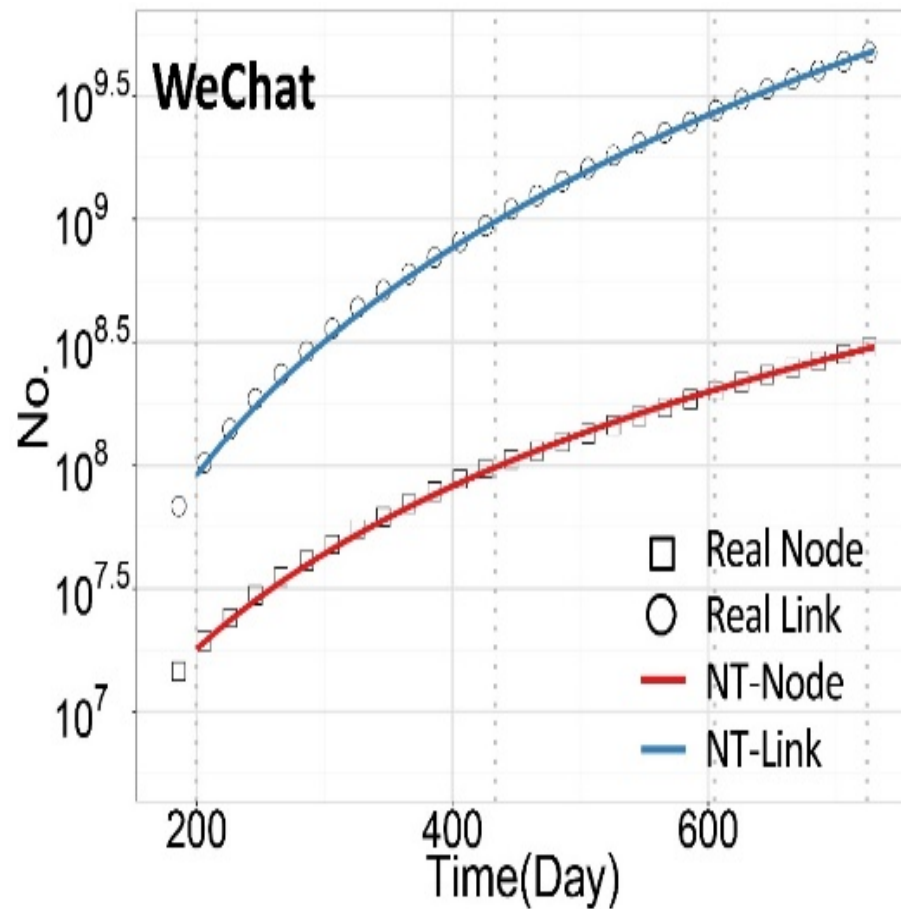
$$\frac{dn(t)}{dt} = \frac{\beta}{t^\theta} n(t)(N - n(t))$$

- ## Links e(t)

$$\frac{de(t)}{dt} = \frac{\beta'}{t^\theta} n(t)\left(\alpha(n(t) - 1)^\gamma - \frac{e(t)}{n(t)}\right) + 2\frac{dn(t)}{dt}$$

# NetTide-Node Model

$$\frac{dn(t)}{dt} = \frac{\beta}{t^\theta} \, n(t) \, (N - n(t))$$

**#nodes(t)**

Total population

- Intuition:
  - **Rich-get-richer**
  - Limitation
  - Fizzling nature

  } = SI; ~Bass

111

# NetTide-Node Model

$$\frac{dn(t)}{dt} = \frac{\beta}{t^{\theta}} \, n(t) \, (N - n(t))$$

#nodes(t)

**Total population**

- Intuition:
  - Rich-get-richer
  - **Limitation**
  - Fizzling nature

} = SI; ~Bass

112

# NetTide-Node Model

$$\frac{dn(t)}{dt} = \frac{\beta}{\boxed{t^{\theta}}} \, n(t) \, (N - n(t))$$

#nodes(t)

Total population

- Intuition:
  - Rich-get-richer
  - Limitation  } = SI; ~Bass
  - **Fizzling nature**

# Results: Accuracy

# Results: Accuracy

# Results: Accuracy

# Results: Accuracy

# Results: Accuracy

# Results: Forecast
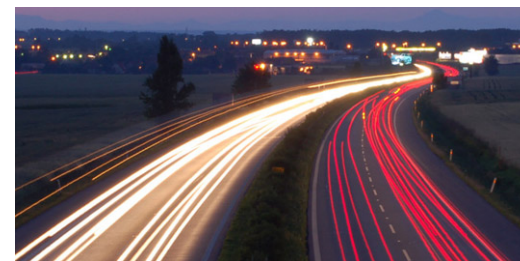
WeChat from 100 million to 300 million

730 days ahead



119

# Part 2: Conclusions

- Time-evolving / heterogeneous graphs -> tensors

- PARAFAC finds patterns

- Surprising temporal patterns (P.L. growth)

