

# Large Graph Mining: Patterns, Cascades, Fraud Detection, and Algorithms

*Christos Faloutsos*

CMU

# Thank you!

- Prof. Chin-Wan Chung



# Thank you!

- Prof. Chin-Wan Chung



# Roadmap

- ➔ • Introduction – Motivation
  - Why study (big) graphs?
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
- Part#3: Cascades and immunization
- Conclusions

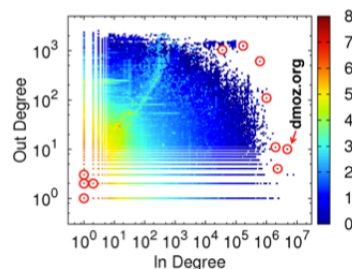


# Graphs - why should we care?

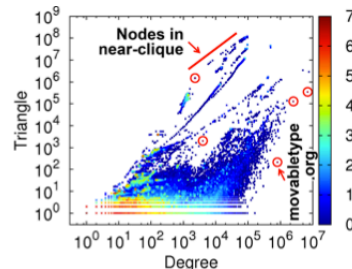


~1B nodes (web sites)  
~6B edges (http links)  
'YahooWeb graph'

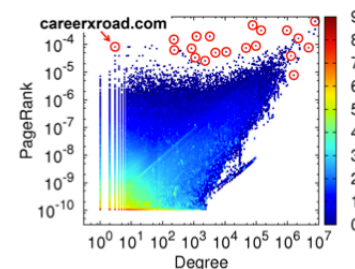
# Graphs - why should we care?



YahooWeb:  
(a) In-degree vs. Out-degree



(b) Degree vs. Triangles

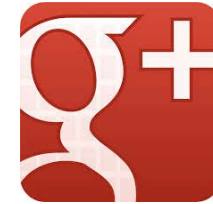


(c) Degree vs. PageRank

~1B nodes (web sites)  
~6B edges (http links)  
'YahooWeb graph'



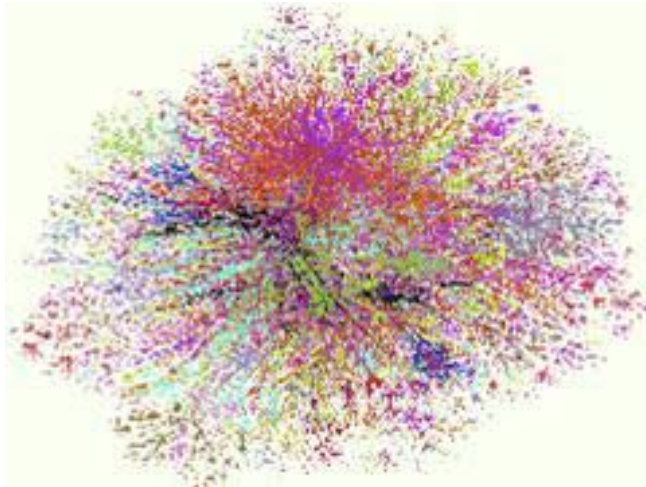
# Graphs - why should we care?



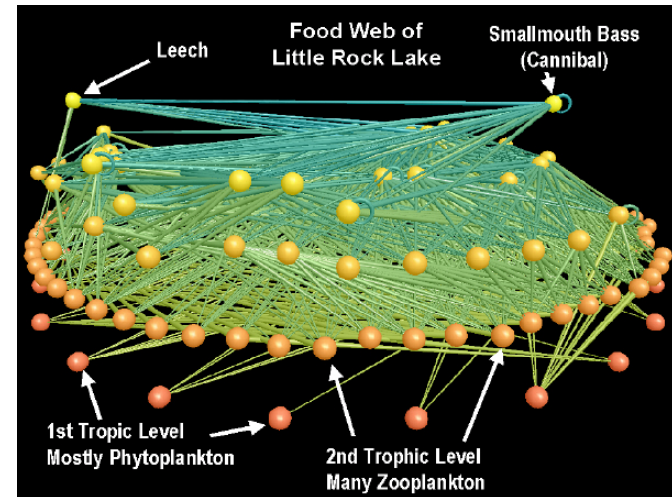
>\$10B; ~1B users



# Graphs - why should we care?



Internet Map  
[lumeta.com]

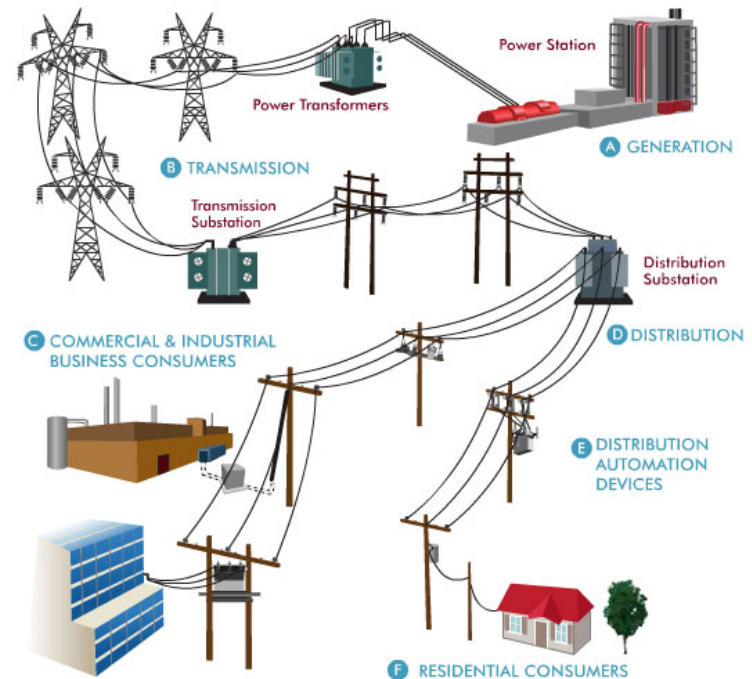


Food Web  
[Martinez '91]





# Graphs - why should we care?

- Power-grid!
  - Nodes: (plants/consumers)
  - Edges: power lines



# Graphs - why should we care?

- web-log ('blog') news propagation 
- computer network security: email/IP traffic and anomaly detection
- Recommendation systems 
- ....
- Many-to-many db relationship -> graph

# Motivating problems

- P1: patterns? Fraud detection?



- P2: patterns in time-evolving graphs / tensors

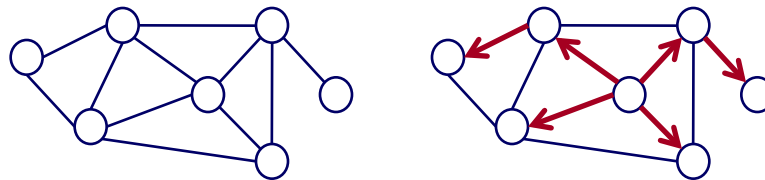
destination



source

time

- P3: cascades – whom to immunize?



# Roadmap

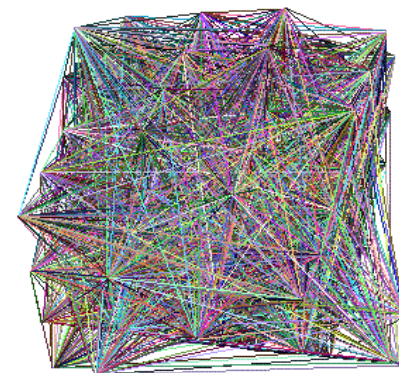


- Introduction – Motivation
  - Why study (big) graphs?
- ➔ • Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
- Part#3: Cascades and immunization
- Conclusions



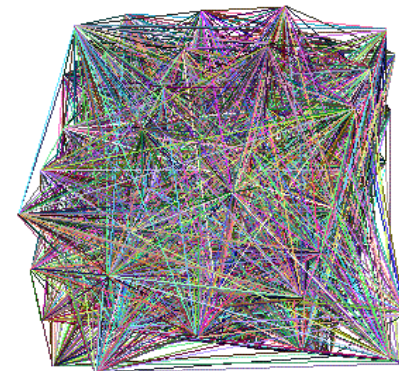
# Part 1: Patterns, & fraud detection

# Laws and patterns

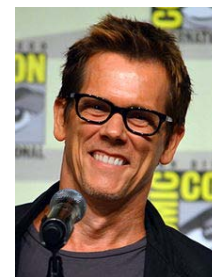


- Q1: Are real graphs random?

# Laws and patterns



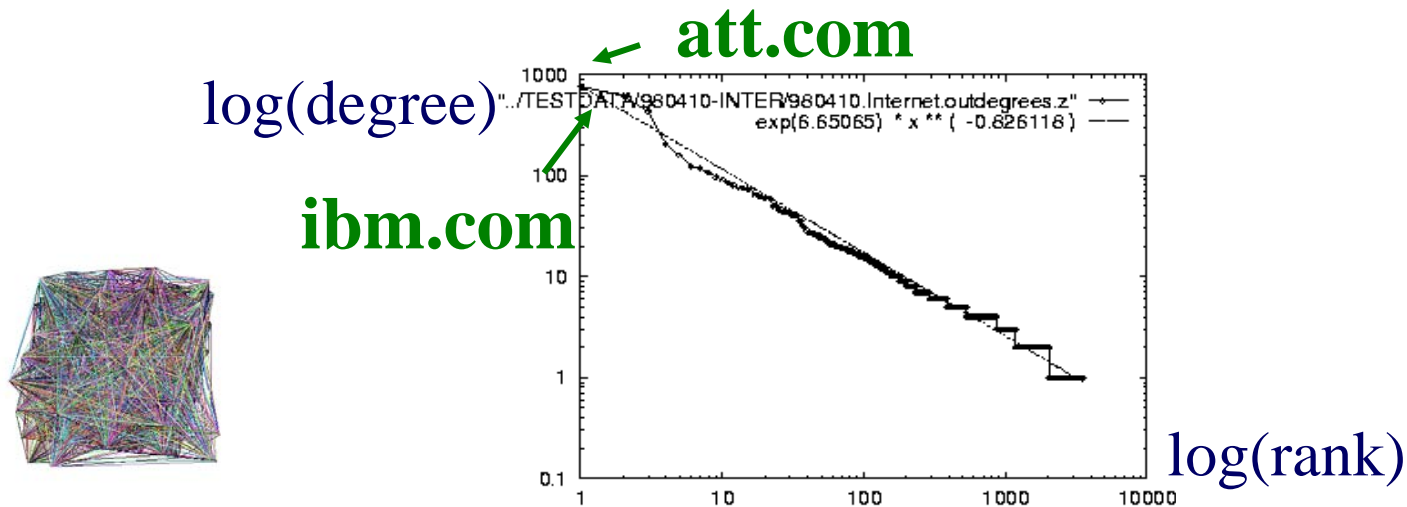
- Q1: Are real graphs random?
- A1: NO!!
  - Diameter (‘6 degrees’; ‘Kevin Bacon’)
  - in- and out- degree distributions
  - other (surprising) patterns
- So, let’s look at the data



# Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

internet domains

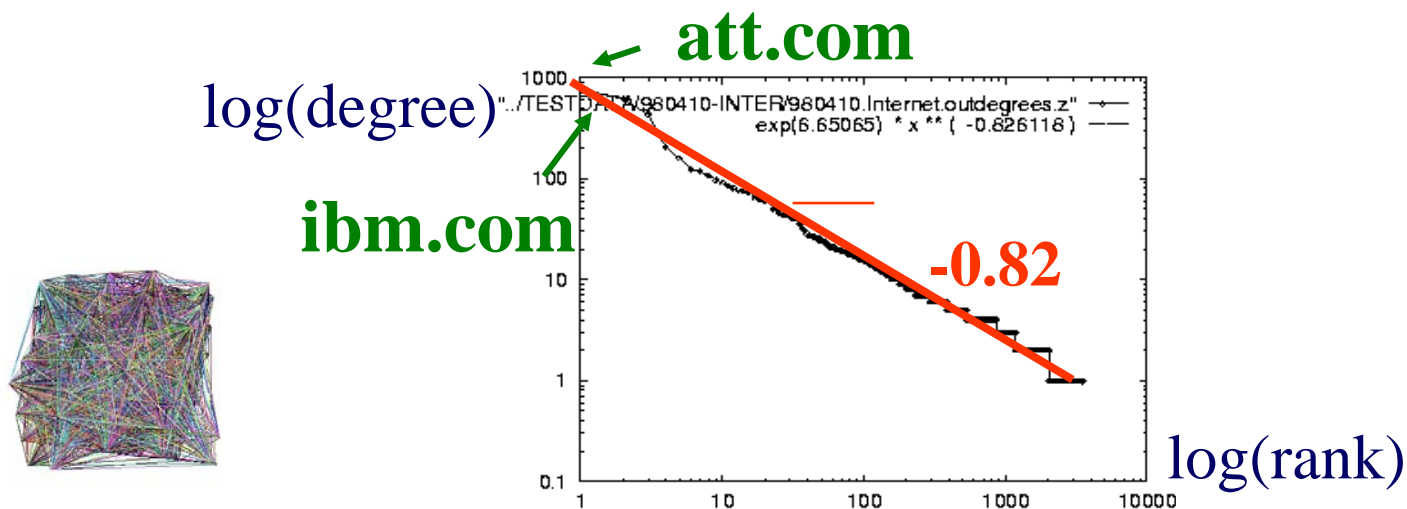




# Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

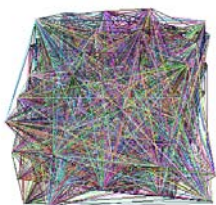
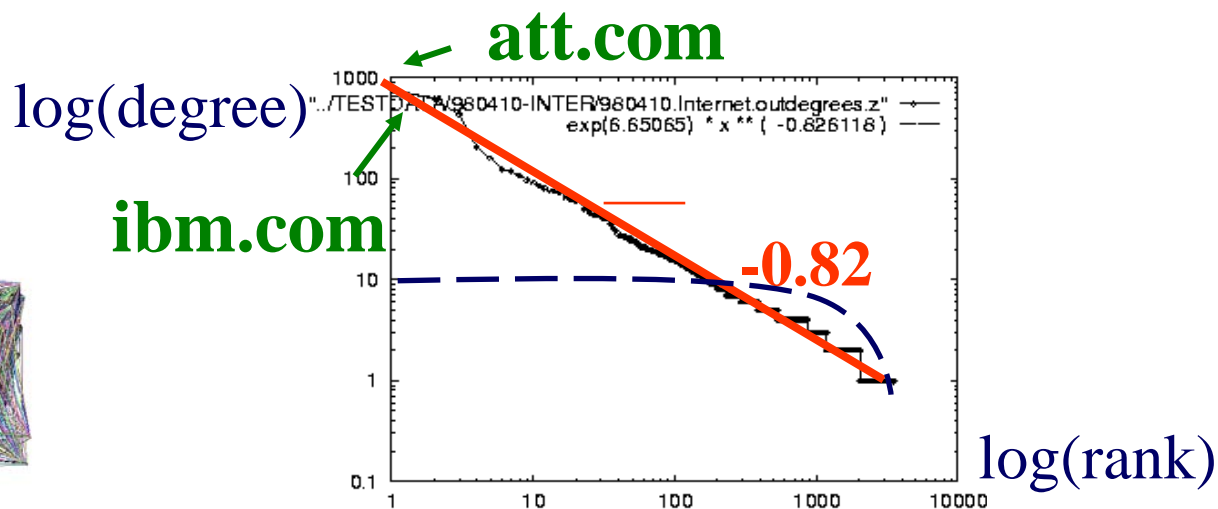
internet domains



