# Mining Large Graphs:
# Patterns, Anomalies, and Fraud Detection

*Christos Faloutsos*

CMU

# Thank you!
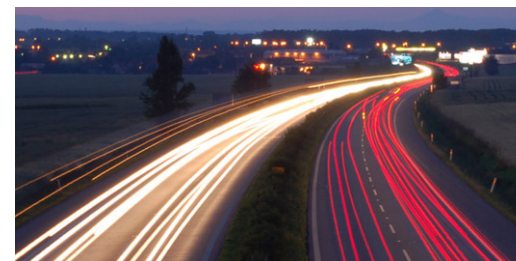
- Nina Balcan

- Kilian Weinberger

# Roadmap



→ • Introduction – Motivation
  – Why study (big) graphs?

• Part#1: Patterns in graphs

• Part#2: time-evolving graphs; tensors

• Conclusions

# Graphs - why should we care?



~1B nodes (web sites)
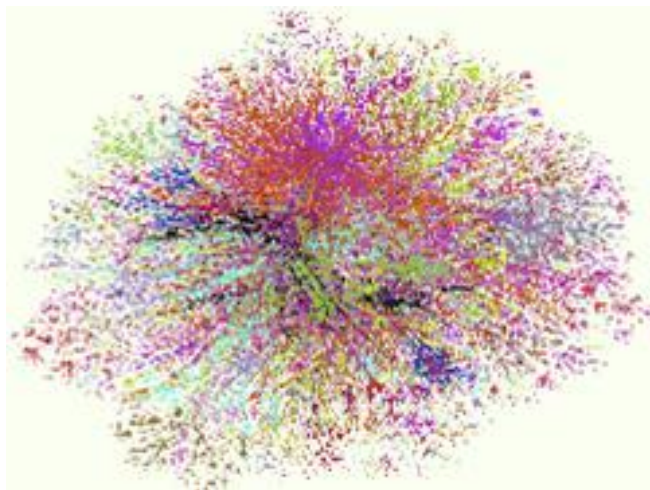~6B edges (http links)
'YahooWeb graph'
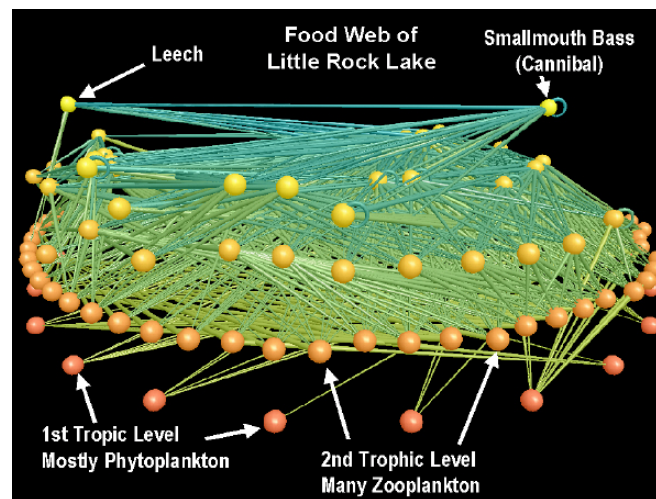
# Graphs - why should we care?

>$10B; ~1B users

# Graphs – why should we care?
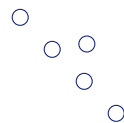


Internet Map
[lumeta.com]



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

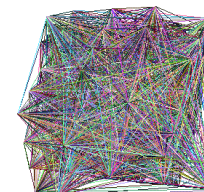- P1: patterns? Fraud detection?

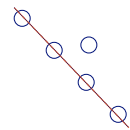- P2: patterns in time-evolving graphs / tensors

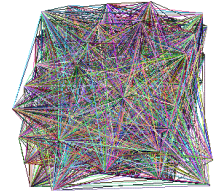destination

source    time

# Motivating problems

- P1: patterns? Fraud detection?

  Patterns ⚓ anomalies

- P2: patterns in time-evolving graphs / tensors

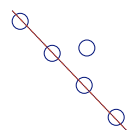  destination

  source    time
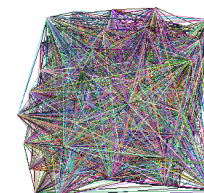
# Motivating problems

- P1: patterns? Fraud detection?

    Patterns ⋈ anomalies*

- P2: patterns in time-evolving graphs / tensors

    destination

    source    time
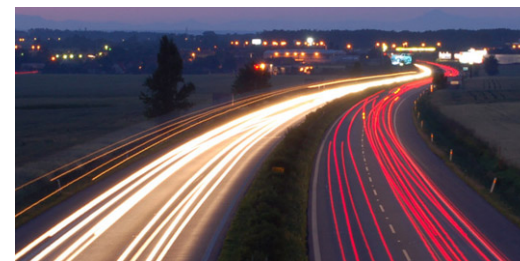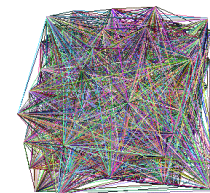
***Robust Random Cut Forest Based Anomaly Detection on Streams** Sudipto Guha, Nina Mishra , Gourav Roy,*

# Roadmap



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

(c) 2016, C. Faloutsos

# 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]
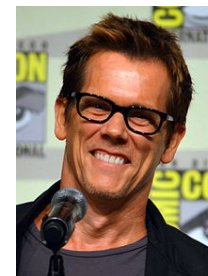
**internet domains**

**att.com**

log(degree)

**ibm.com**

```
1000  "../TESTDATA/980410-INTER/980410.Internet.outdegrees.z"
              exp(6.65065) * x ** ( -0.826118)
 100
  10
   1
 0.1
      1    10   100  1000 10000
```

log(rank)

# Solution# S.1

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

**internet domains**

**att.com**

log(degree)

**ibm.com**

**-0.82**

log(rank)

"../TESTDATA/980410-INTER/980410.Internet.outdegrees.z"
exp(6.65065) * x ** ( -0.826118 )
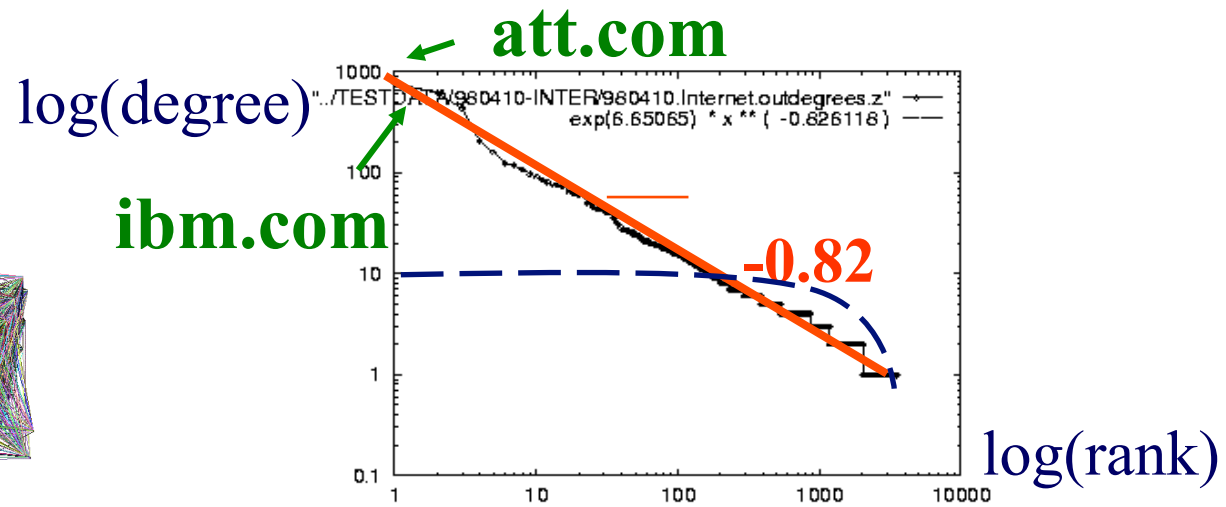
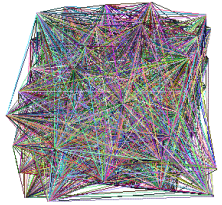(c) 2016, C. Faloutsos

# Solution# S.1

- Q: So what?

**internet domains**



**att.com**

log(degree)

**ibm.com**

**-0.82**

log(rank)

# Solution# S.1

= friends of friends (F.O.F.)

- Q: So what?
- A1: # of two-step-away pairs:

**internet domains**

**att.com**

log(degree)

**ibm.com**

"../TESTDATA/980410-INTER/980410.Internet.outdegrees.z"
exp(6.65065) * x ** ( -0.826118 )

**-0.82**

log(rank)

# Solution# S.1

- Q: So what?
- A1: # of two-step-away pairs: 100^2 * N= 10 Trillion

*= friends of friends (F.O.F.)*

**internet domains**



**att.com**

log(degree)

**ibm.com**

**-0.82**

log(rank)

# Solution# S.1

- Q: So what?

= friends of friends (F.O.F.)