# Roadmap



- Introduction – Motivation
  - Why study (big) graphs?
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
➡ - Acknowledgements and Conclusions

# Thanks

# Cast

Akoglu,
Leman

Araujo,
Miguel

Beutel,
Alex

Chau,
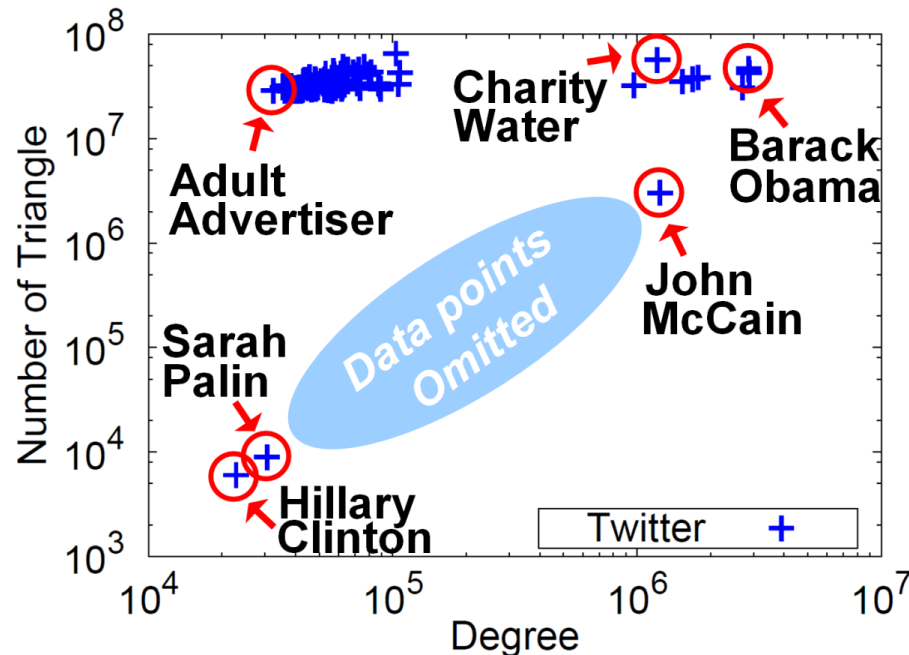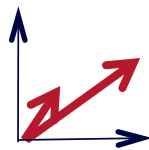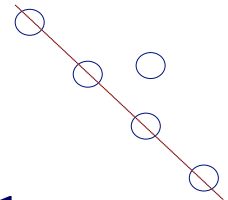Polo

Eswaran,
Dhivya

Hooi,
Bryan

Kang, U

Koutra,
Danai

Papalexakis,
Vagelis

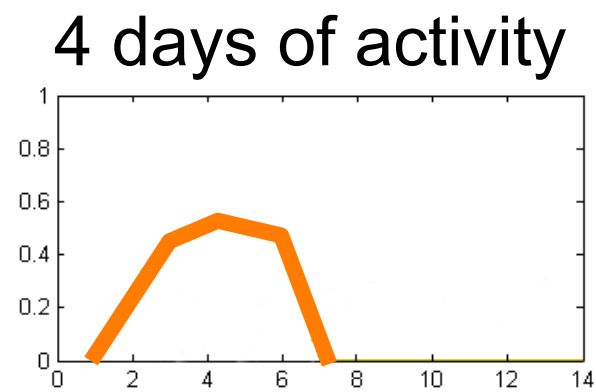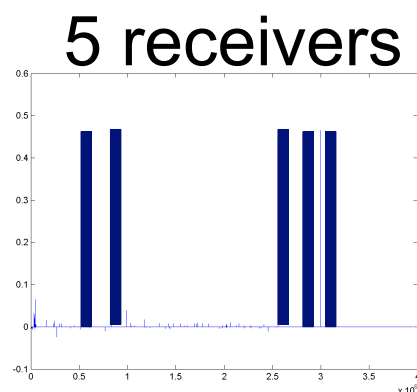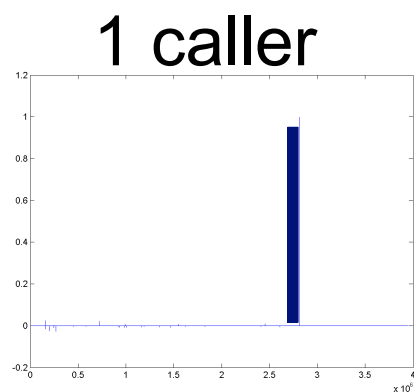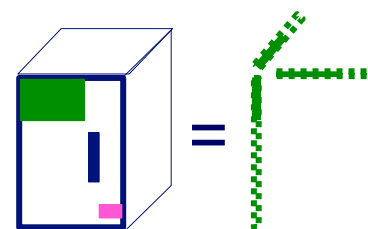Shah,
Neil