# Solution# S.1

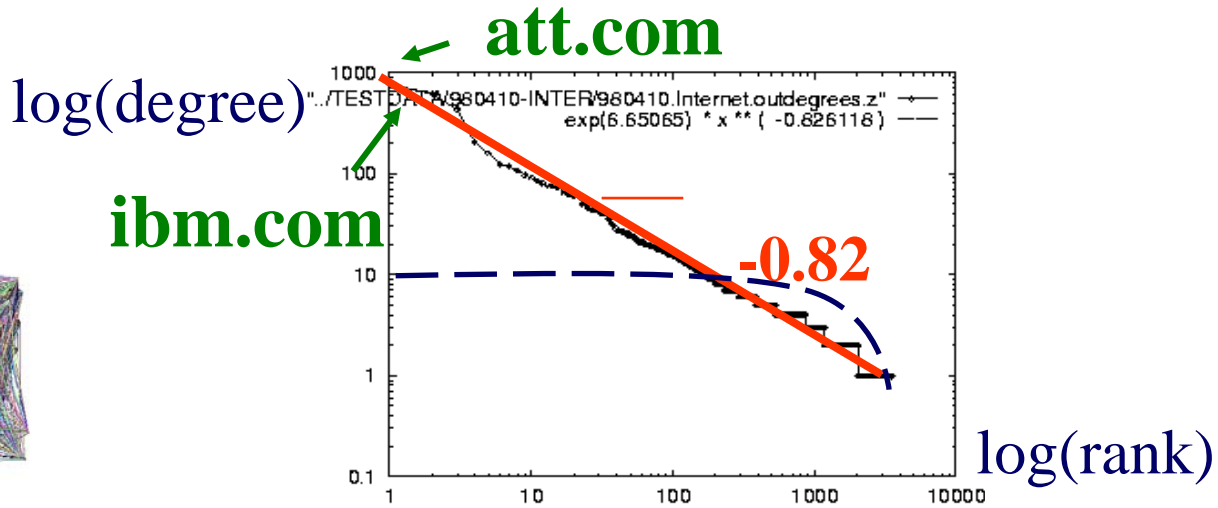
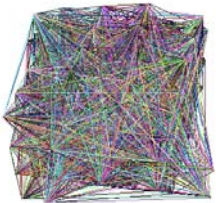
- Q: So what?

internet domains



# Solution# S.1

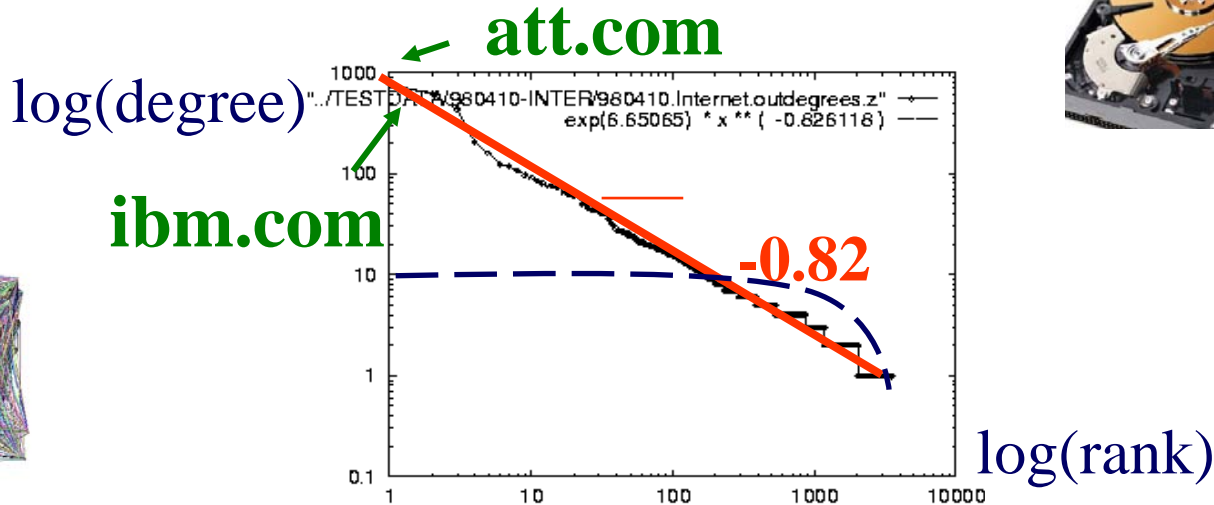
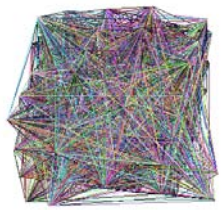
- Q: So what?
- A1: # of two-step-away pairs: **internet domains**  
= friends of friends (F.O.F.)



# Solution# S.1

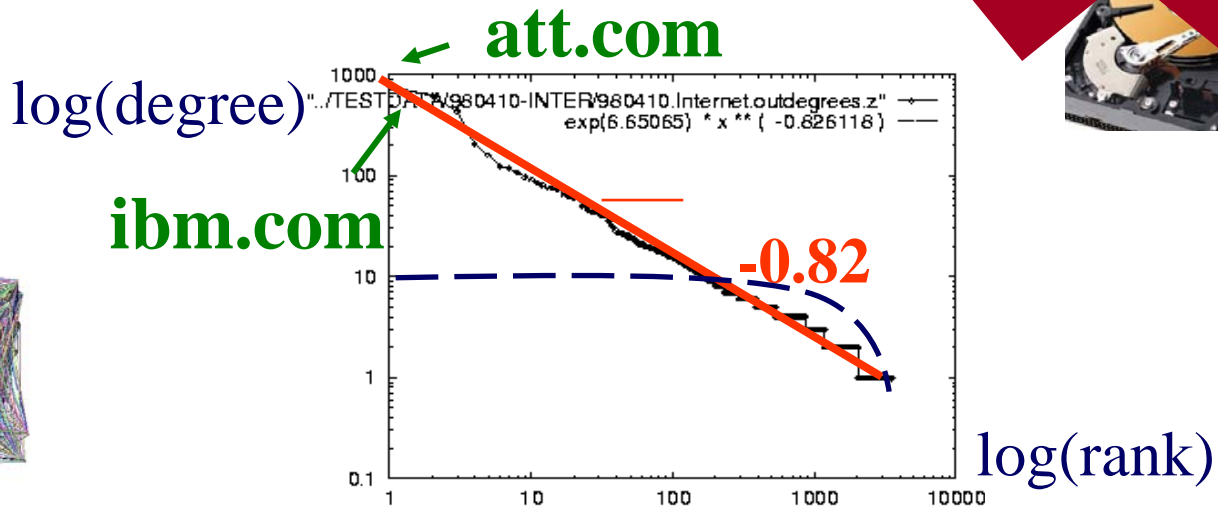
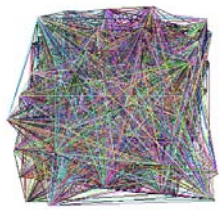
- Q: So what?
- A1: # of two-step-away pairs:  $100^2 * N = 10$  Trillion internet domains

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



# Solution# S.1

- Q: So what?
- A1: # of two-step-away pairs:  $100^2 \times 100^2 = 10^8$  Trillion internet domains



# Gaussian trap

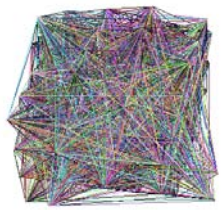
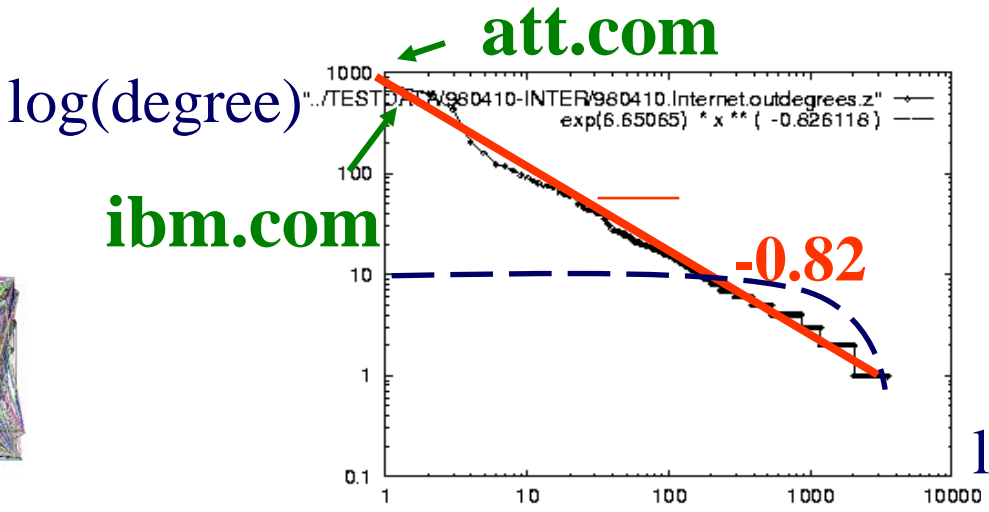
## Solution# S.1



- Q: So what? = friends of friends (F.O.F.)
- A1: # of two-step-away pairs:  $O(d_{max}^2) \sim 10M^2$  internet domains



~0.8PB -> a data center(!)



## Solution# S.1



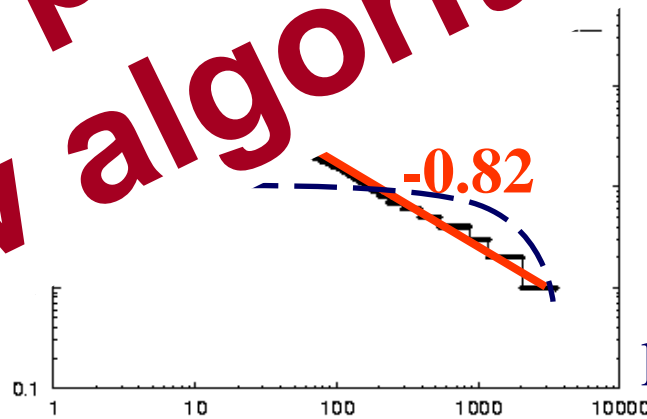
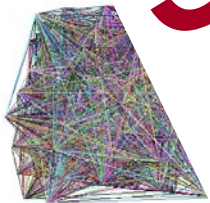
- Q: So what?
- A1: # of two-step-aww  
inter

? ) ~ 10M<sup>2</sup>



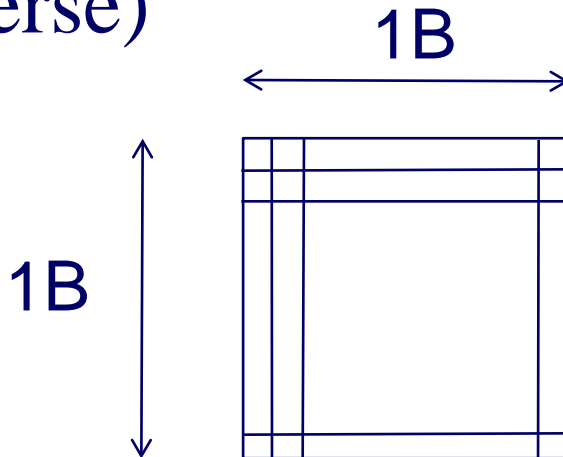
~0.8PB ->  
a data center(!)

**Such patterns ->  
New algorithms**



# Observation – big-data:

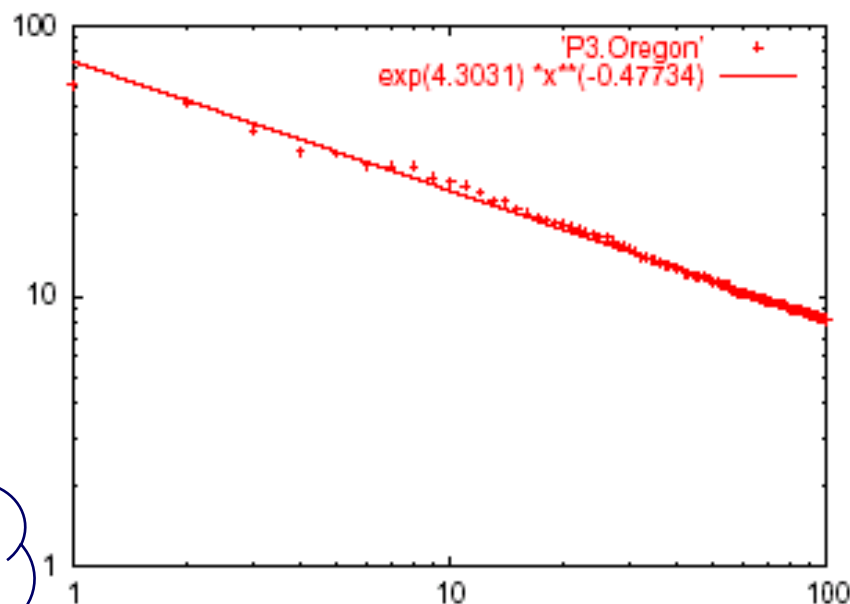
- $O(N^2)$  algorithms are ~intractable -  $N=1B$
- $N^2$  seconds = 31B years ( $>2x$  age of universe)





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

Eigenvalue



Exponent = slope

$$E = -0.48$$

May 2001

$$\mathbf{A} \mathbf{x} = \lambda \mathbf{x}$$

Rank of decreasing eigenvalue

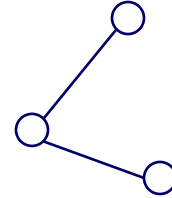
- A2: power law in the eigenvalues of the adjacency matrix ('eig()')

# Roadmap



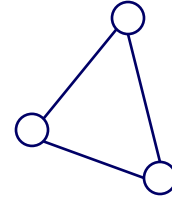
- Introduction – Motivation
- Part#1: Patterns in graphs
  - ➔ – Patterns: Degree; Triangles
  - Anomaly/fraud detection
  - Graph understanding
- Part#2: time-evolving graphs; tensors
- Part#3: Cascades and immunization
- Conclusions

# Solution# S.3: Triangle ‘Laws’



- Real social networks have a lot of triangles

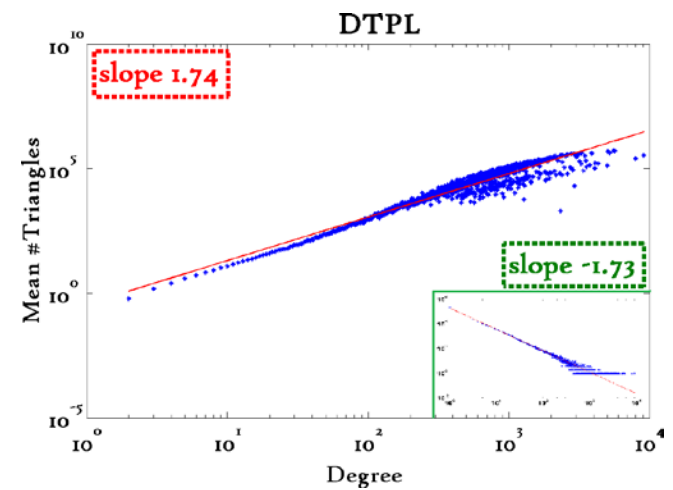
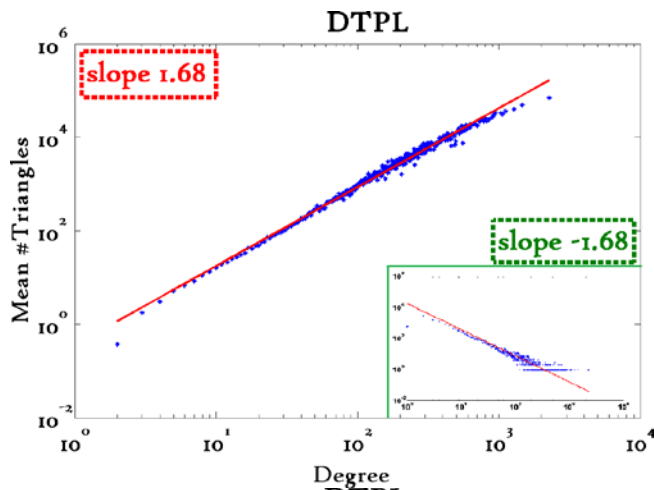
# Solution# S.3: Triangle ‘Laws’



- Real social networks have a lot of triangles
  - Friends of friends are friends
- Any patterns?
  - 2x the friends, 2x the triangles ?