- A1: # of two-step-away pairs: 100^2 ... Trillion

**internet domains**

**att.com**

log(degree)

**ibm.com**

```
"../TESTDATA/980410-INTER/980410.Internet.outdegrees.z"
exp(6.65065) * x ** ( -0.826118)
```
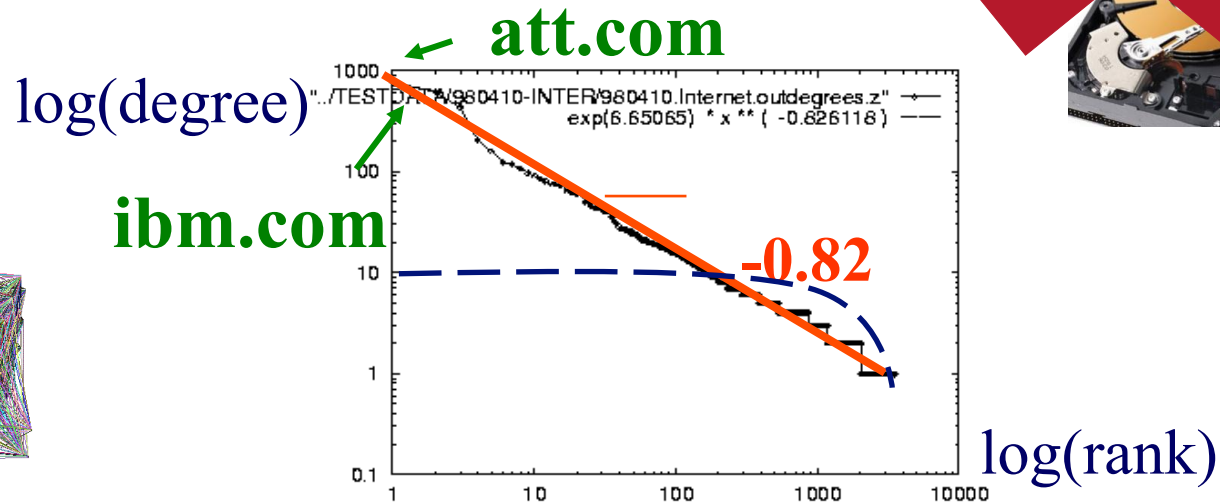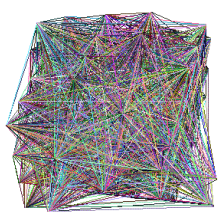
**-0.82**

log(rank)
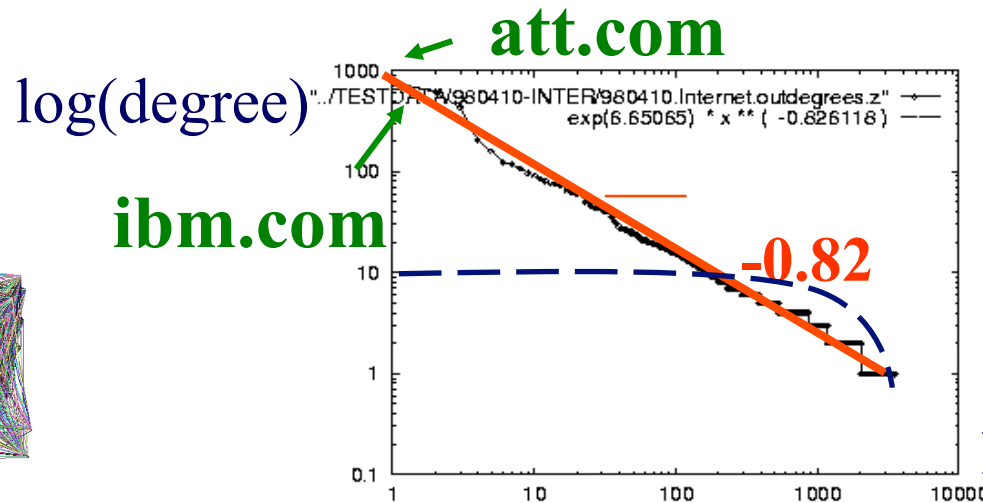
# Solution# S.1

- Q: So what?

= friends of friends (F.O.F.)

- A1: # of two-step-away pairs: O(d_max ^2) ~ 10M^2

**internet domains**

**att.com**

log(degree)

**ibm.com**

-0.82

~0.8PB ->
a data center(!)

DCO @ CMU

(c) 2016, C. Faloutsos

# Solution# S.1

**Gaussian trap**

- Q: So what?
- A1: # of two-step-aw~ ~?) ~ 10M^2
  inte~

**Such patterns -> New algorithms**

$\Downarrow$

~0.8PB -> a data center(!)

**-0.82**

0.1    1    10    100    1000    10000

(c) 2016, C. Faloutsos

# **Observation – big-data:**

- O($N^2$) algorithms are ~intractable  - N=1B

- $N^2$ seconds = 31B years (>2x age of universe)

1B

1B

(c) 2016, C. Faloutsos

# Observation – big-data:

- O($N^2$) algorithms are ~intractable  - N=1B

31M

- $N^2$ seconds = 31B years
- 1,000 machines

1B

(c) 2016, C. Faloutsos

# **Observation – big-data:**

- O($N^2$) algorithms are ~intractable  - N=1B

31K

- $N^2$ seconds = 31B years

- 1M machines

1B

(c) 2016, C. Faloutsos
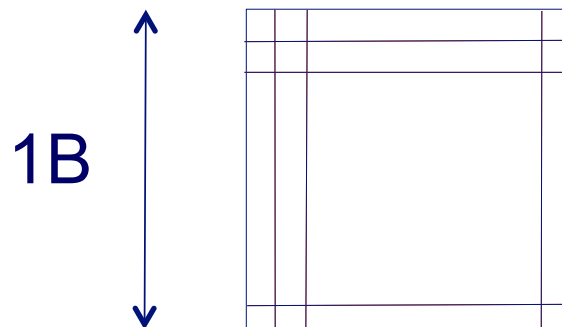
# Observation – big-data:

- O($N^2$) algorithms are ~intractable  - N=1B

3

- $N^2$ seconds = 31B years
- 10B machines ~ $10Trillion

1B

(c) 2016, C. Faloutsos

# Observation – big-data:

- O($N^2$) algorithms are ~intractable  - N=1B

## And parallelism might not help

- $N^2$ seconds = 31B years    3

- 10B machines ~ $10Trillion

1B

(c) 2016, C. Faloutsos

# Solution# S.2: Eigen Exponent *E*

Eigenvalue



Exponent = slope

*E = -0.48*

**A x = λ x**

Rank of decreasing eigenvalue

- A2: power law in the eigenvalues of the adjacency matrix ('eig()')

(c) 2016, C. Faloutsos

# Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
  - → Patterns: Degree; Triangles
  - Anomaly/fraud detection
  - Graph understanding
- 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 ?

(c) 2016, C. Faloutsos

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



Reuters

slope 1.68
slope -1.68
DTPL

SN

slope 1.74
slope -1.73
DTPL

Epinions

slope 1.61
slope -1.59
DTPL

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

# Triangle Law: Computations
## [Tsourakakis ICDM 2008]



But: triangles are expensive to compute
(3-way join; several approx. algos) – $O(d_{max}^2)$
Q: Can we do that quickly?
A:

**details**

# Triangle Law: Computations
## [Tsourakakis ICDM 2008]

But: triangles are expensive to compute
  (3-way join; several approx. algos) – $O(d_{max}^2)$
Q: Can we do that quickly?
A: Yes!

$A x = \lambda x$

**#triangles = 1/6 Sum ( $\lambda_i^3$ )**
  (and, because of skewness (S2) ,
    we only need the top few eigenvalues! - O(E)

# Triangle counting for large graphs?



Anomalous nodes in Twitter(~ 3 billion edges)