Shin,
Kijung

Song,
Hyun Ah

# **CONCLUSION#1 – Big data**

- **Patterns** ⚹ **Anomalies**

- **Large** datasets reveal patterns/outliers that are invisible otherwise

# CONCLUSION#2 – tensors
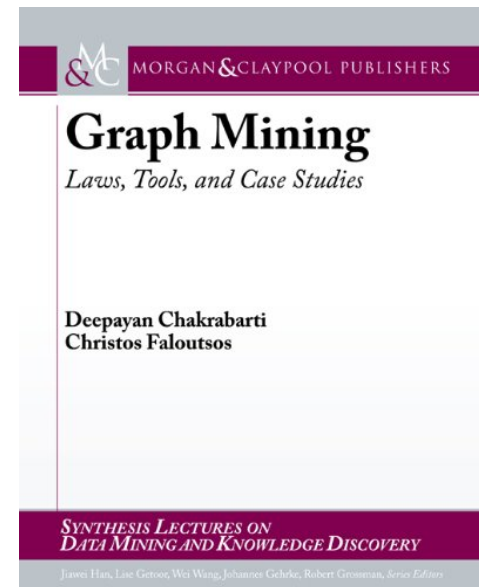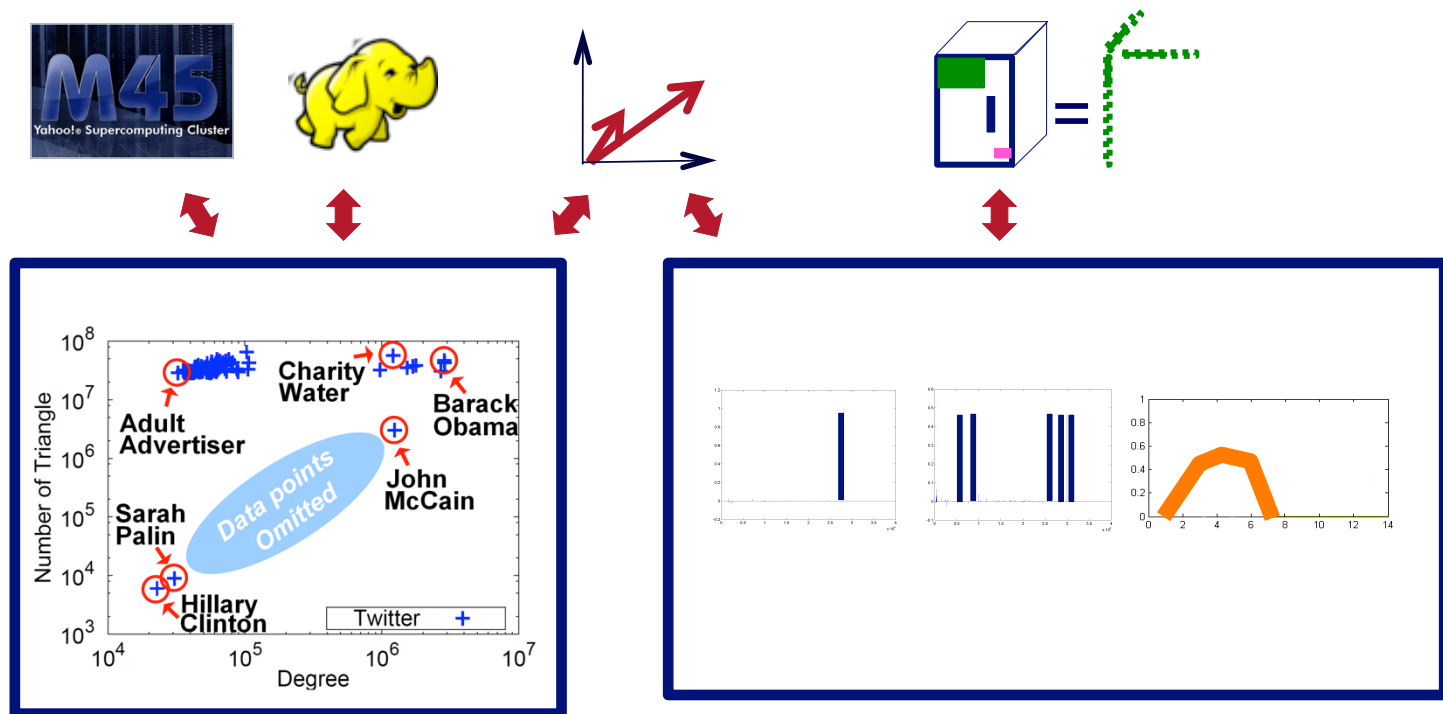
- powerful tool



1 caller      5 receivers      4 days of activity

# References

- D. Chakrabarti, C. Faloutsos: *Graph Mining – Laws, Tools and Case Studies*, Morgan Claypool 2012

- http://www.morganclaypool.com/doi/abs/10.2200/S00449ED1V01Y201209DMK006

# TAKE HOME MESSAGE:

# Cross-disciplinarity

# Thank you!

## Cross-disciplinarity