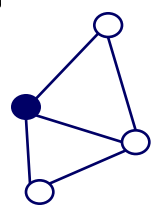
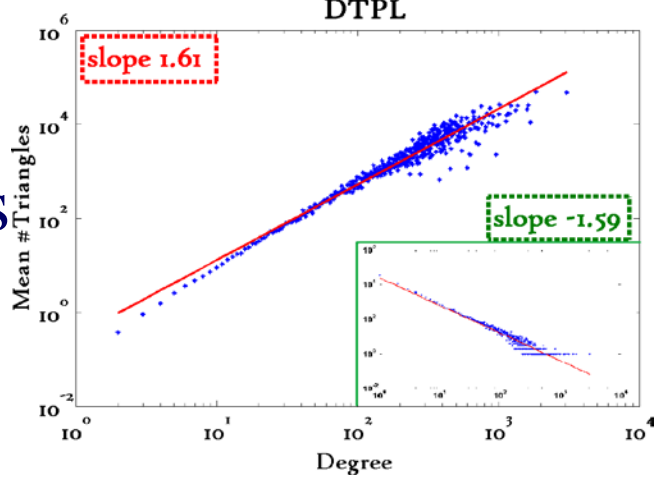
# Triangle Law: #S.3 [Tsourakakis ICDM 2008]

Reuters



SN

Epinions



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

# Triangle Law: Computations

[Tsourakakis ICDM 2008]



But: triangles are expensive to compute

(3-way join; several approx. algos) –  $O(d_{\max}^2)$

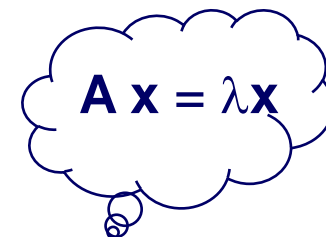
Q: Can we do that quickly?

A: Yes!

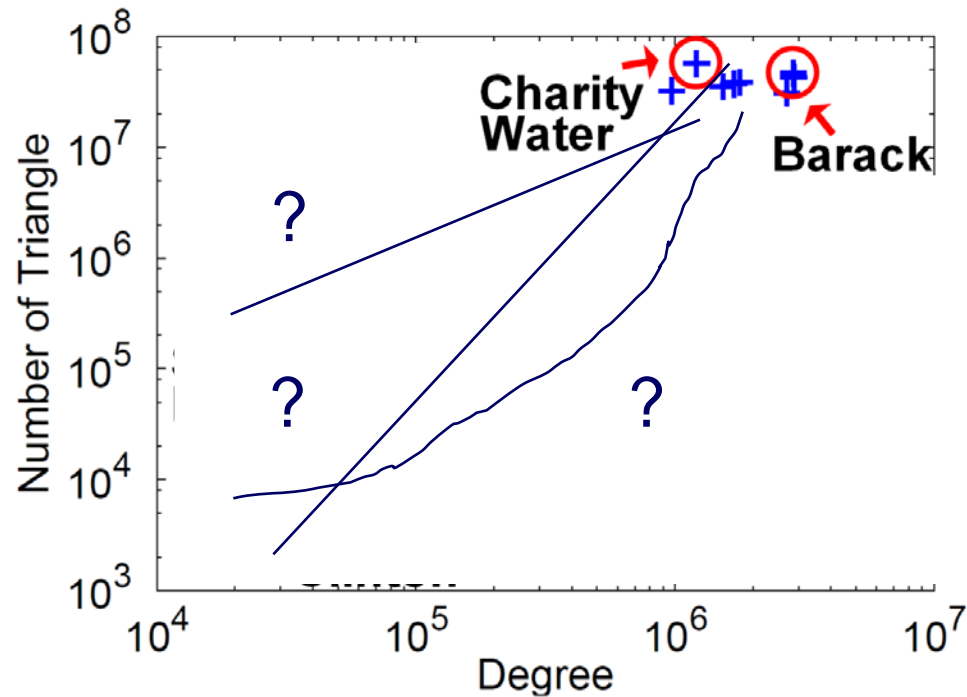
**#triangles = 1/6 Sum (  $\lambda_i^3$  )**

(and, because of skewness (S2) ,

we only need the top few eigenvalues! -  $O(E)$



# Triangle counting for large graphs?



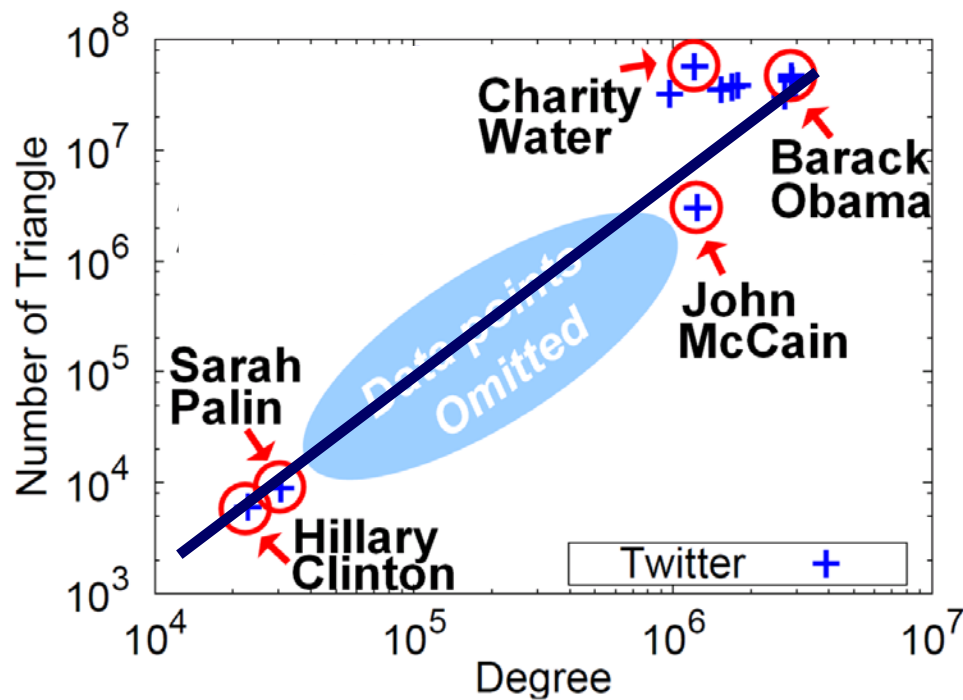
Anomalous nodes in Twitter (~ 3 billion edges)

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





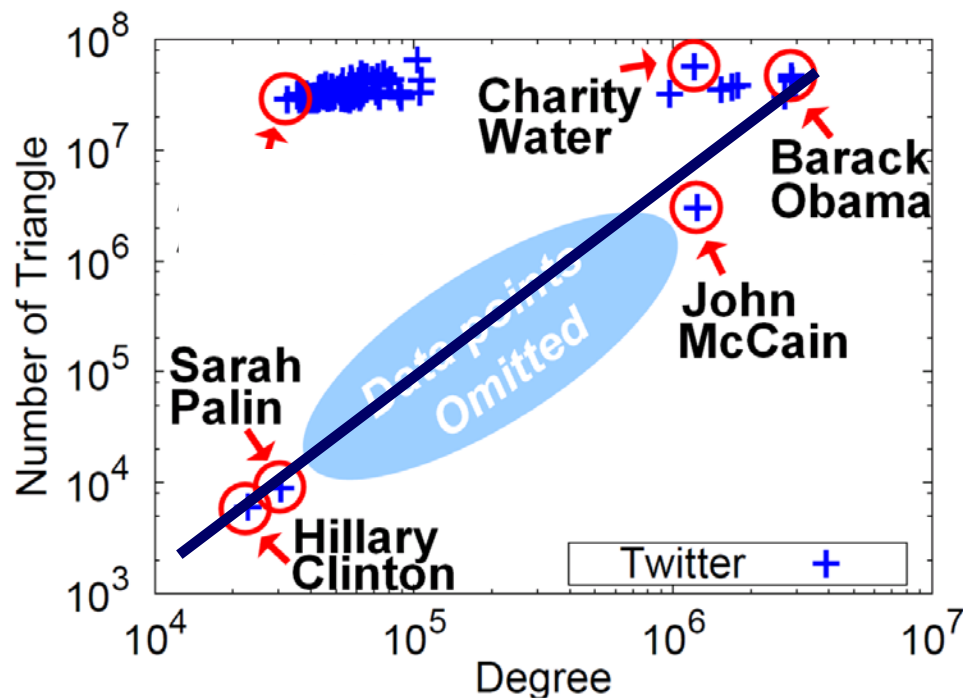
# Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

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

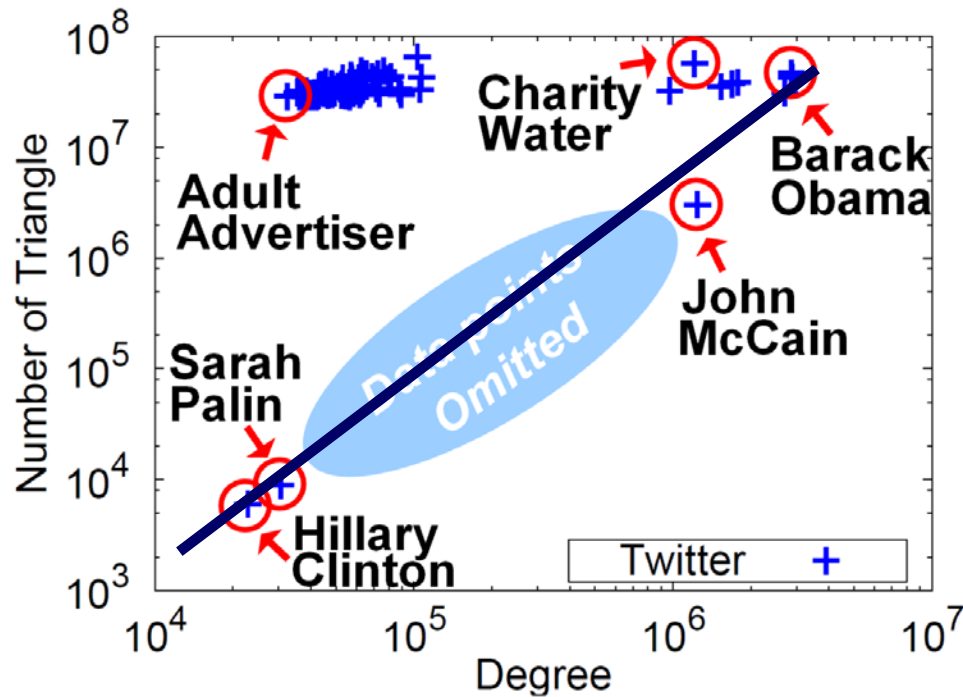
# Triangle counting for large graphs?



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

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

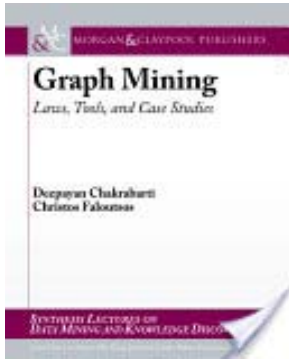
# MORE Graph Patterns

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

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



- Deepayan Chakrabarti and Christos Faloutsos, [Graph Mining: Laws, Tools, and Case Studies](#) Oct. 2012, Morgan Claypool.



# Roadmap



- Introduction – Motivation
- Part#1: Patterns in graphs
  - Patterns
  - ➔ – Anomaly / fraud detection
  - Graph understanding
- Part#2: time-evolving graphs; tensors
- Part#3: Cascades and immunization
- Conclusions

# Fraud

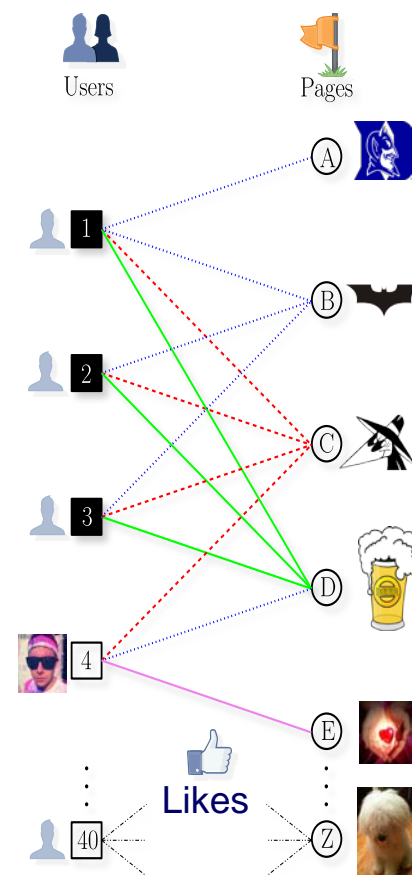
- Given
  - Who ‘likes’ what page, and when
- Find
  - Suspicious users and suspicious products



**CopyCatch: Stopping Group Attacks by Spotting Lockstep Behavior in Social Networks**, Alex Beutel, Wanhong Xu, Venkatesan Guruswami, Christopher Palow, Christos Faloutsos *WWW, 2013*.

# Fraud

- Given
  - Who ‘likes’ what page, and when
- Find
  - Suspicious users and suspicious products

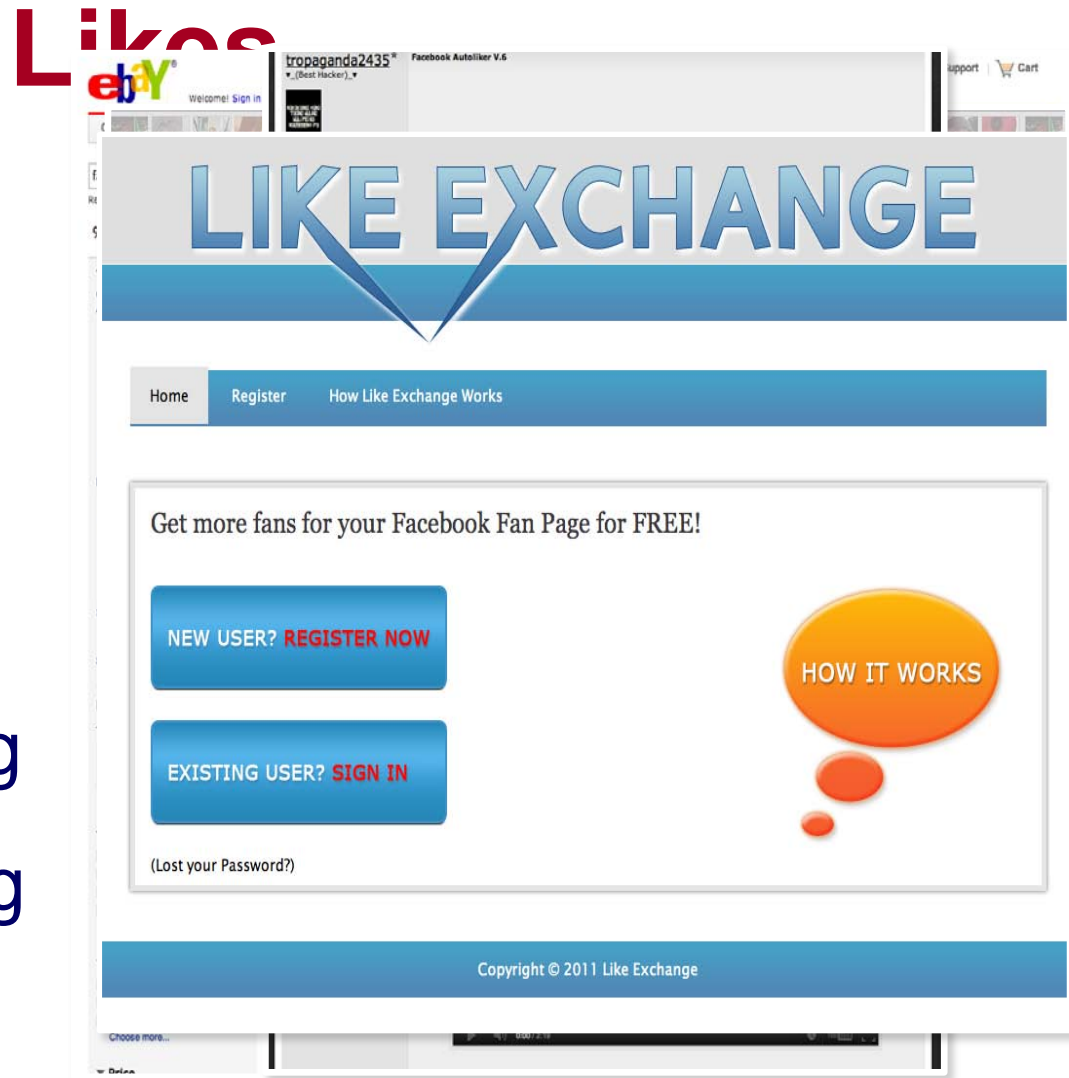


**CopyCatch: Stopping Group Attacks by Spotting Lockstep Behavior in Social Networks**, Alex Beutel, Wanhong Xu, Venkatesan Guruswami, Christopher Palow, Christos Faloutsos *WWW, 2013*.