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

(c) 2016, C. Faloutsos

# 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]

# MORE Graph Patterns

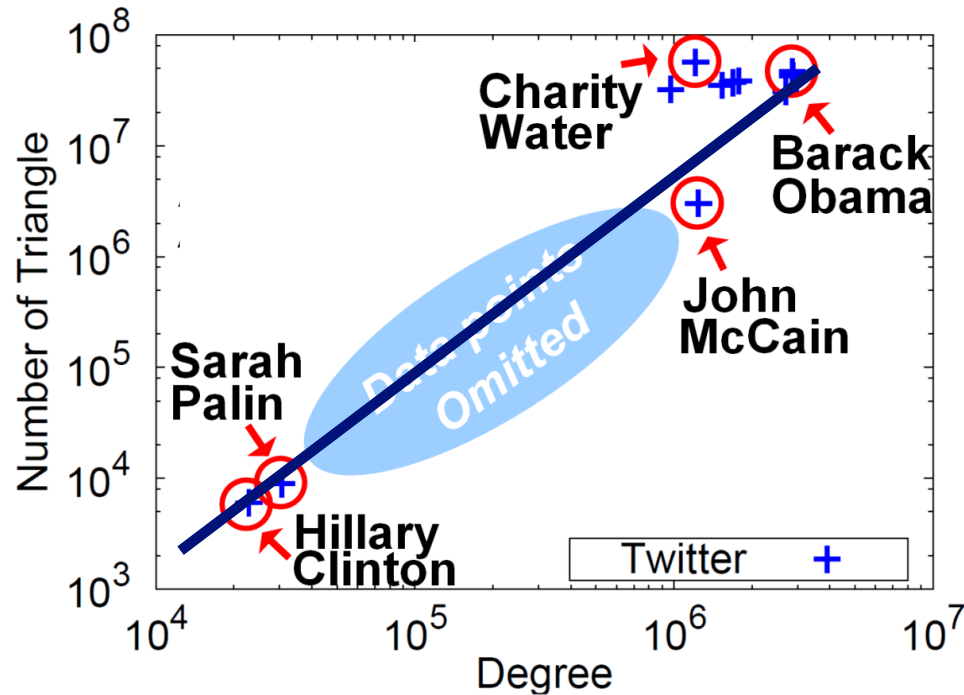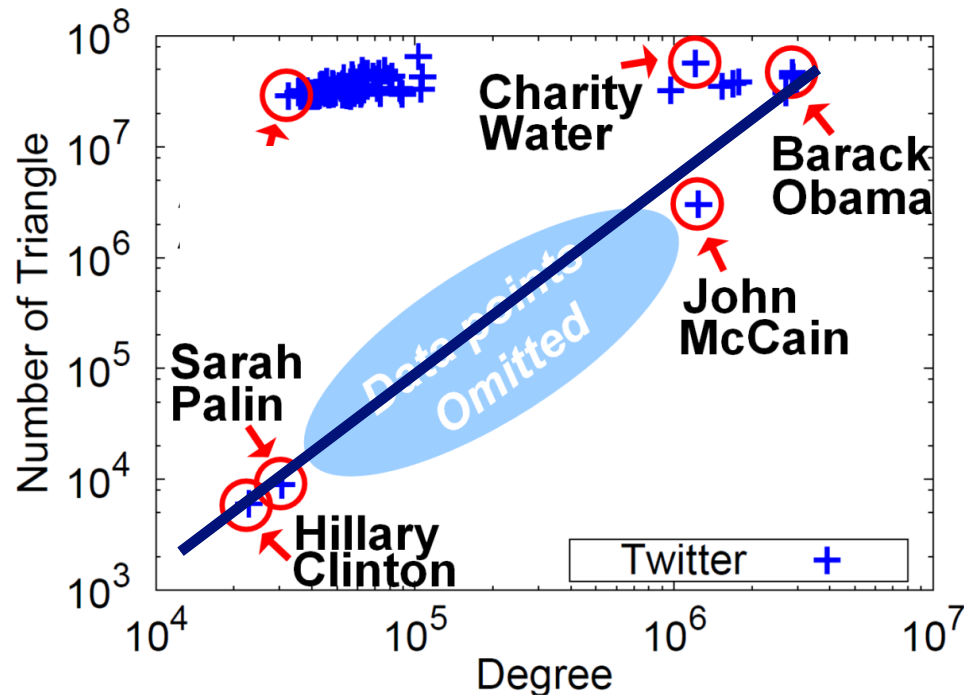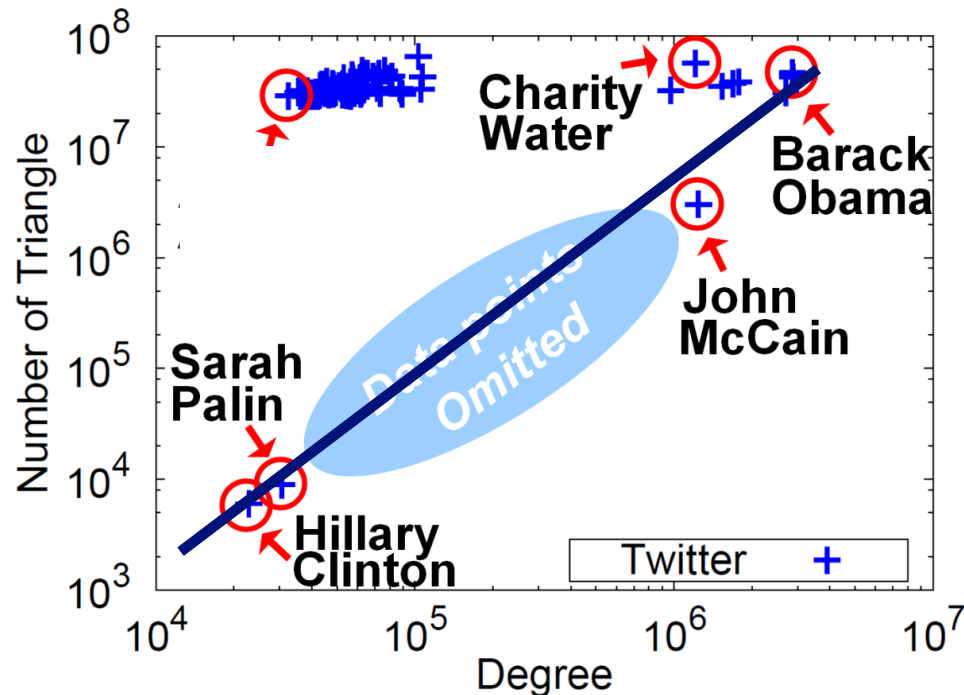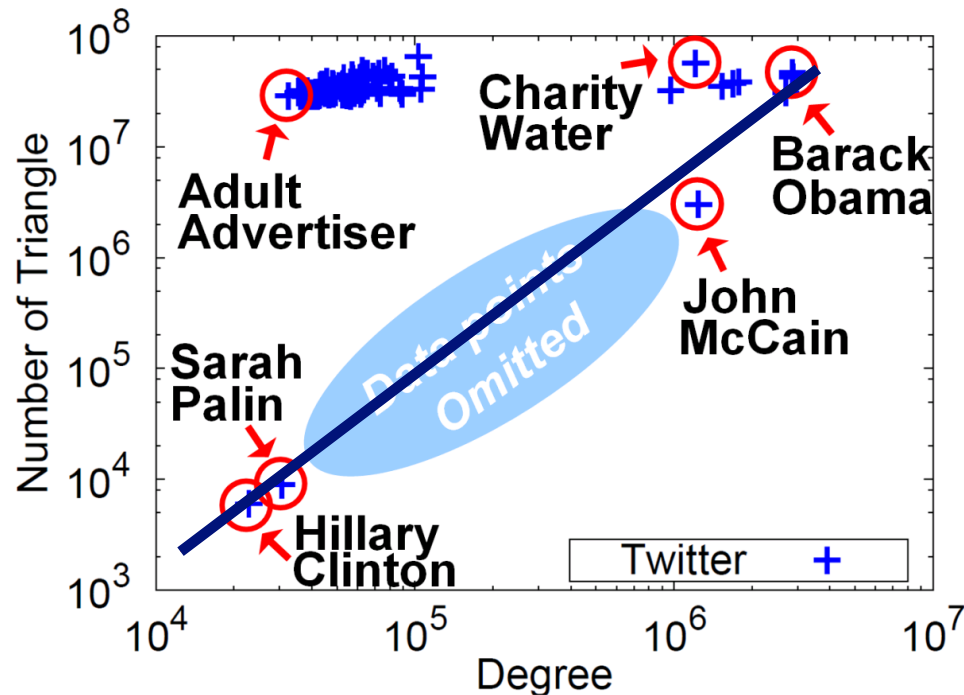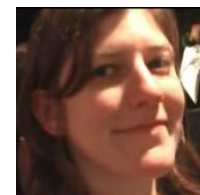| | Unweighted | Weighted |
|---|---|---|
| **Static** | ✓ **L01.** Power-law degree distribution [Faloutsos et al. `99, Kleinberg et al. `99, Chakrabarti et al. `04, Newman `04]<br>✓ **L02.** Triangle Power Law (TPL) [Tsourakakis `08]<br>✓ **L03.** Eigenvalue Power Law (EPL) [Siganos et al. `03]<br>**L04.** 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.
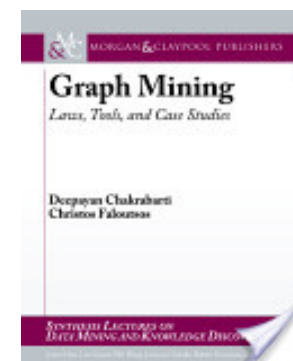
# MORE Graph Patterns

| | 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)
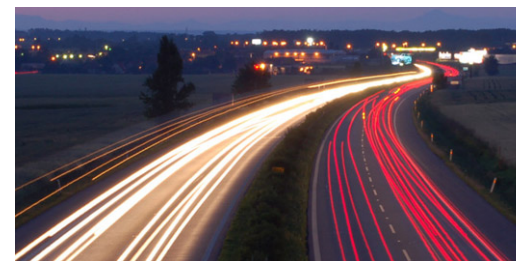


• Deepayan Chakrabarti and Christos Faloutsos, *Graph Mining: Laws, Tools, and Case Studies* Oct. 2012, Morgan Claypool.



**http://www.cs.cmu.edu/~christos/TALKS/16-06-19-ICML/**

# Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
  - Patterns
  - Anomaly / fraud detection
    - CopyCatch
    - Spectral methods ('fBox')          Patterns &harr; anomalies
    - Belief Propagation
- Part#2: time-evolving graphs; tensors
- Conclusions

# Fraud

- Given
  - Who 'likes' what page, and when

- Find
  - Suspicious users and suspicious products

# Fraud



Users     Pages

- Given
  - Who 'likes' what page, and when

- Find
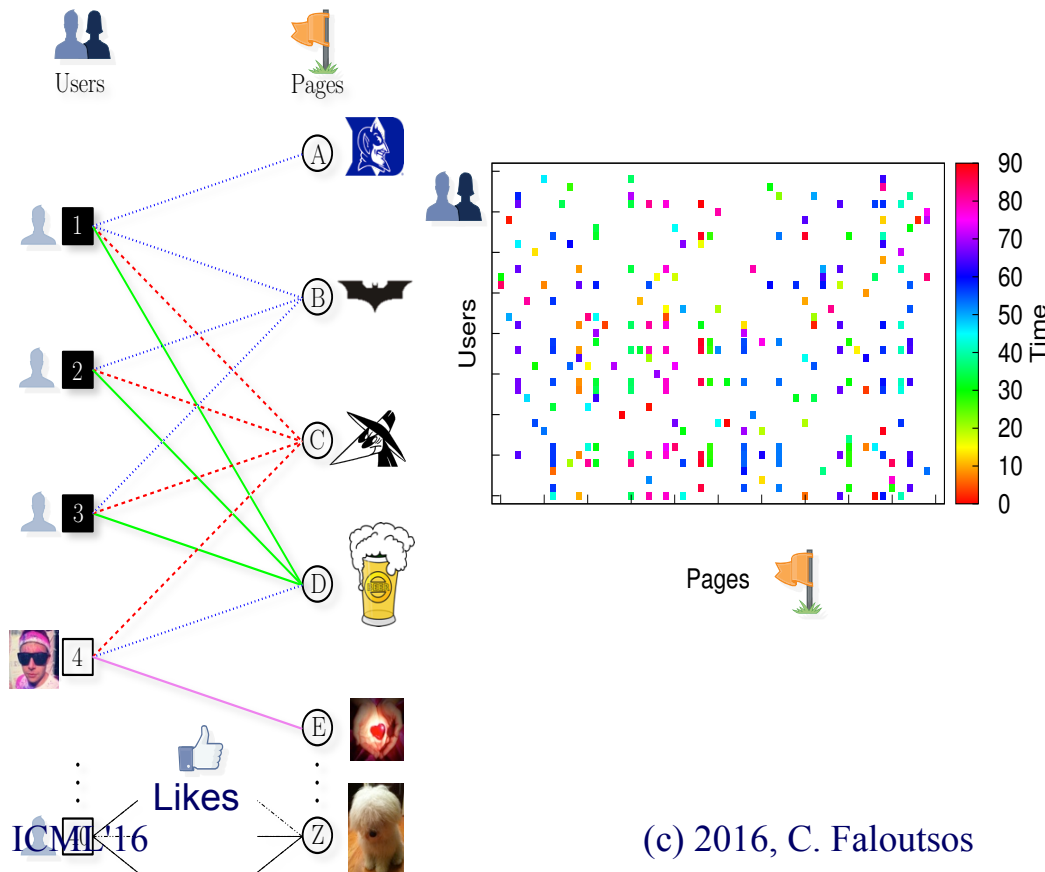  - Suspicious users and suspicious products