# Ill-gotten Facebook Pages

- Popular Page = \$
- Fake 'likes' through unethical means:
  - Fake accounts
  - Malware
  - Credential stealing
  - Social Engineering

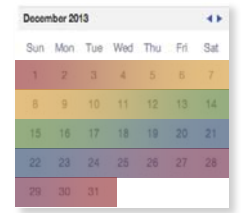
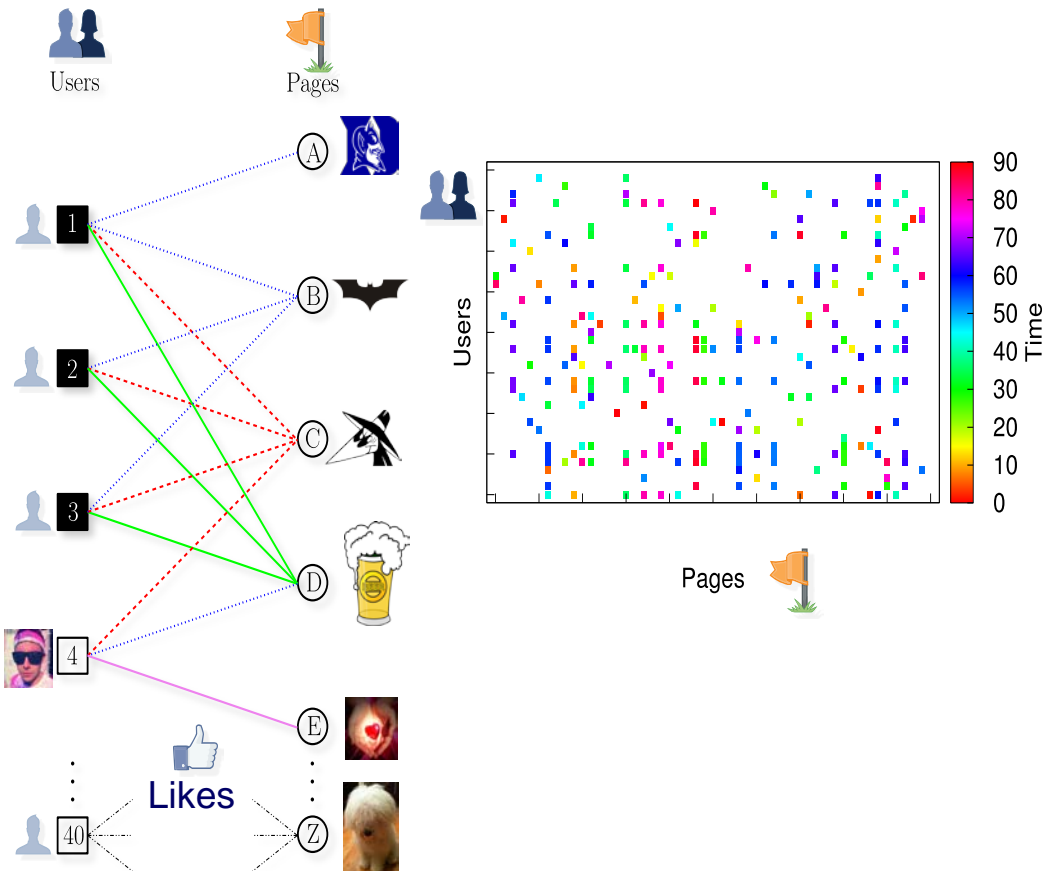


# Graph Patterns and Lockstep Behavior

Our intuition Behavior



- Lockstep behavior: Same Likes, same time



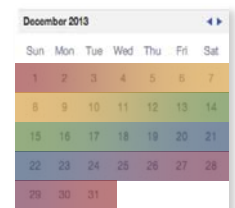
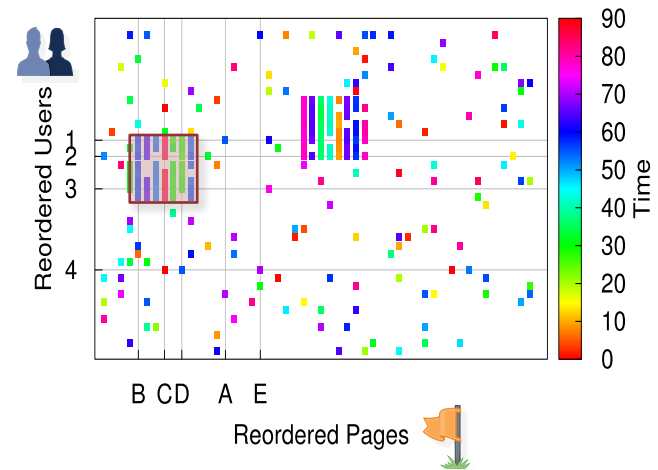
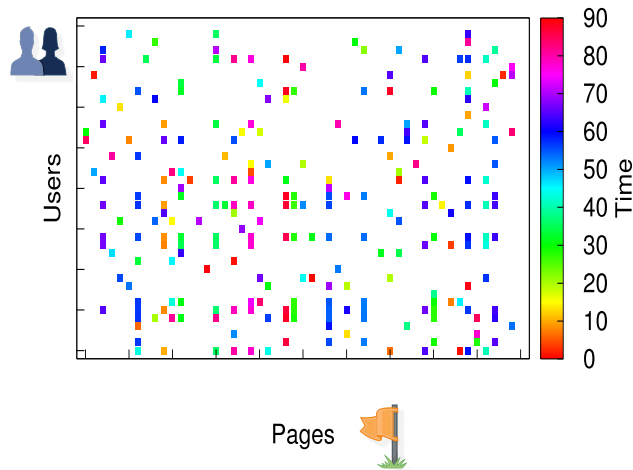
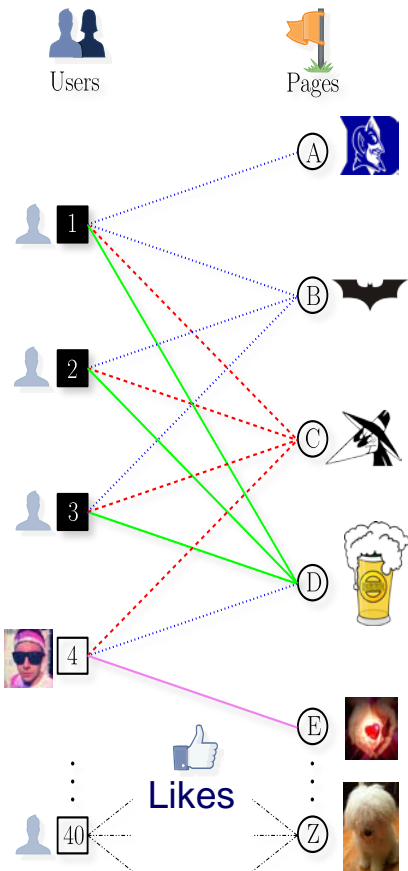
# Graph Patterns and Lockstep Behavior

Our intuition

## Behavior



- Lockstep behavior: Same Likes, same time



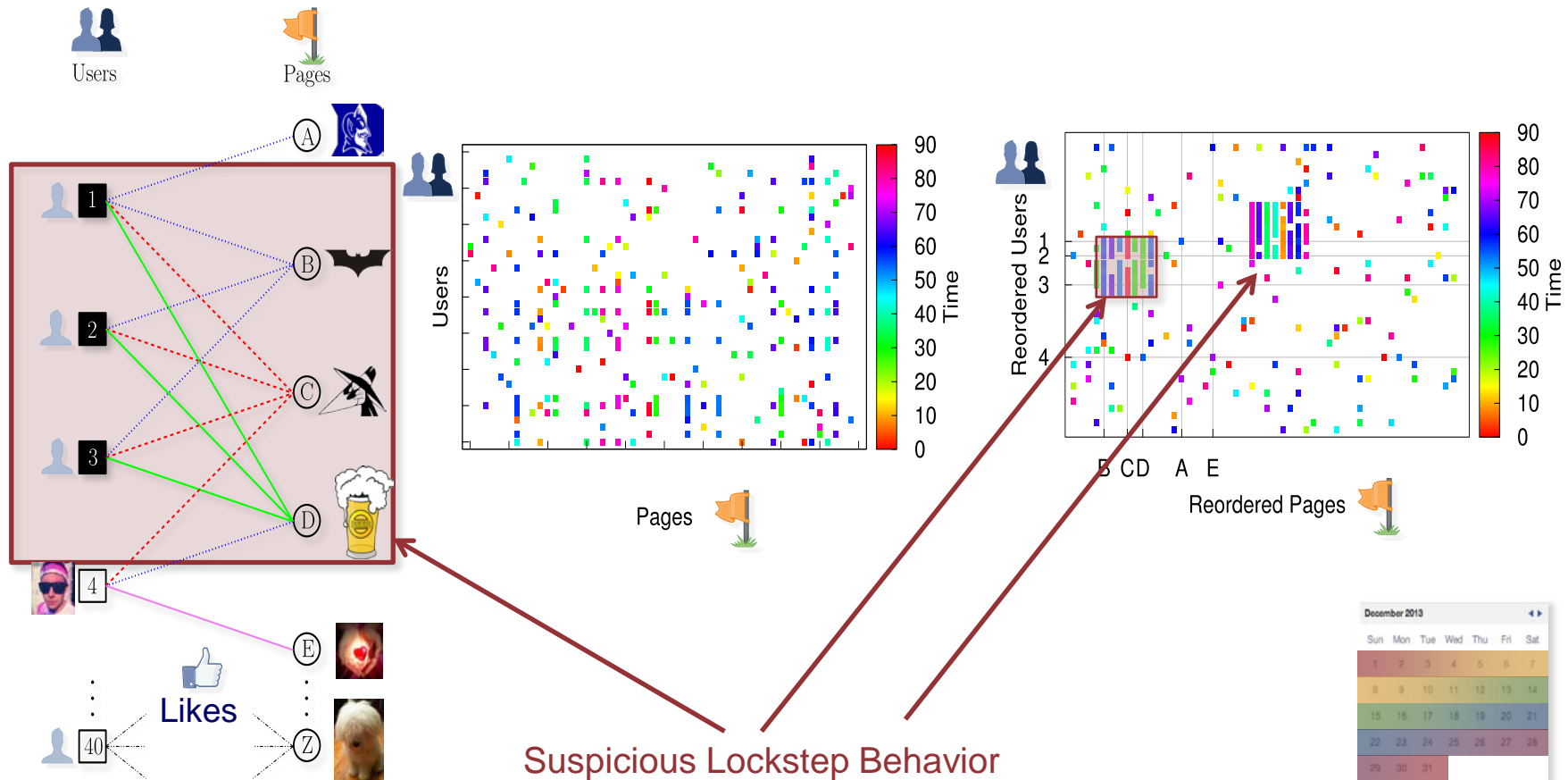
# Graph Patterns and Lockstep Behavior

Our intuition

## Behavior



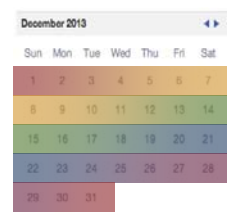
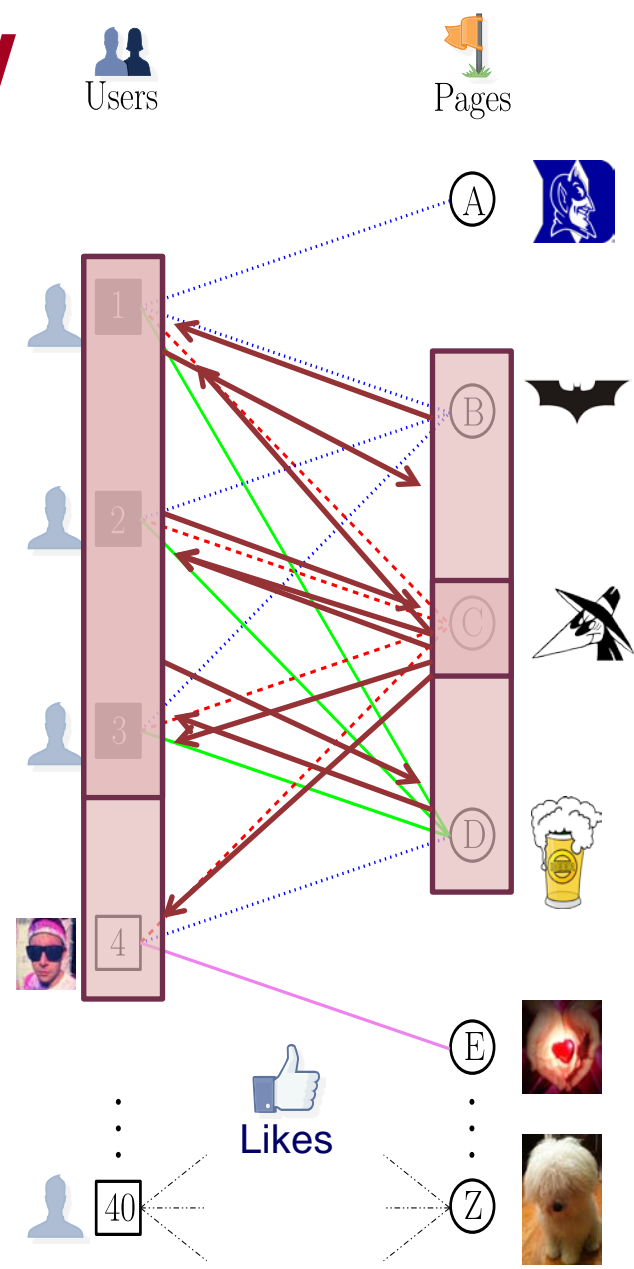
- Lockstep behavior: Same Likes, same time



Suspicious Lockstep Behavior

# MapReduce Overview

- Use Hadoop to search for many clusters in parallel:
  - Start with randomly seed
  - Update set of Pages and center Like times for each cluster
  - Repeat until convergence

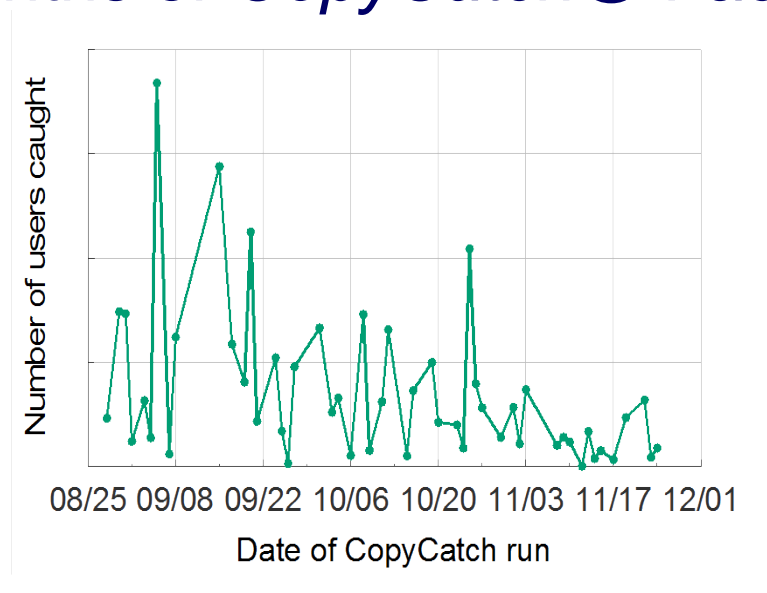


# Deployment at Facebook

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

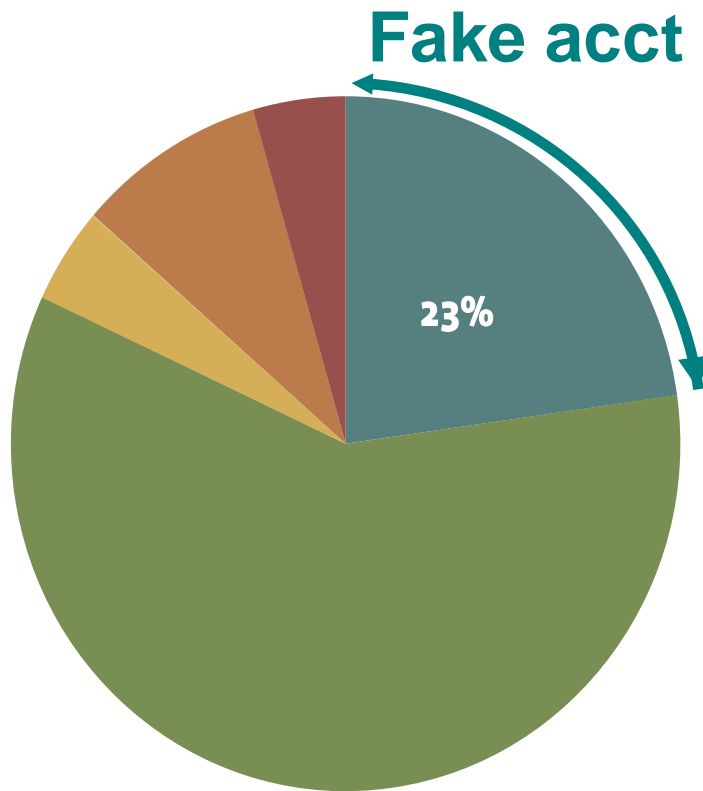
3 months of *CopyCatch*@ Facebook

#users  
caught



time

# Deployment at Facebook



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

Manually labeled 22 randomly selected *clusters* from February 2013

# Roadmap

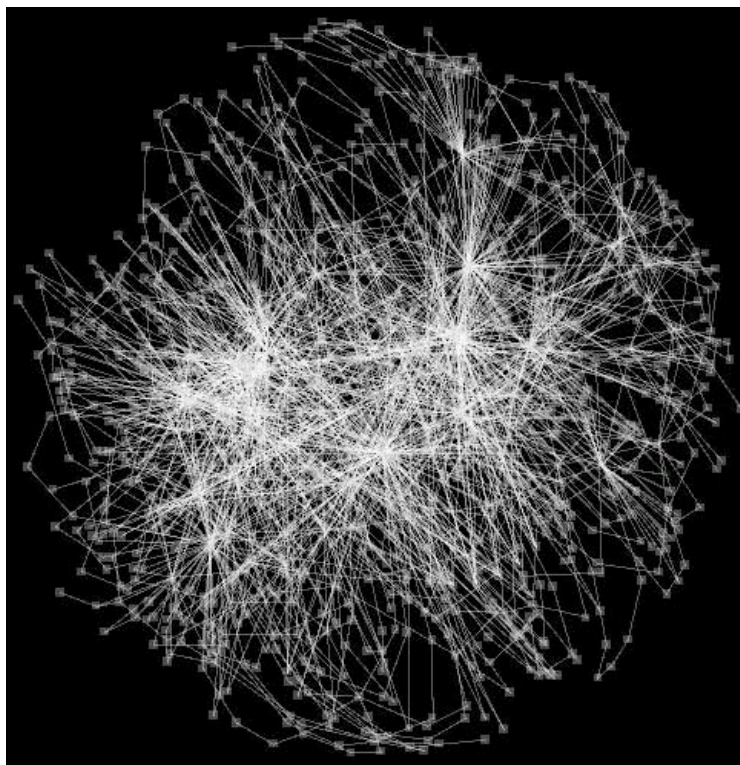


- Introduction – Motivation
- Part#1: Patterns in graphs
  - Patterns
  - Anomaly / fraud detection
  - ➔ – Graph understanding
- Part#2: time-evolving graphs; tensors
- Part#3: Cascades and immunization
- Conclusions





# Wikipedia - editors



- Nodes: editors
- Edge A->B: 'A' changed 'B'

**Any pattern?**

# *VoG: Summarizing and Understanding Large Graphs*

Danai Koutra,



U Kang,



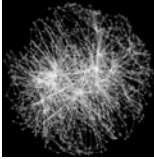
Jilles Vreeken,



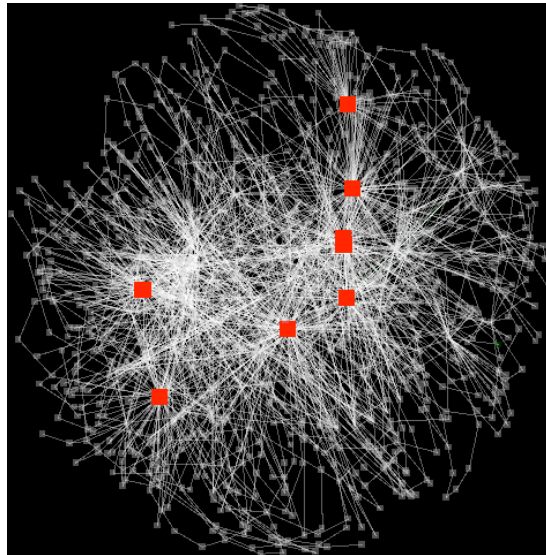
Christos Faloutsos.

*SDM 2014*, Philadelphia, PA, April 2014.

**Code:** [www.cs.cmu.edu/~dkoutra/CODE/vog.tar](http://www.cs.cmu.edu/~dkoutra/CODE/vog.tar)

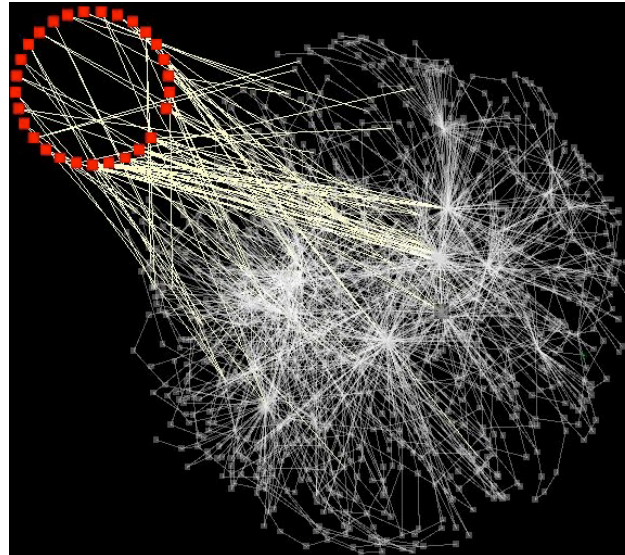


# VoG: Summarizing Wiki-controversy



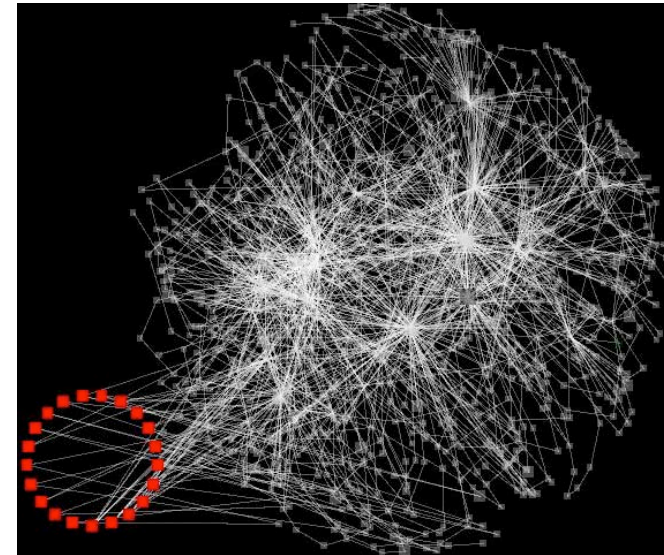
*top-8 star  
structures:  
admins, heavy  
wiki users, bots*

WWW, Seoul

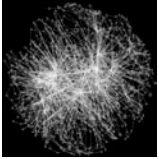


*warring factions  
changing each-  
other's edits.  
(Kiev vs Kiyv)*

(c) 2014, C. Faloutsos



*Ditto, between  
vandals*



# VoG: Summarizing Graphs using Rich Vocabularies

## Main Ideas:

(1) Use 'vocabulary' of subgraph types



(2) Minimum Description Length (MDL) and above vocabulary, to summarize graph

# Summary of Part#1

- \*many\* patterns in real graphs
  - Power-laws everywhere
  - Gaussian trap
    - $\text{Avg} \ll \text{Max}$



# Roadmap

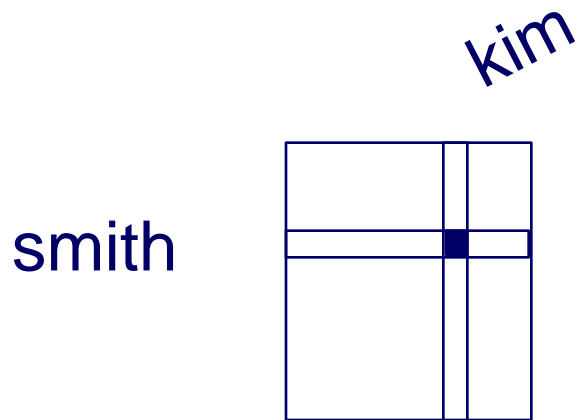


- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
  - ➔ – P2.1: time-evolving graphs
  - P2.2: with side information (‘coupled’ M.T.F.)
  - Speed
- Part#3: Cascades and immunization
- Conclusions

# Part 2: Time evolving graphs; tensors

# Graphs over time -> tensors!

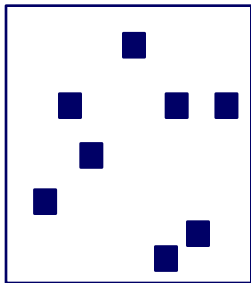
- Problem #2.1:
  - Given who calls whom, and when
  - Find patterns / anomalies





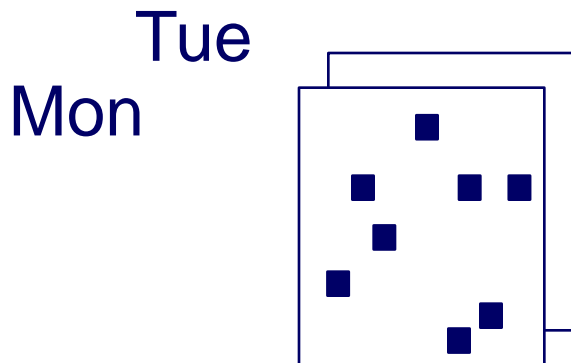
# Graphs over time -> tensors!

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



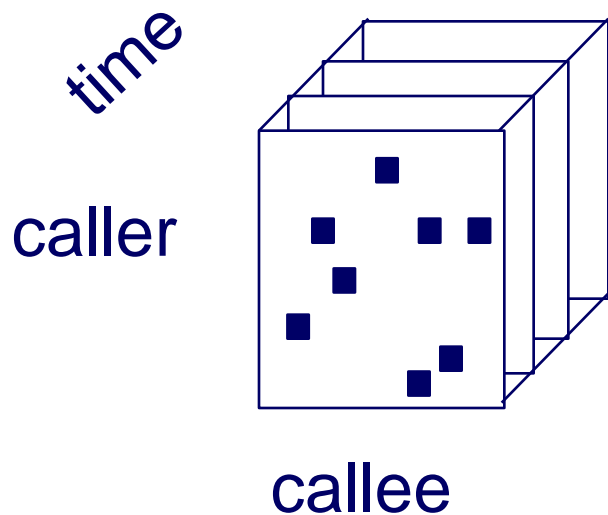
# Graphs over time -> tensors!

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



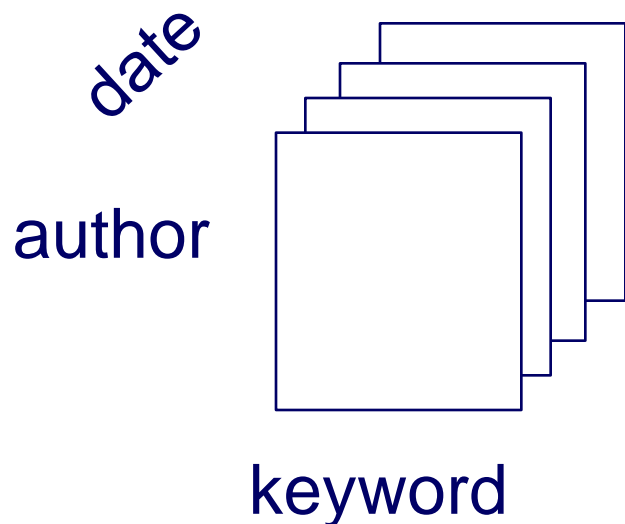
# Graphs over time -> tensors!

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



# Graphs over time -> tensors!

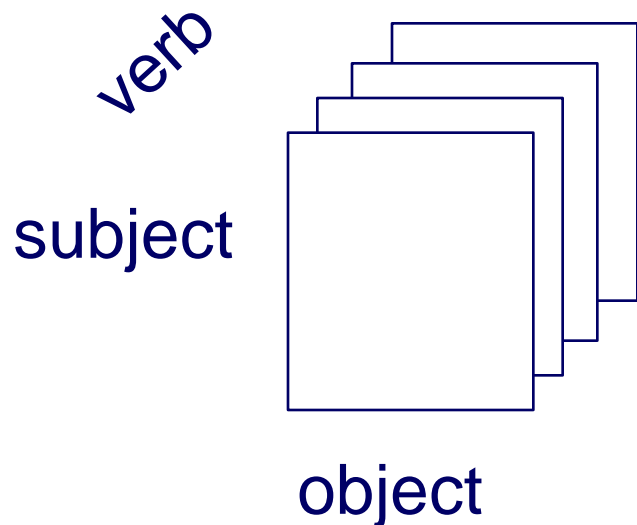
- Problem #2.1':
  - Given author-keyword-date
  - Find patterns / anomalies



**MANY** more settings,  
with  $>2$  'modes'

# Graphs over time -> tensors!

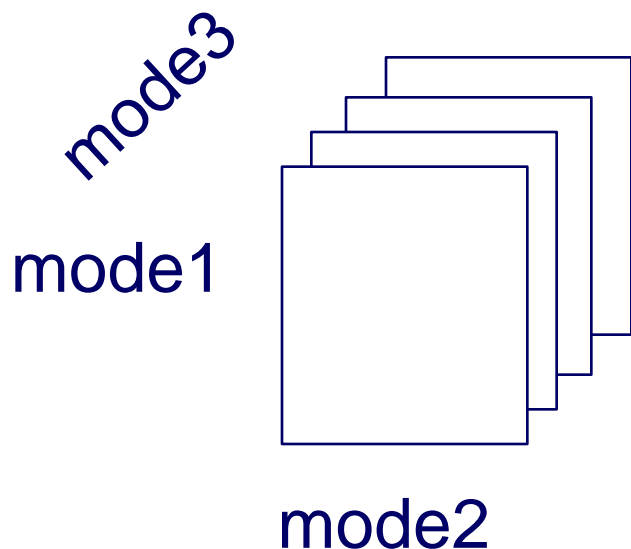
- Problem #2.1’’:
  - Given subject – verb – object facts
  - Find patterns / anomalies



**MANY** more settings,  
with  $>2$  ‘modes’

# Graphs over time -> tensors!

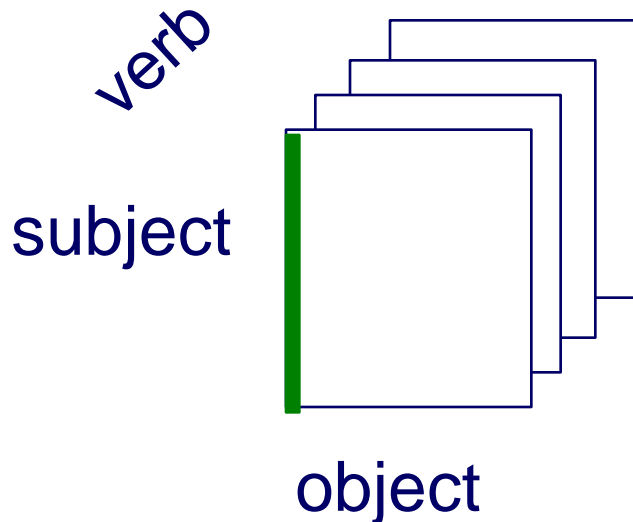
- Problem #2.1''':
  - Given <triplets>
  - Find patterns / anomalies