# Graph Patterns and Lockstep Behavior

## Our intuition

- Lockstep behavior: Same Likes, same time

Users        Pages

(c) 2016, C. Faloutsos

45
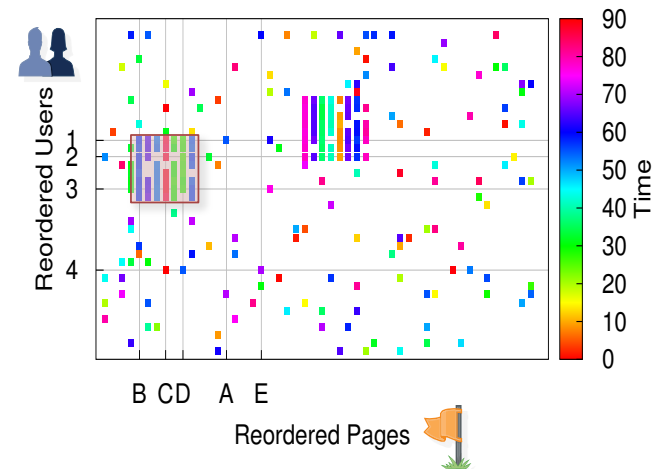
# Graph Patterns and Lockstep Behavior

## Our intuition

- Lockstep behavior: Same Likes, same time



Users   Pages

Likes

(c) 2016, C. Faloutsos

46
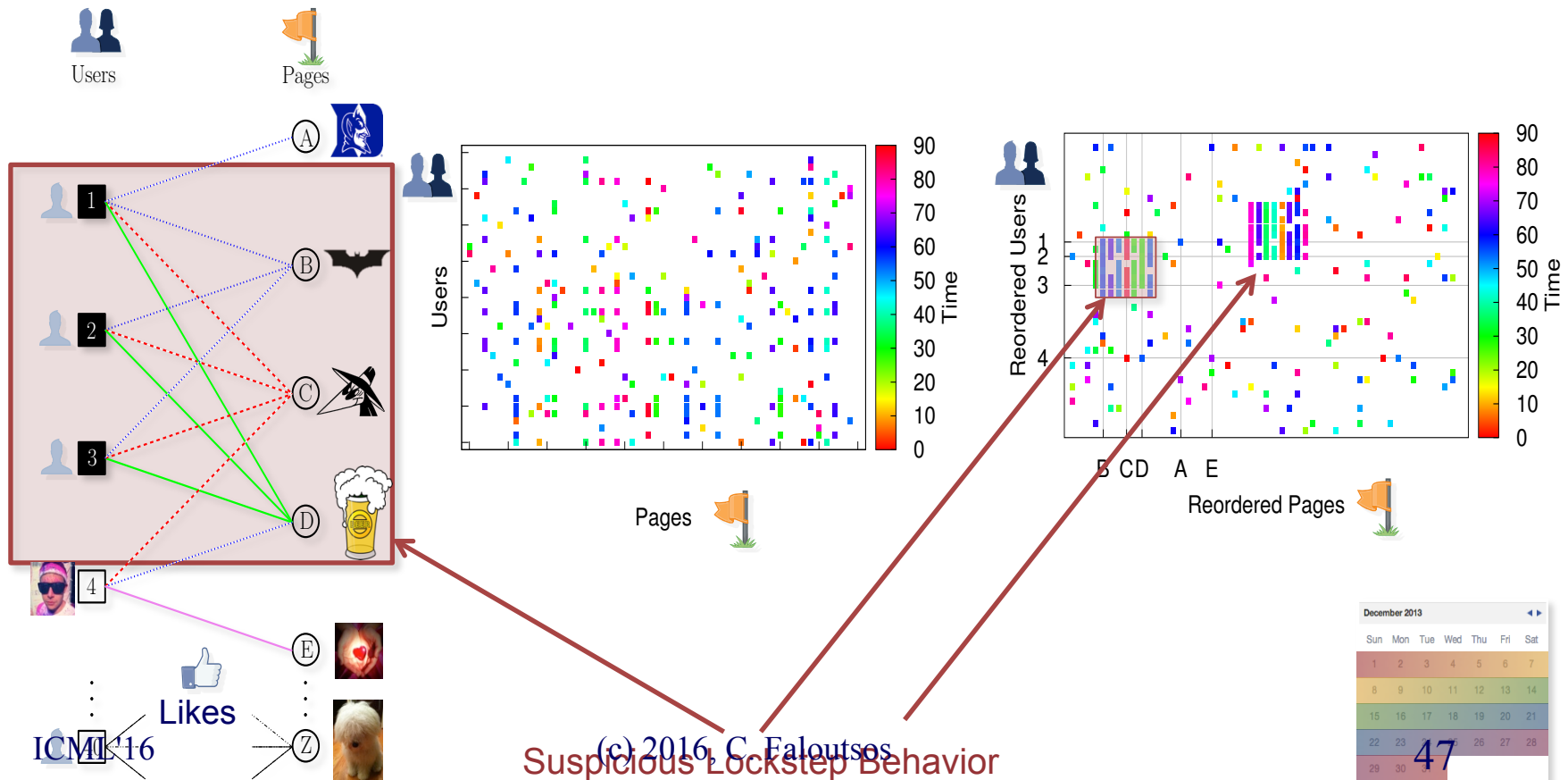
# Graph Patterns and Lockstep Behavior

## Our intuition

- Lockstep behavior: Same Likes, same time

Users    Pages



Suspicious Lockstep Behavior

47

# MapReduce Overview

Users      Pages

- Use Hadoop to search for many clusters in parallel:

  1. **Start** with randomly seed

  2. **Update** set of Pages and center Like times for each cluster

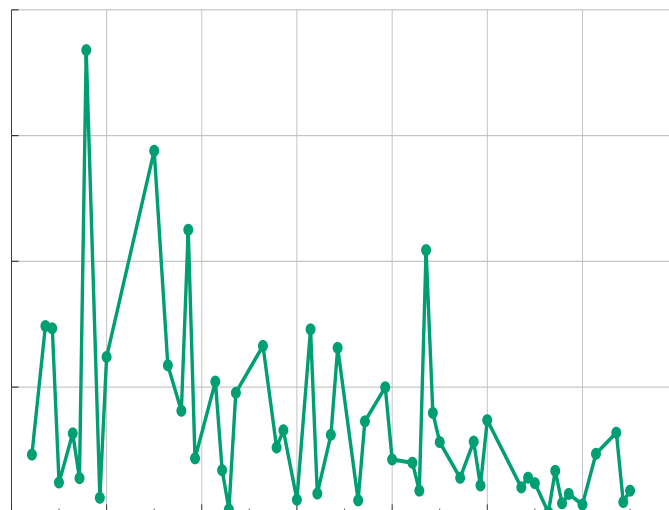  3. **Repeat** until convergence

Likes
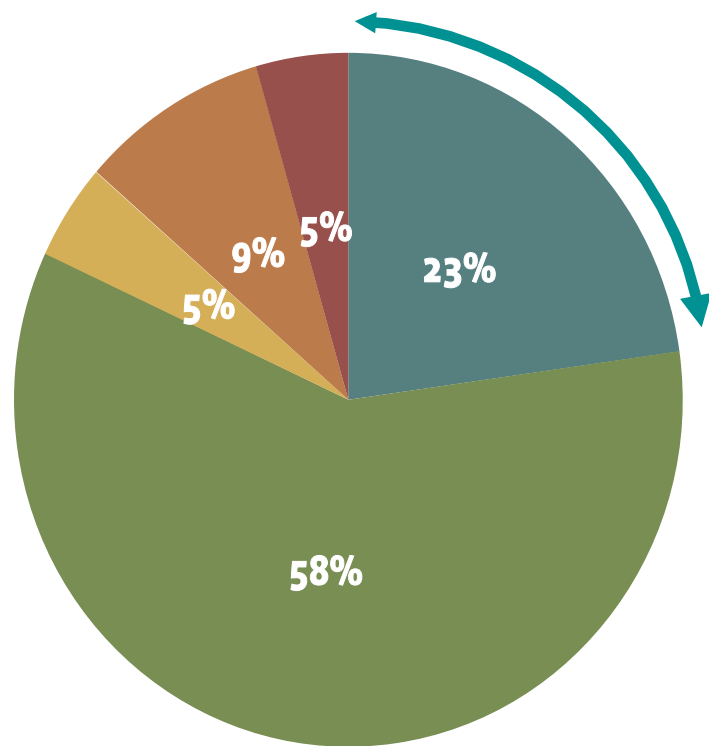
# Deployment at Facebook

- *CopyCatch* runs regularly (along with many other security mechanisms, and a large Site Integrity team)

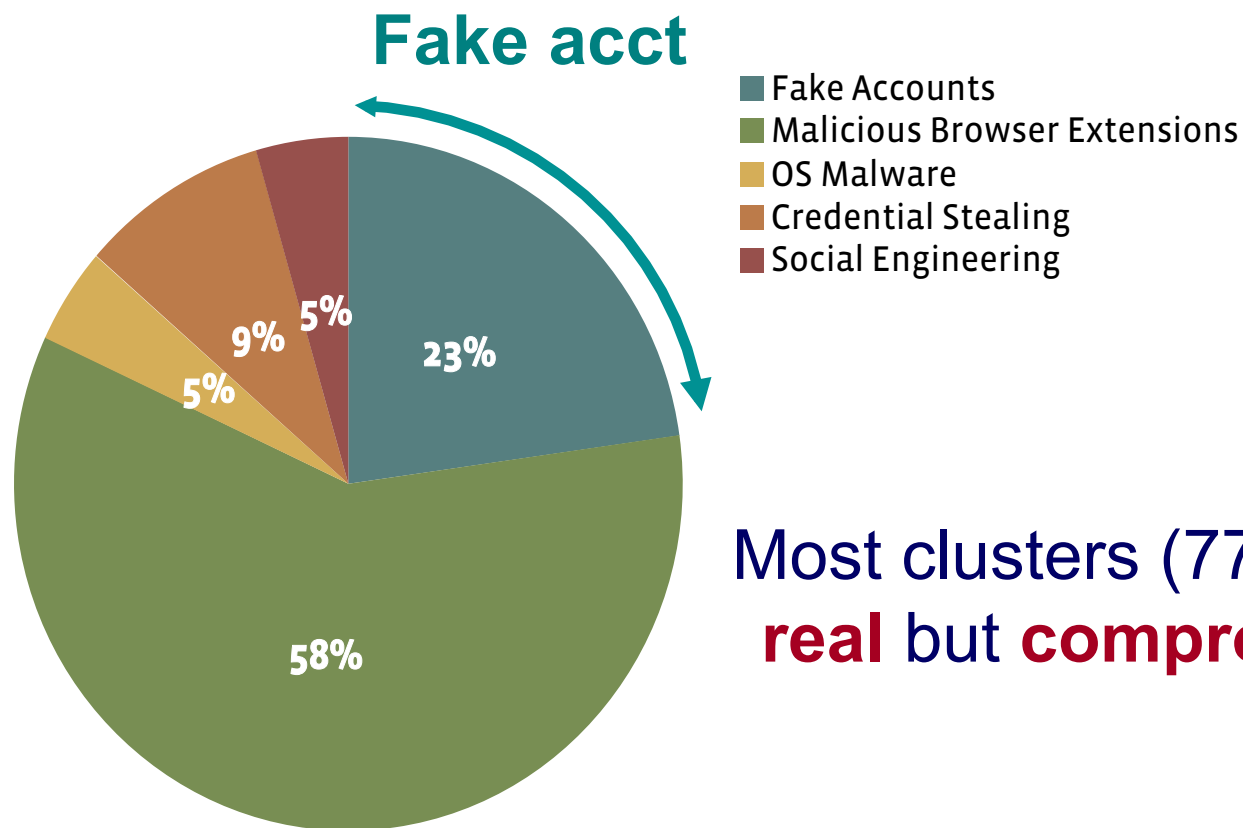3 months of *CopyCatch* @ Facebook

#users caught

(c) 2016, C. Faloutsos

time

# Deployment at Facebook



Manually labeled 22 randomly selected
*clusters* from February 2013

(c) 2016, C. Faloutsos

# Deployment at Facebook

**Fake acct**

Fake Accounts
Malicious Browser Extensions
OS Malware
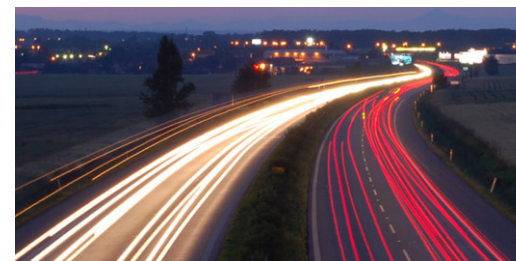Credential Stealing
Social Engineering

23%

58%

5%

9%

5%

Most clusters (77%) come from **real** but **compromised** users

Manually labeled 22 randomly selected *clusters* from February 2013

# Roadmap



- Introduction – Motivation
- Part#1: Patterns in graphs
  - Patterns
  - Anomaly / fraud detection
    - CopyCatch
    - Spectral methods ('fBox')
    - Belief Propagation
- Part#2: time-evolving graphs; tensors
- Conclusions

# Problem: Social Network Link Fraud

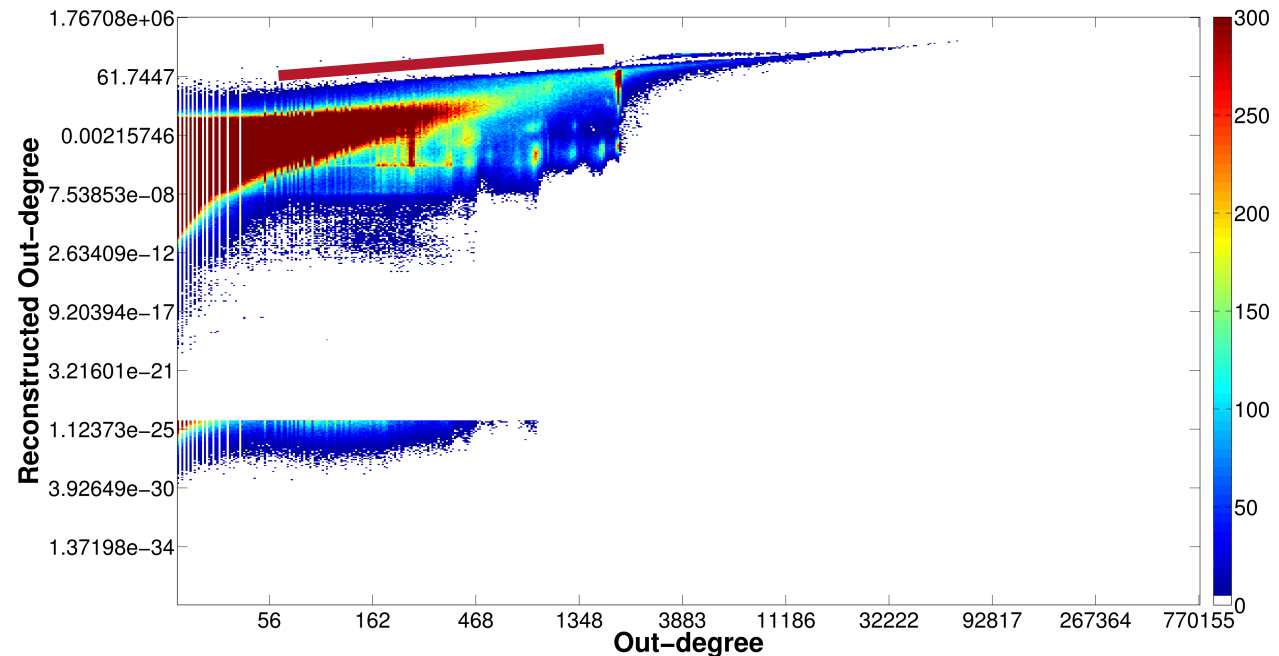Target: find "stealthy" attackers missed by other algorithms

Clique

41.7M nodes
1.5B edges

Bipartite
core

(c) 2016, C. Faloutsos

# Problem: Social Network Link Fraud

Target: find "stealthy" attackers missed by other algorithms

**Takeaway:** use *reconstruction error* between true/latent representation!

Neil Shah, Alex Beutel, Brian Gallagher and Christos Faloutsos. *Spotting Suspicious Link Behavior with fBox: An Adversarial Perspective.* ICDM 2014, Shenzhen, China.
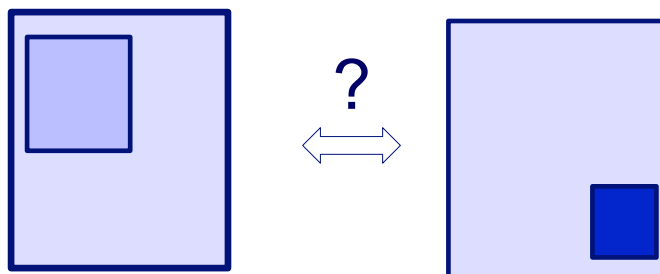
# Roadmap



- Introduction – Motivation
- Part#1: Patterns in graphs
  - Patterns
  - Anomaly / fraud detection
    - CopyCatch
    - Spectral methods ('fBox', **suspiciousness**)
    - Belief Propagation
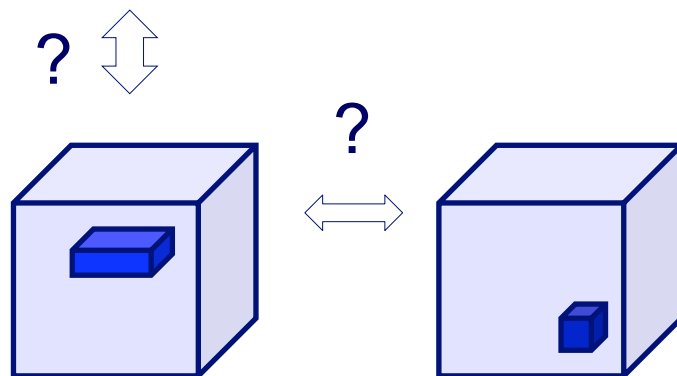- Part#2: time-evolving graphs; tensors
- Conclusions

# 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.
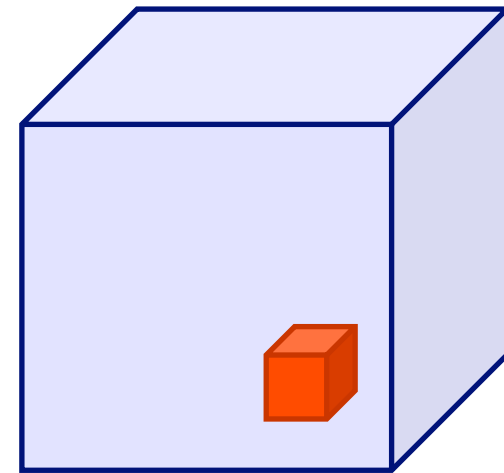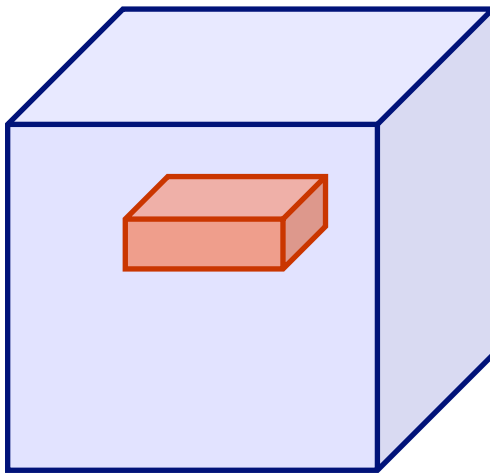
# Suspicious Patterns in Event Data

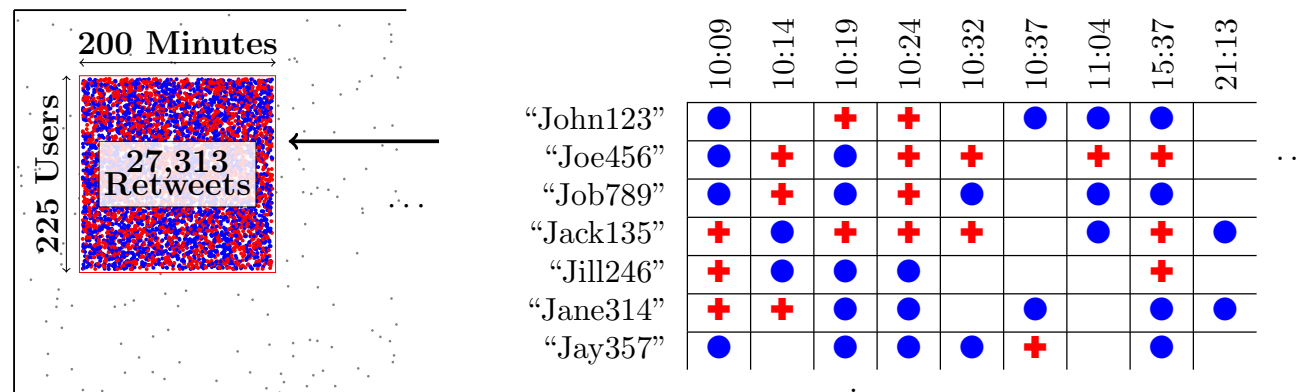## Which is more suspicious?