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

# Graphs & side info

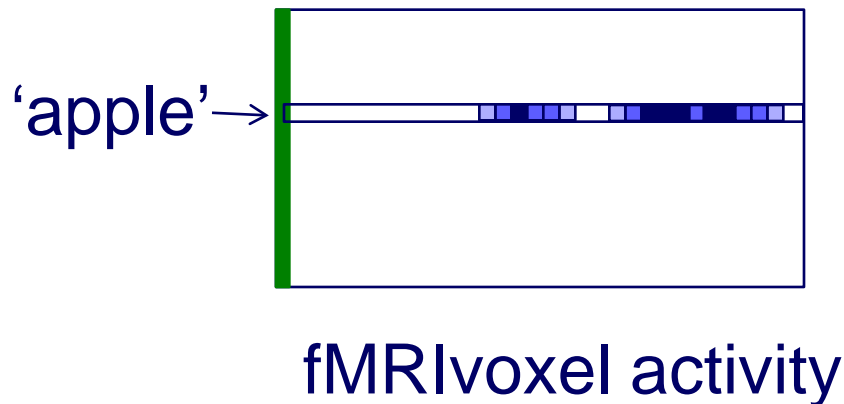
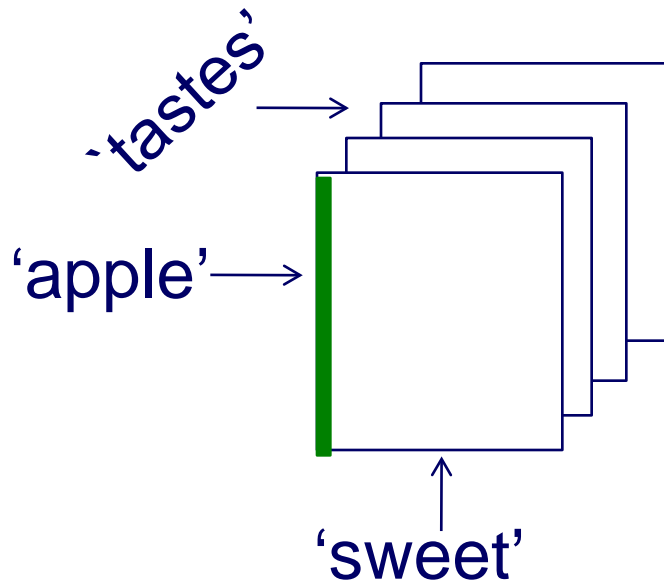
- Problem #2.2: coupled (eg., side info)
  - Given subject – verb – object facts
    - And voxel-activity for each subject-word
  - Find patterns / anomalies



fMRI voxel activity

# 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



'apple tastes sweet'



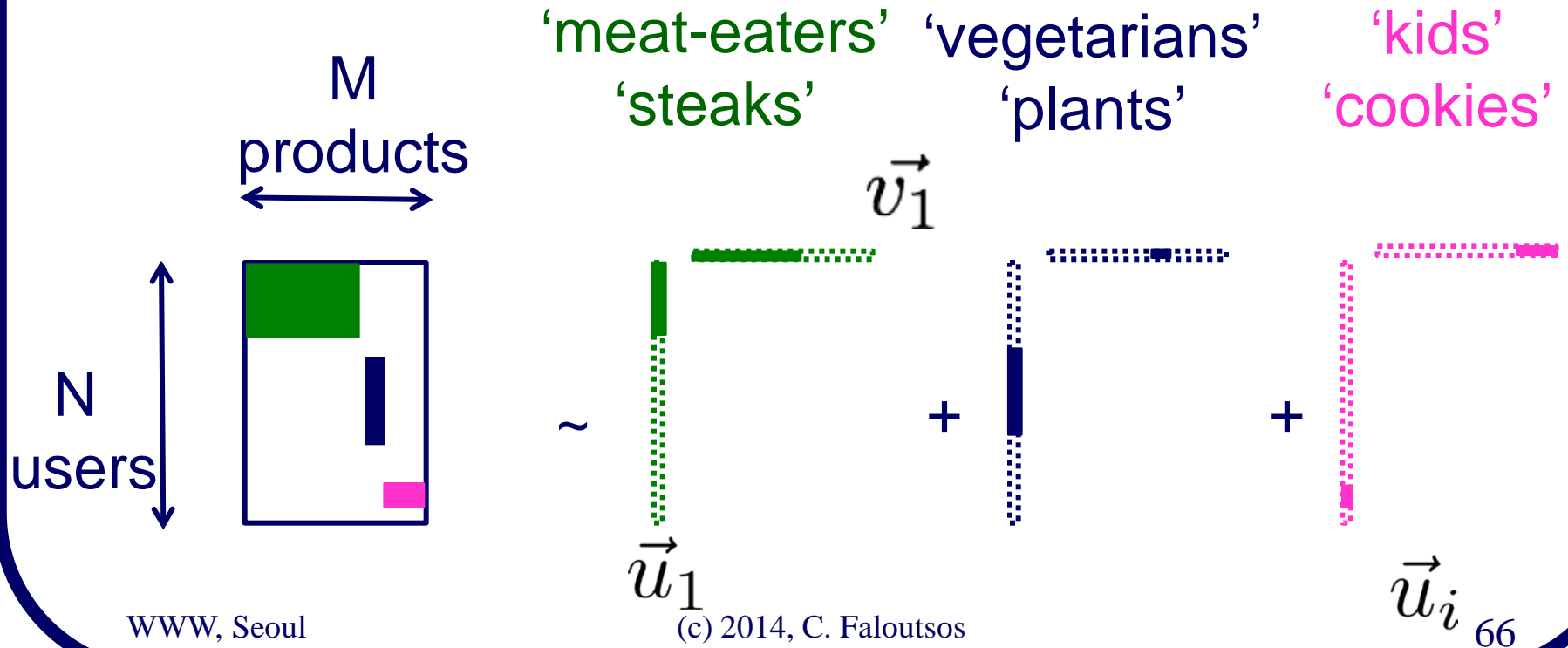
# Roadmap



- 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
- Part#3: Cascades and immunization
- Conclusions

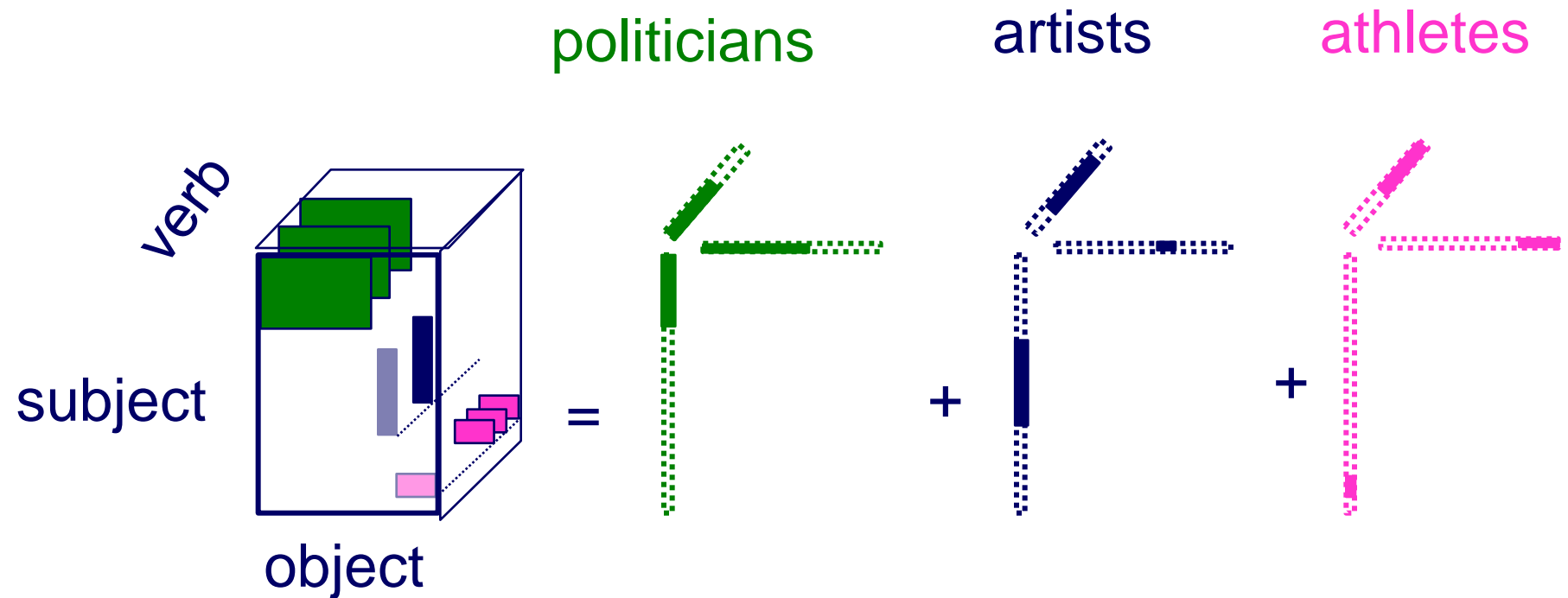
# Answer to both: tensor factorization

- Recall: (SVD) matrix factorization: finds blocks



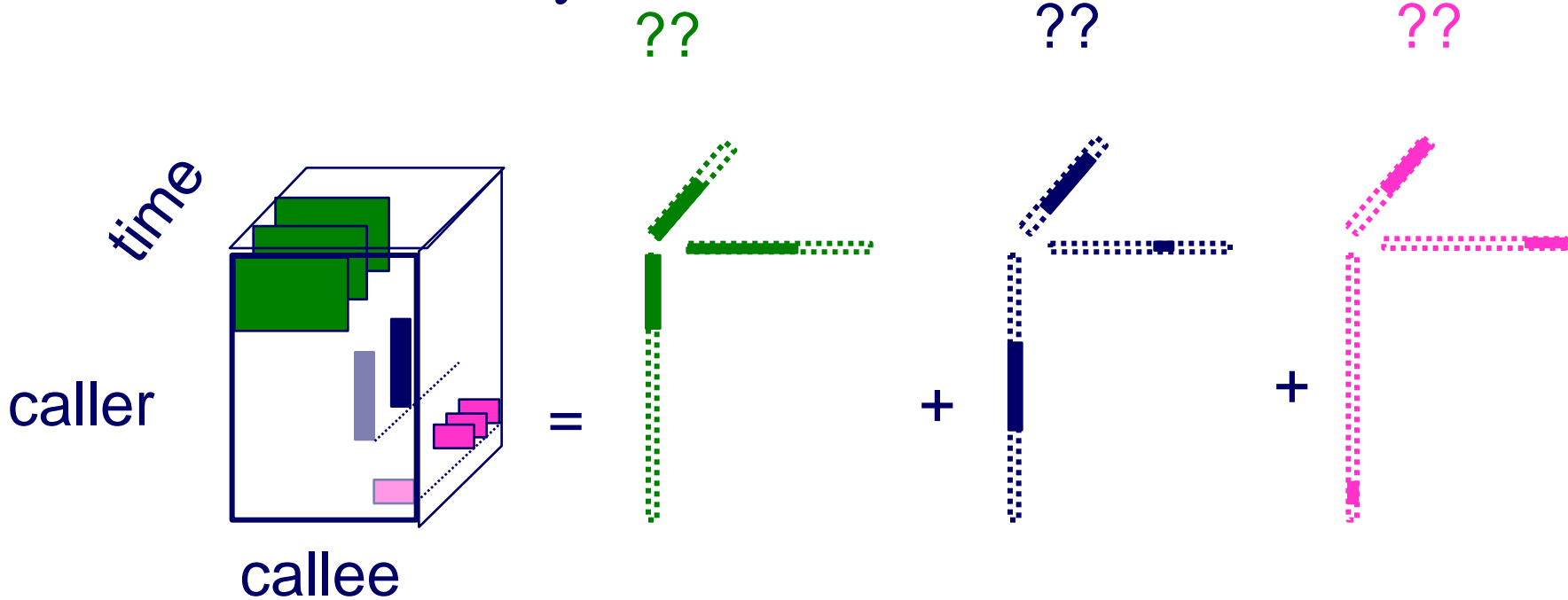
# Answer to both: tensor factorization

- PARAFAC decomposition

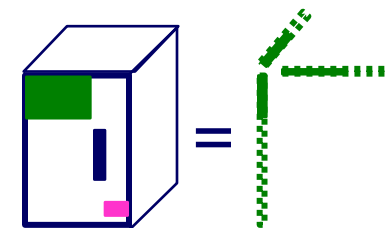


# Answer: tensor factorization

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

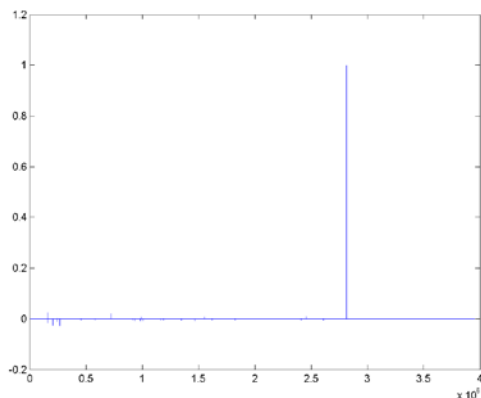


# Anomaly detection in time-evolving graphs

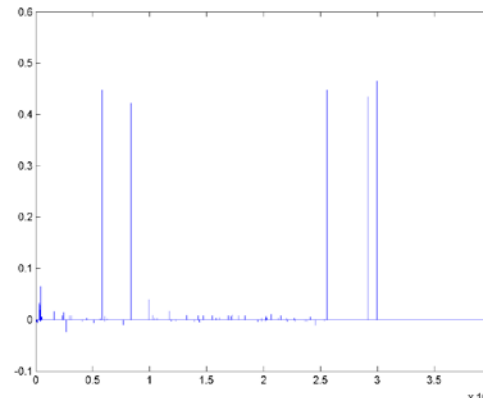


- Anomalous communities in phone call data:
  - European country, 4M clients, data over 2 weeks

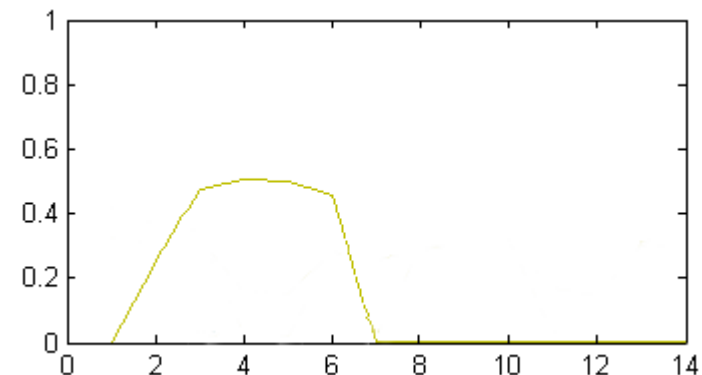
1 caller



5 receivers

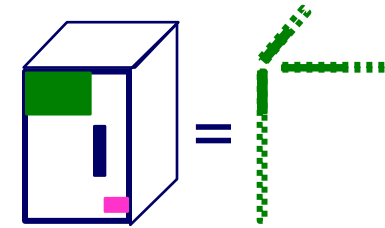


4 days of activity



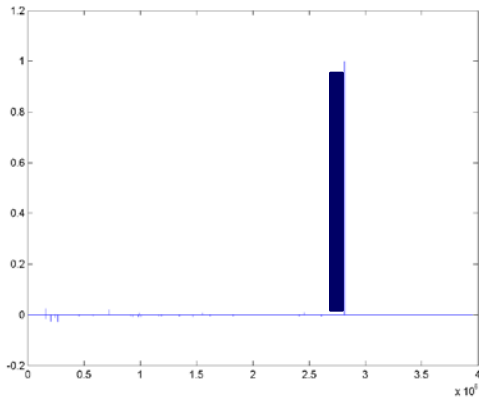
**~200 calls to EACH receiver on EACH day!**