20,000 Users
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 * $D_{KL}(p || p_{background})$**

Analyze

# Suspicious Patterns in Event Data



200 Minutes

225 Users

27,313 Retweets

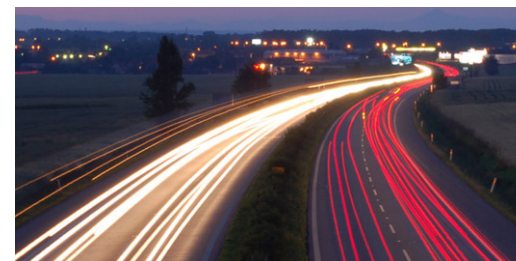|  | 10:09 | 10:14 | 10:19 | 10:24 | 10:32 | 10:37 | 11:04 | 15:37 | 21:13 |
|---|---|---|---|---|---|---|---|---|---|
| "John123" | ● |  | ✚ | ✚ |  | ● | ● | ● |  |
| "Joe456" | ● | ✚ | ● | ✚ | ✚ |  | ✚ | ✚ |  |
| "Job789" | ● | ✚ | ● | ✚ | ● |  | ● | ● |  |
| "Jack135" | ✚ | ● | ✚ | ✚ | ✚ |  | ● | ✚ | ● |
| "Jill246" | ✚ | ● | ● | ● |  |  |  | ✚ |  |
| "Jane314" | ✚ | ✚ | ● | ● |  | ● |  | ● | ● |
| "Jay357" | ● |  | ● | ● | ● | ✚ |  | ● |  |

Retweeting: "Galaxy Note Dream Project: Happy Happy Life Traveling the World"

|  | # | User × tweet × IP × minute | Mass $c$ | Suspiciousness |
|---|---|---|---|---|
| **CROSSSPOT** | 1 | 14×1×2×1,114 | 41,396 | 1,239,865 |
|  | 2 | 225×1×2×200 | 27,313 | 777,781 |
|  | 3 | 8×2×4×1,872 | 17,701 | 491,323 |
| HOSVD | 1 | 24×6×11×439 | 3,582 | 131,113 |
|  | 2 | 18×4×5×223 | 1,942 | 74,087 |
|  | 3 | 14×2×1×265 | 9,061 | 381,211 |

# Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
  - Patterns
  - Anomaly / fraud detection
    - CopyCatch
    - Spectral methods ('fBox')
    - (Matrix re-ordering + education -> 'groupNteach')
    - Belief Propagation
- Part#2: time-evolving graphs; tensors
- Conclusions

# Problem dfn:
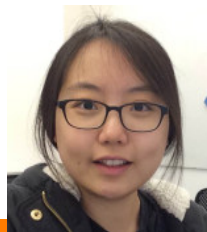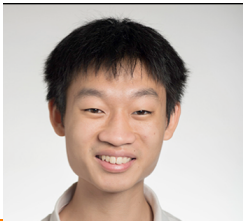
e.g.

|  | ears | fins | stripes | lungs | gills | carnivore | ... |
|--------|------|------|---------|-------|-------|-----------|-----|
| salmon |      | ●    |         |       | ●     |           |     |
| tiger  | ●    |      | ●       | ●     |       | ● (v)     |     |
| jaguar | ●    |      |         | ●     |       | ●         |     |
| tuna   |      | ●    |         |       | ●     |           |     |
| lion   | ●    |      |         | ●     |       | ●         |     |
| ⋮      |      |      |         |       |       |           |     |

|  | ears | lungs | carnivore | stripes | fins | gills | ... |
|--------|------|-------|-----------|---------|------|-------|-----|
| lion   | ●    | ●     | ●         |         |      |       |     |
| tiger  | ●    | ●     | ●         | ●       |      |       |     |
| jaguar | ●    | ●     | ●         |         |      |       |     |
| tuna   |      |       |           |         | ●    | ●     |     |
| salmon |      |       |           |         | ●    | ●     |     |
| ⋮      |      |       |           |         |      |       |     |

# Problem definition

- **Given** a large binary matrix of facts of *(object, property)* pairs

- **Find** *groupings* of the facts and the *order* of transmission

- To **optimize** 'student effort' (**->** incremental learning curve, '*ALOC*')



Bryan Hooi, Hyun Ah Song, et al, "*Matrices, Compression, Learning Curves: Formulation, and the GroupNTeach Algorithms*", PAKDD 2016
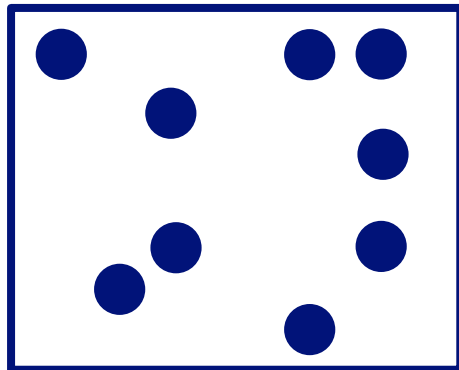
# Details:

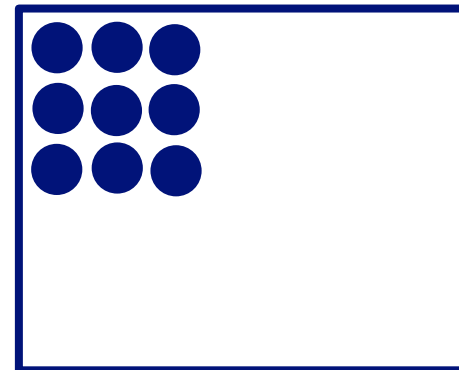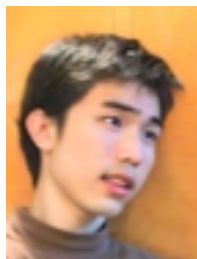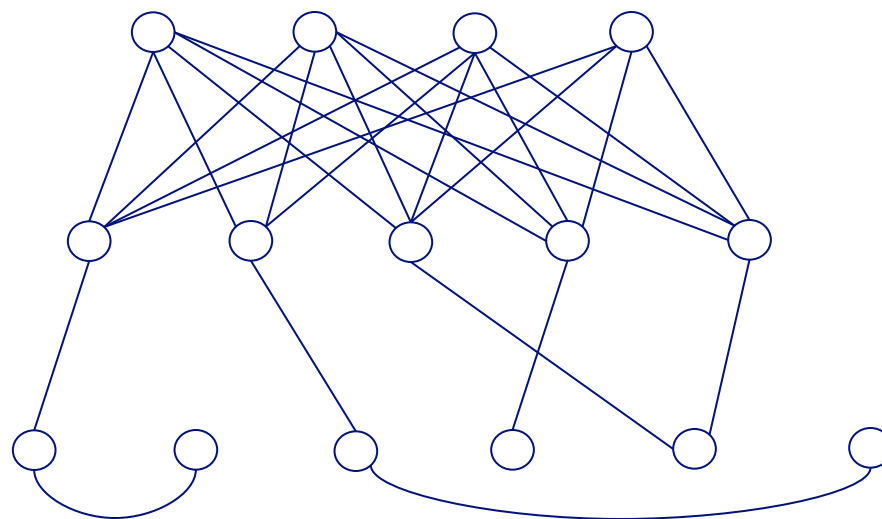Given a large binary matrix of objects and properties, re-order rows and columns,

**G1. Metric** for better encoding of matrix for student learning?

**G2**. How do we construct **language** to describe it?

**G3.** How do we **optimize** this metric?

**vs.**

# Pictorial Problem definition



**vs.** **vs.**

*ALOC* *ALOC* *ALOC*

#dots learned

#bits transmitted

Score

(c) 2016, C. Faloutsos

# Results



Side effects

Side effects

Drugs

Drugs

**Teaching order**

Anti-depressants    Hyper-tension    Pain

# Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
  - Patterns
  - Anomaly / fraud detection
    - CopyCatch
    - Spectral methods ('fBox')
    - 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

(c) 2016, C. Faloutsos

# E-bay Fraud detection

(c) 2016, C. Faloutsos

# E-bay Fraud detection - NetProbe

(c) 2016, C. Faloutsos

# 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
    - CopyCatch
    - Spectral methods ('fBox')
    - Belief Propagation; fast computation & unification
- Part#2: time-evolving graphs; tensors
- Conclusions

**Carnegie Mellon**

# 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

*ECML PKDD, 5-9 September 2011, Athens, Greece*

# 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*')
- 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** | $[\mathbf{I} + a\,\mathbf{D} - c'\underline{\mathbf{A}}]$ | $\times$ | $\mathbf{b_h}$ | $=$ | $\boldsymbol{\phi_h}$ |

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

adjacency
matrix

final
labels/
beliefs

prior
labels/
beliefs

# Results: Scalability



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

# Summary of Part#1

- *many* patterns in real graphs
  - Power-laws everywhere
  - Gaussian trap
    - Avg << Max
  - 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; tensors
  - P2.1: time-evolving graphs
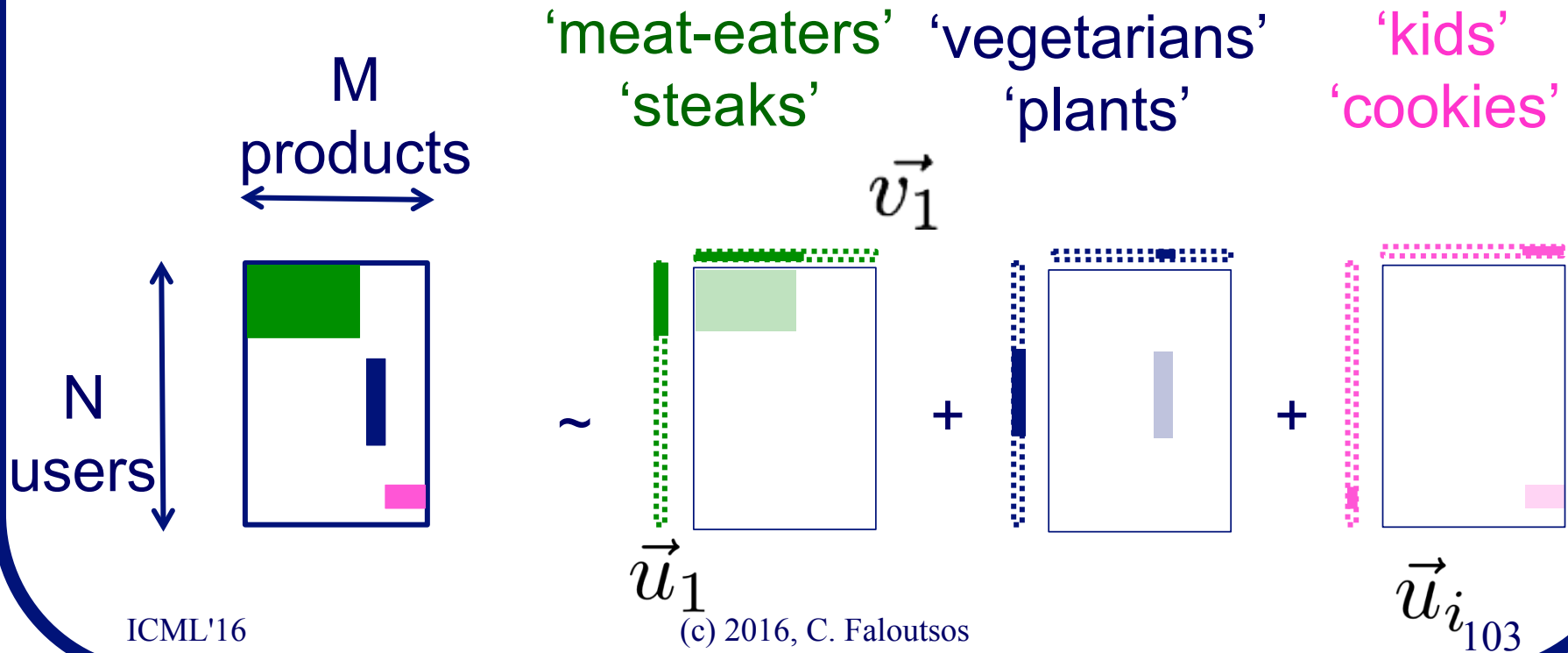  - [P2.2: with side information ('coupled' M.T.F.)
  - Speed]
- Conclusions

# 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

(c) 2016, C. Faloutsos

# 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)

**SKIP**

# Graphs & side info

- Problem #2.2: coupled (eg., side info)
  - Given subject – verb – object facts
    - And voxel-activity for each subject-word
  - Find patterns / anomalies

verb

subject

object

fMRI voxel activity

`apple tastes sweet'

(c) 2016, C. Faloutsos

100

# Graphs & side info

- Problem #2.2: coupled (eg., side info)
  - Given subject – verb – object facts
    - And voxel-activity for each subject-word
  - Find patterns / anomalies

'tastes'

'apple'

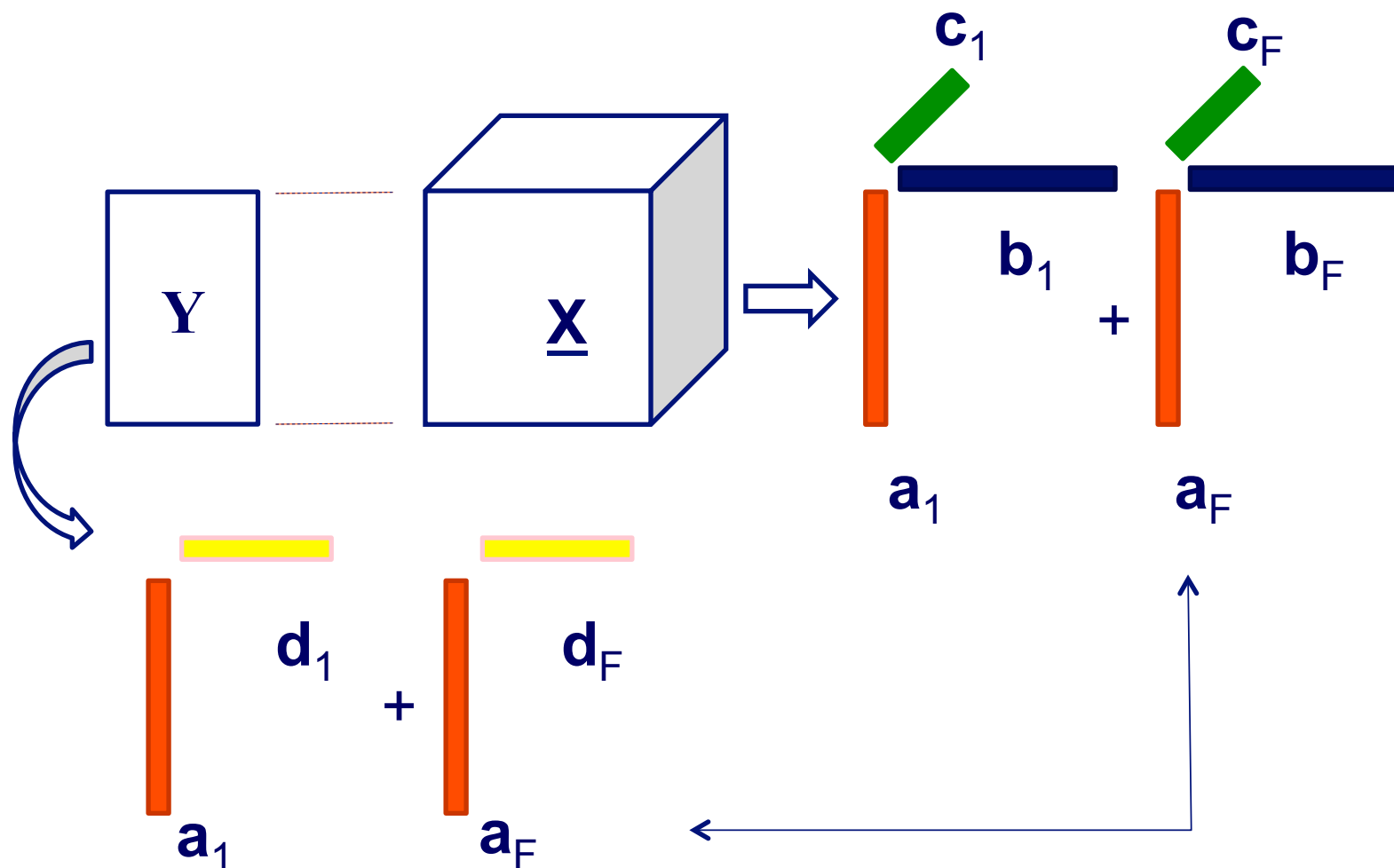'sweet'

`apple tastes sweet'

'apple'

fMRI voxel activity

# Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
    - P2.1: time-evolving graphs
    - [P2.2: with side information ('coupled' M.T.F.)
    - Speed]
- Conclusions

# Answer to both: tensor factorization

- Recall: (SVD) matrix factorization: finds blocks



'meat-eaters'
'steaks'

'vegetarians'
'plants'

'kids'
'cookies'

M products

N users

$\sim$   $\vec{u}_1$   $\vec{v}_1$   $+$   $+$   $\vec{u}_{i_{103}}$

# Answer to both: tensor factorization

- PARAFAC decomposition

# Answer: tensor factorization

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



(c) 2016, C. Faloutsos

# 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

**SKIP**

- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
  - P2.1: Discoveries @ phonecall network
  - [P2.2: Discoveries in neuro-semantics
  - Speed]
- Conclusions

# Coupled Matrix-Tensor Factorization (CMTF)

SKIP

# Neuro-semantics

- **Brain Scan Data***
  - 9 persons
  - 60 nouns
- **Questions**
  - 218 questions
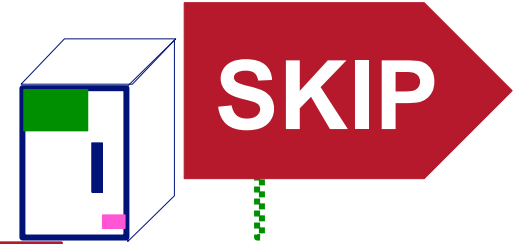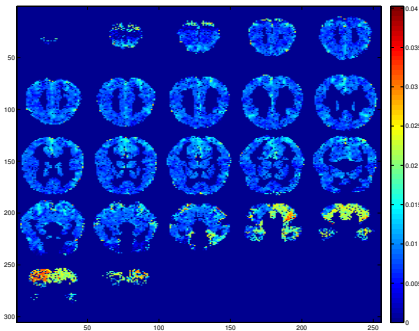  - 'is it alive?', 'can you eat it?'

# Neuro-semantics

- **Brain Scan Data***
  - 9 persons
  - 60 nouns
- **Questions**
  - 218 questions
  - 'is it alive?', 'can you eat it?'

## Patterns?

(c) 2016, C. Faloutsos

# **Neuro-semantics**

**Nouns**

- beetle
- pants
- bee

**Questions**

- can it cause you pain?
- do you see it daily?
- is it conscious?

**Nouns**

- bear
- cow
- coat

**Questions**

- does it grow?
- is it alive?
- was it ever alive?

**Nouns**

- glass
- tomato
- bell

**Questions**

- can you pick it up?
- can you hold it in one hand?
- is it smaller than a golfball?'

**Nouns**

- bed
- house
- car

**Questions**

- does it use electricity?
- can you sit on it?
- does it cast a shadow?