# Anomaly detection in time-evolving graphs

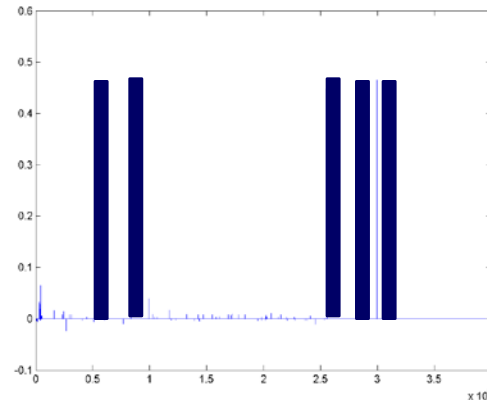


- Anomalous communities in phone call data:
  - European country, 4M clients, data over 2 weeks

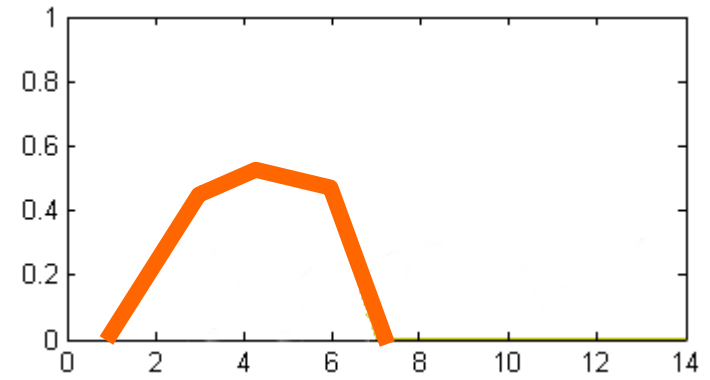
1 caller



5 receivers

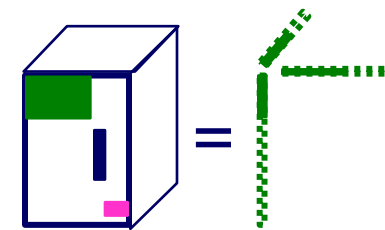


4 days of activity

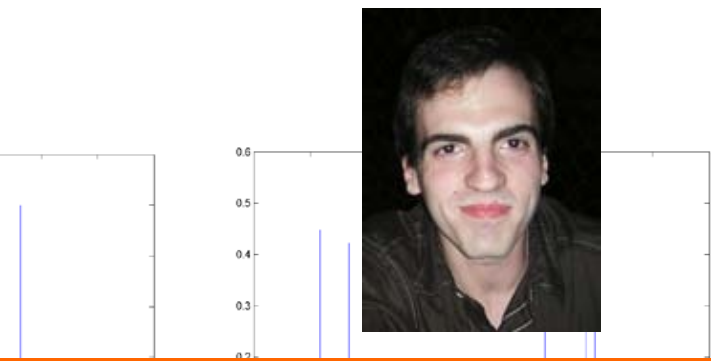
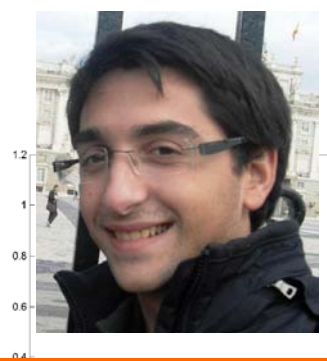


~200 calls to EACH receiver on EACH day!

# 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

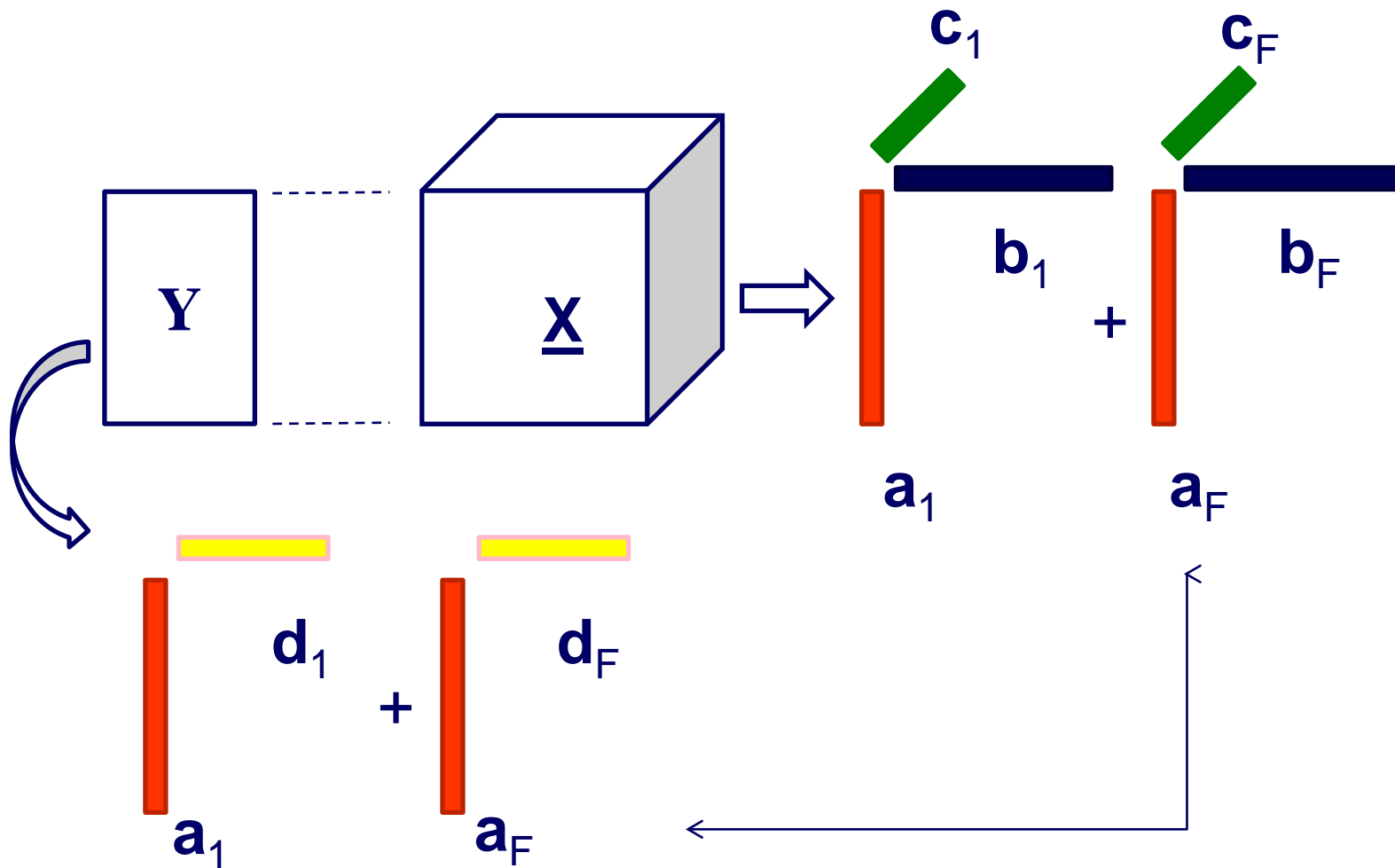


- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
  - P2.1: Discoveries @ phonecall network
  - P2.2: Discoveries in neuro-semantic
  - Speed
- Part#3: Cascades and immunization
- Conclusions





# Coupled Matrix-Tensor Factorization (CMTF)



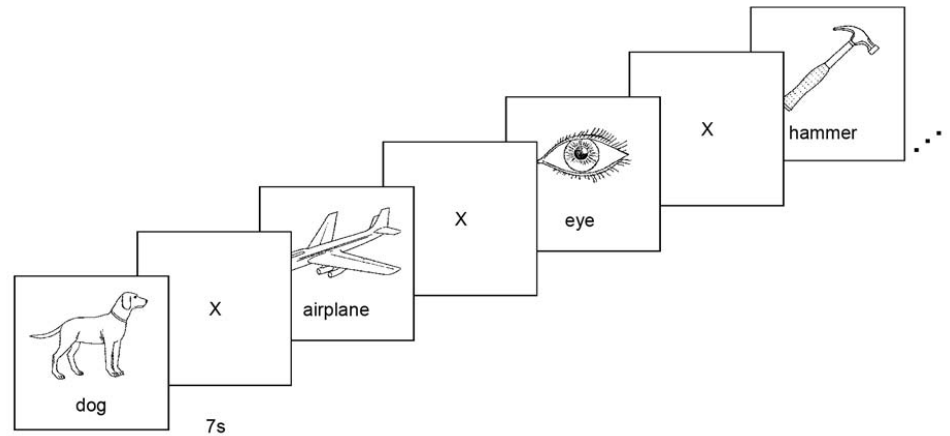
# Neuro-semantic

- **Brain Scan Data\***

- 9 persons
- 60 nouns

- **Questions**

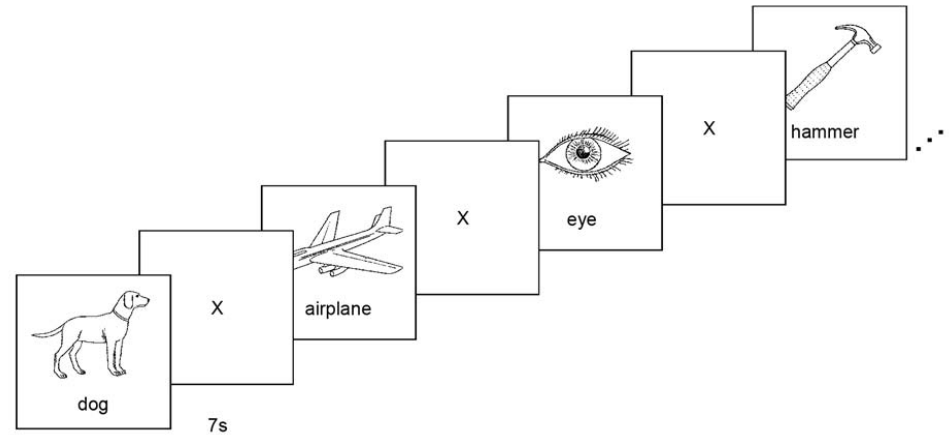
- 218 questions
- 'is it alive?', 'can you eat it?'



\*Mitchell et al. *Predicting human brain activity associated with the meanings of nouns*. Science, 2008.  
Data@ [www.cs.cmu.edu/afs/cs/project/theo-73/www/science2008/data.html](http://www.cs.cmu.edu/afs/cs/project/theo-73/www/science2008/data.html)

# Neuro-semantic

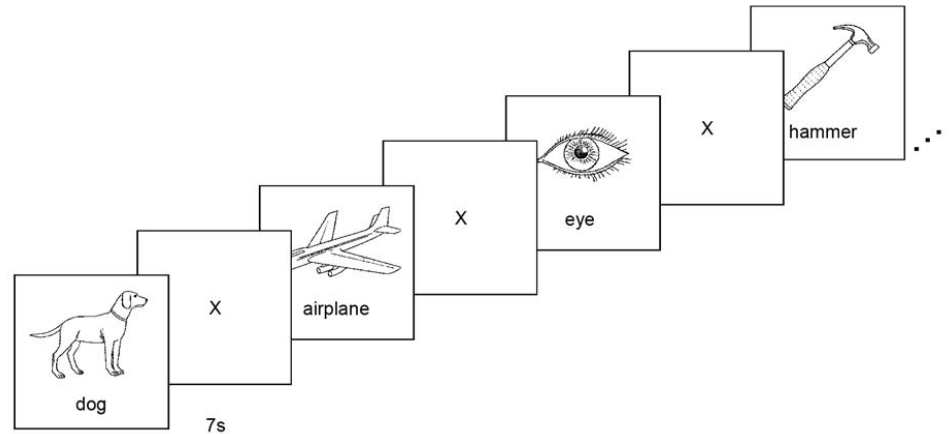
- **Brain Scan Data\***
  - 9 persons
  - 60 nouns
- **Questions**
  - 218 questions
  - ‘is it alive?’, ‘can you eat it?’



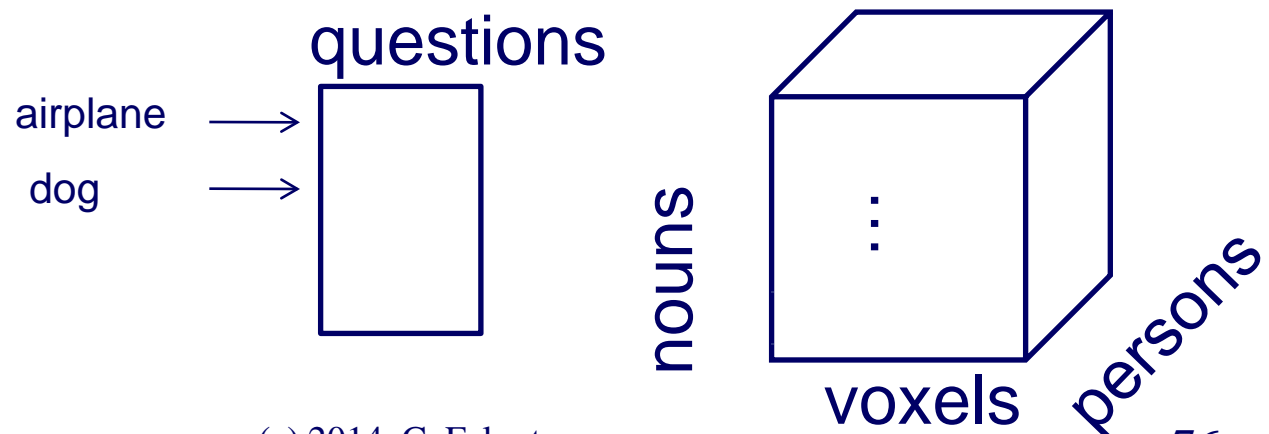
## Patterns?

# Neuro-semantic

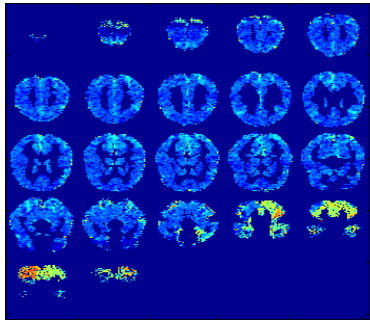
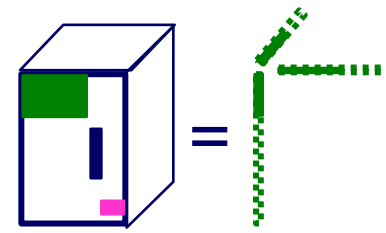
- **Brain Scan Data\***
  - 9 persons
  - 60 nouns
- **Questions**
  - 218 questions
  - 'is it alive?', 'can you eat it?'



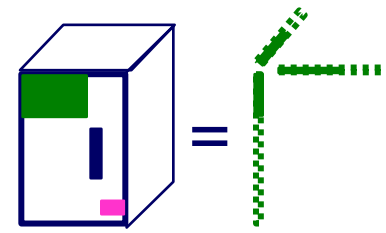
## Patterns?