Premotor Cortex

**Group1**   **Group 2**   **Group 3**   **Group 4**
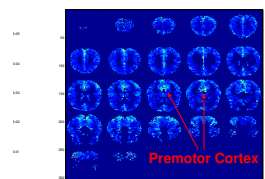
# Neuro-semantics

SKIP

## Small items -> Premotor cortex

**Nouns**

glass
tomato
bell

**Questions**

can you pick it up?
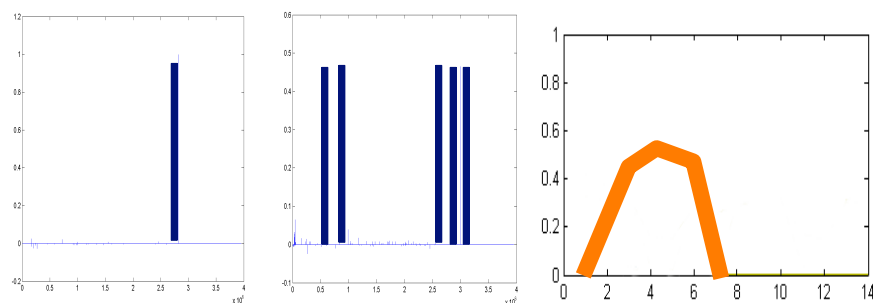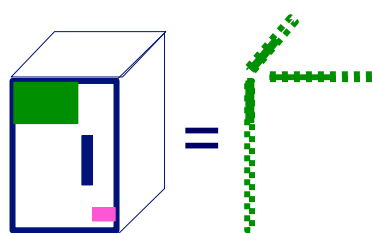can you hold it in one hand?
is it smaller than a golfball?'



**Premotor Cortex**

**Group 3**

# Neuro-semantics

**Small items ->
Premotor cortex**



**Nouns**
glass
tomato
bell

**Questions**
can you pick it up?
can you hold it in one hand?
is it smaller than a golfball?'

Premotor Cortex

**Group 3**





Evangelos Papalexakis, Tom Mitchell, Nicholas Sidiropoulos, Christos Faloutsos, Partha Pratim Talukdar, Brian Murphy, *Turbo-SMT: Accelerating Coupled Sparse Matrix-Tensor Factorizations by 200x*, SDM 2014
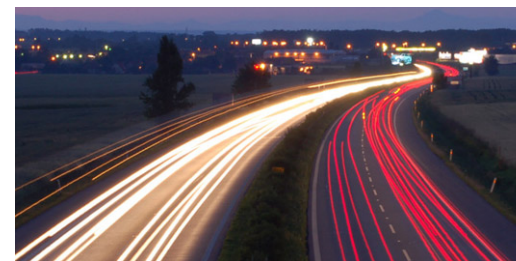
# Part 2: Conclusions

- Time-evolving / heterogeneous graphs -> tensors

- PARAFAC finds patterns

- (GigaTensor/HaTen2 -> fast & scalable)

# Roadmap

- ## Introduction – Motivation
  - ### Why study (big) graphs?
- ## Part#1: Patterns in graphs
- ## Part#2: time-evolving graphs; tensors
- ## ➡ Acknowledgements and Conclusions
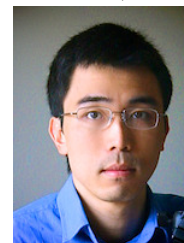
(c) 2016, C. Faloutsos

# Thanks

# Cast

Carnegie Mellon

Akoglu, Leman

Araujo, Miguel

Beutel, Alex

Chau, Polo

Hooi, Bryan

Kang, U

Koutra, Danai

Papalexakis, Vagelis

Shah, Neil

Song, Hyun Ah

# CONCLUSION#1 – Big data
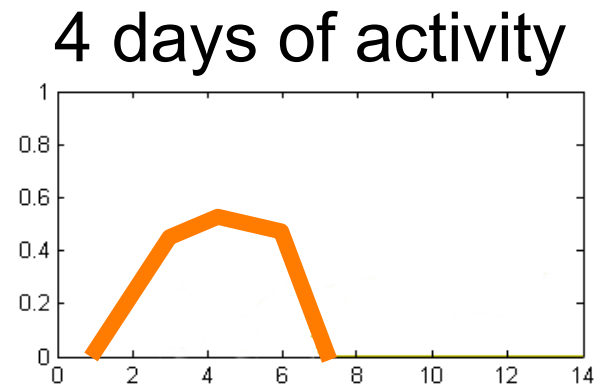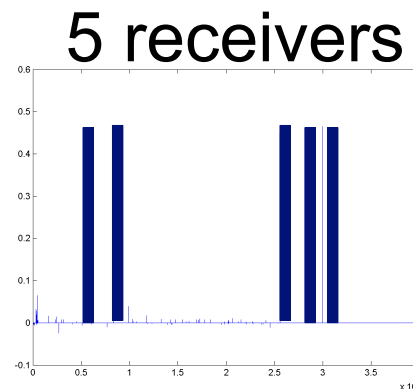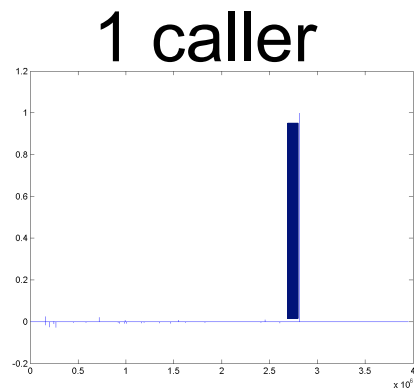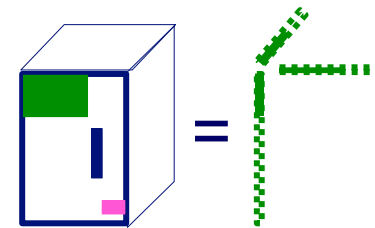
- **Patterns** **Anomalies**

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

(c) 2016, C. Faloutsos
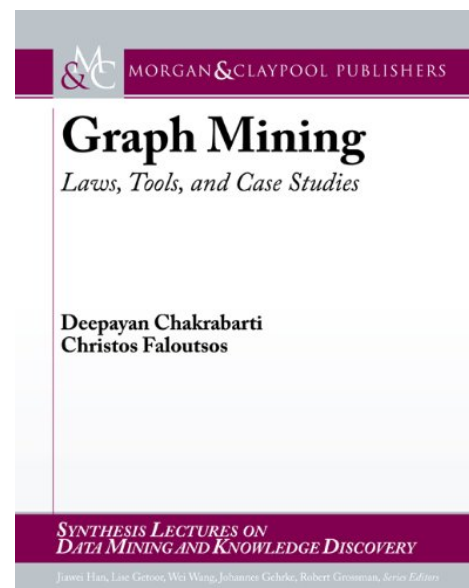
# CONCLUSION#2 – tensors

- powerful tool



1 caller     5 receivers     4 days of activity
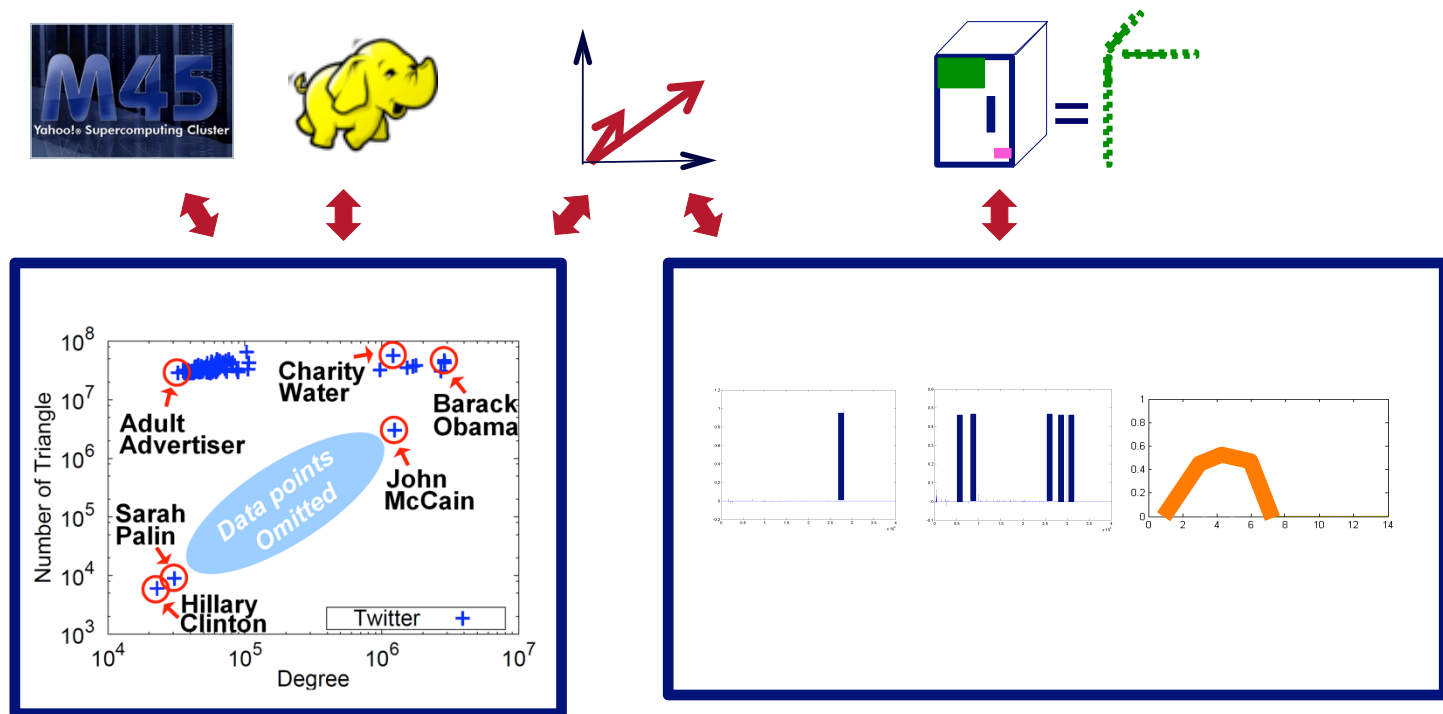
# References

- 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



(c) 2016, C. Faloutsos

**Thank you!**

**Cross-disciplinarity**