# Neuro-semantic



# Neuro-semantic



**Small items ->  
Premotor cortex**

# Neuro-semantic

**Small items ->  
Premotor cortex**



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

# Roadmap



- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
  - P2.1: Discoveries @ phonecall network
  - P2.2: Discoveries in neuro-semantic
  - Speed
- ➔ • Part#3: Cascades and immunization
- Conclusions

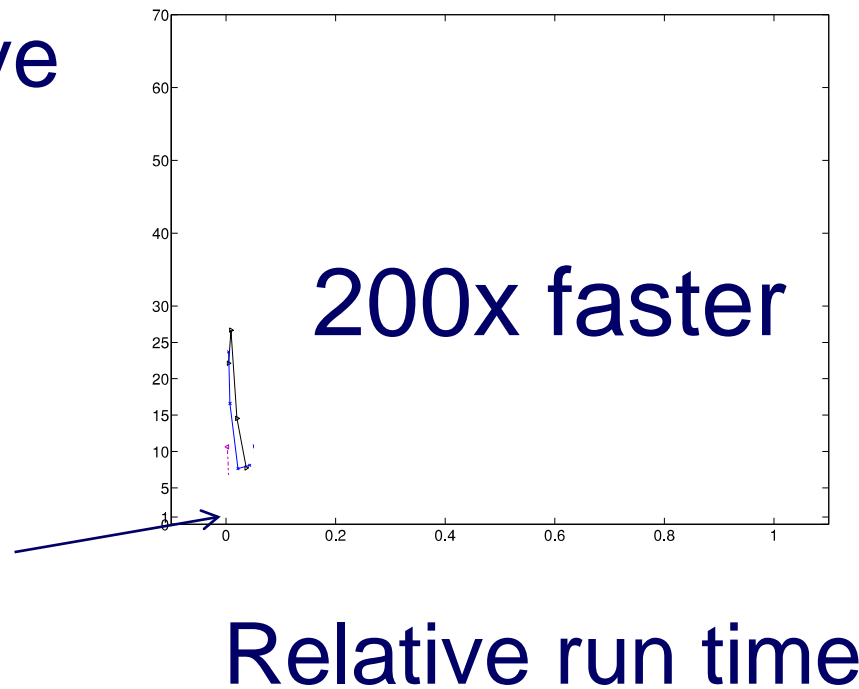


# Speed of tensor/CMTF analysis

- Q1: Can we make it fast?
- Q2: Does it work for large, disk-based data?

# A1: Turbo-SMT

Relative  
cost

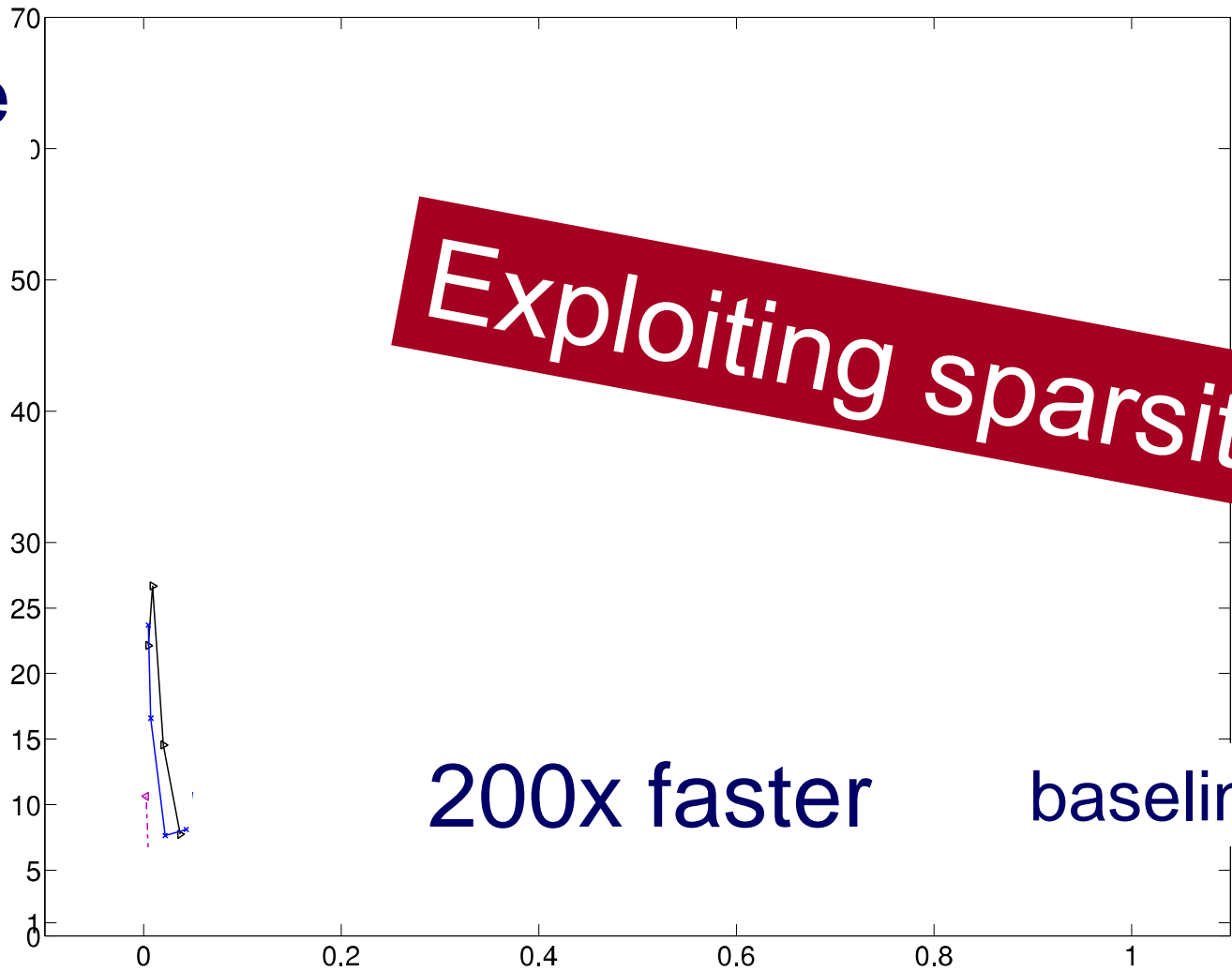


Ideal

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

# A1: Turbo-SMT

Relative cost



200x faster

baseline

Ideal

Relative run time

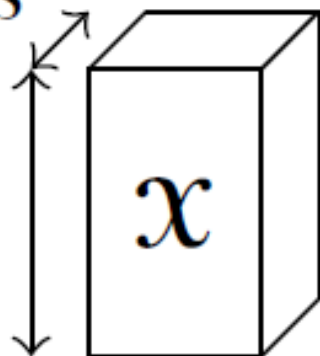
## Q2: spilling to the disk?

Reminder: tensor (eg., Subject-verb-object)

144M non-zeros

48M verbs

26M subjects



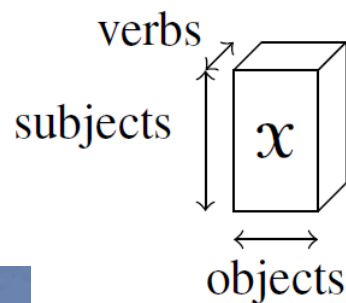
26M objects

NELL (Never Ending  
Language Learner)  
@CMU

# A2: GigaTensor

Reminder: tensor (eg., Subject-verb-object)

26M x 48M x 26M, 144M non-zeros



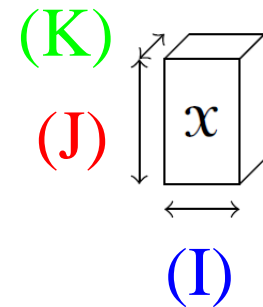
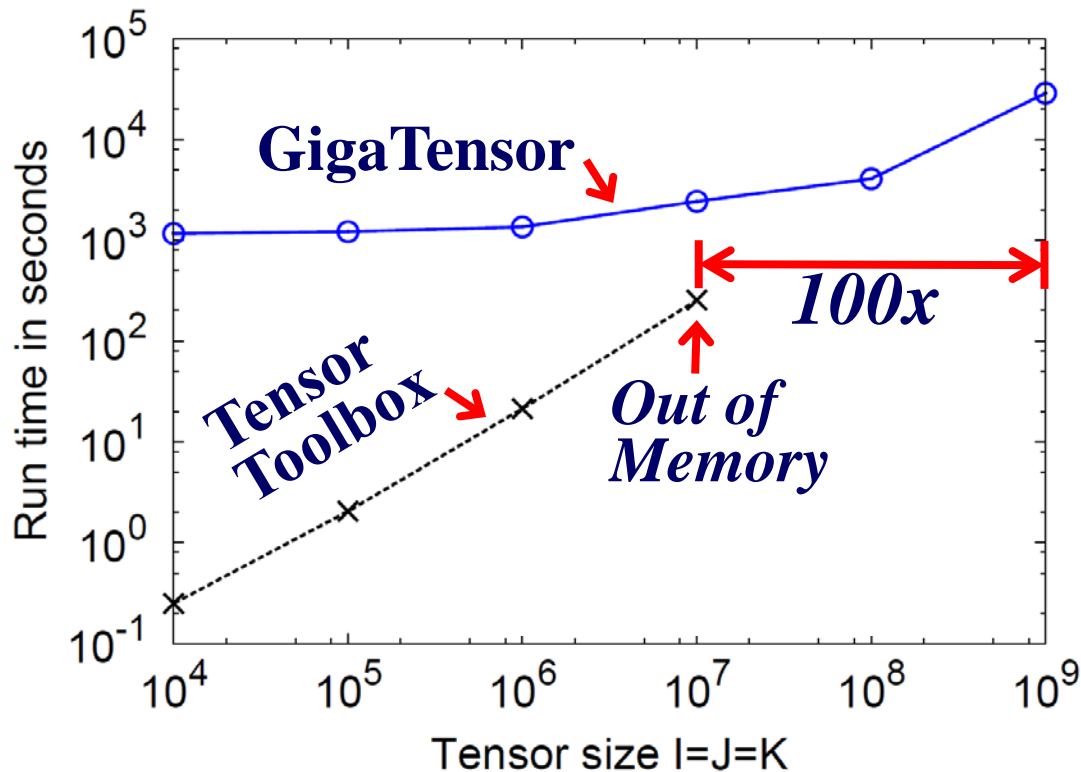
NELL (Never Ending  
Language  
Learner)@CMU



U Kang, Evangelos E. Papalexakis, Abhay Harpale, Christos Faloutsos, *GigaTensor: Scaling Tensor Analysis Up By 100 Times - Algorithms and Discoveries*, KDD'12

# A2: GigaTensor

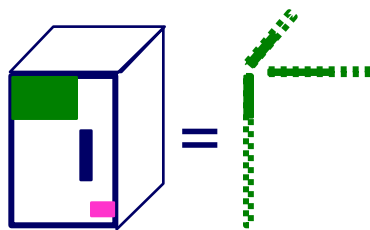
- GigaTensor solves *100x* larger problem



Number of  
nonzero  
=  $I / 50$

## Part 2: Conclusions

- Time-evolving / heterogeneous graphs -> tensors
- PARAFAC finds patterns
- Turbo-SMT; GigaTensor -> fast & scalable



# Roadmap



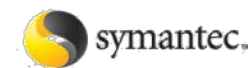
- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
- ➔ • Part#3: Cascades and immunization
- Conclusions



# Part 3: Cascades & Immunization

# Why do we care?

- Information Diffusion
- Viral Marketing
- Epidemiology and Public Health
- Cyber Security
- Human mobility
- Games and Virtual Worlds
- Ecology
- .....



# Roadmap

- A case for cross-disciplinarity
- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: Cascade analysis
  - ➔ – (Fractional) Immunization
  - Epidemic thresholds
- Conclusions



# *Fractional Immunization of Networks*

B. Aditya Prakash,

Lada Adamic,

Theodore Iwashyna (M.D.),

Hanghang Tong,

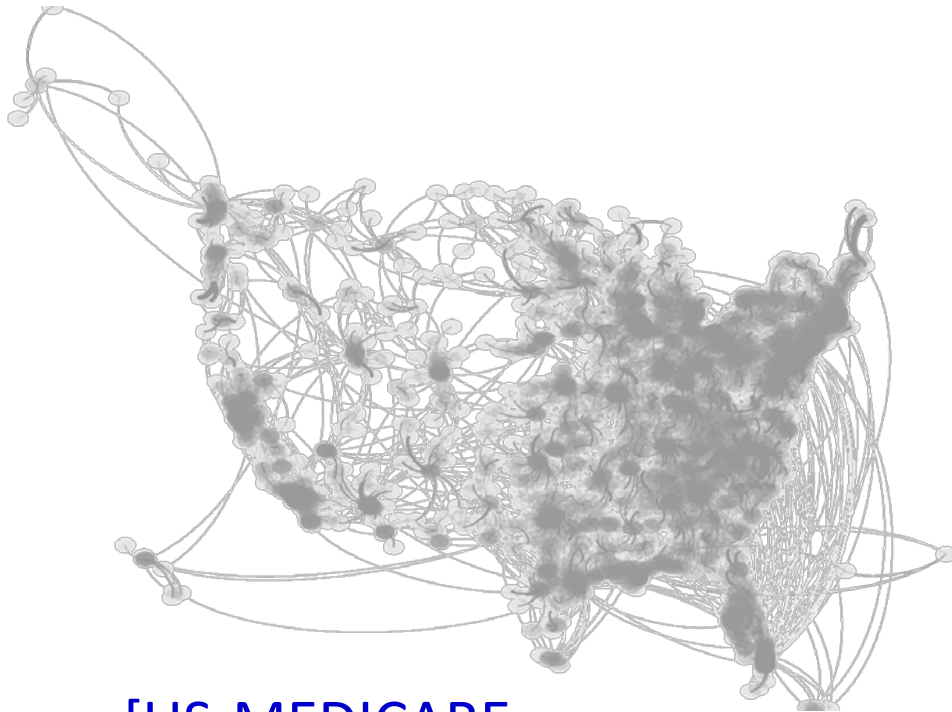
Christos Faloutsos

SDM 2013, Austin, TX



# Whom to immunize?

- Dynamical Processes over networks



[US-MEDICARE  
NETWORK 2005]

WWW, Seoul

- Each circle is a hospital
- ~3,000 hospitals
- More than 30,000 patients transferred

**Problem:** Given  $k$  units of disinfectant, whom to immunize?

(c) 2014, C. Faloutsos

# Whom to immunize?

**~6x  
fewer!**



CURRENT PRACTICE

[US-MEDICARE  
NETWORK 2005]



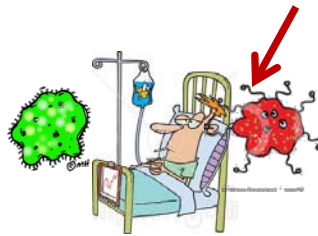
OUR METHOD

Hospital-acquired inf. : 99K+ lives, \$5B+ per year

# Fractional Asymmetric Immunization



Hospital



Drug-resistant Bacteria  
(like XDR-TB)



30%



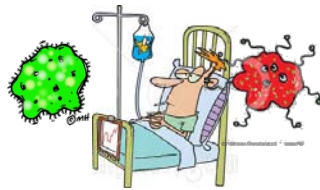
Another  
Hospital



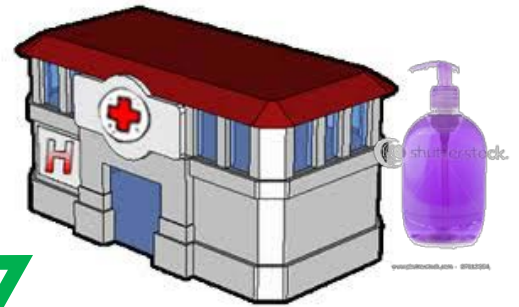
# Fractional Asymmetric Immunization



Hospital



30%



Another Hospital

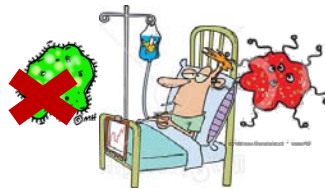




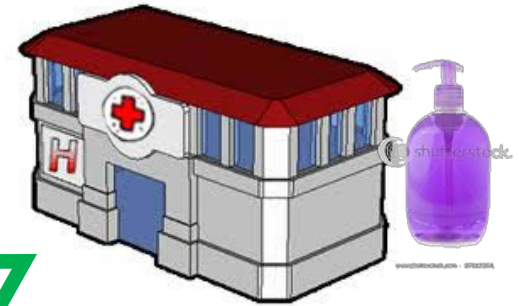
# Fractional Asymmetric Immunization



Hospital



15%



Another Hospital



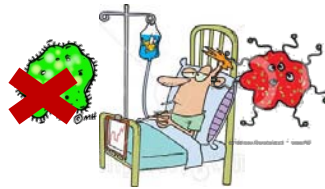
# Fractional Asymmetric Immunization



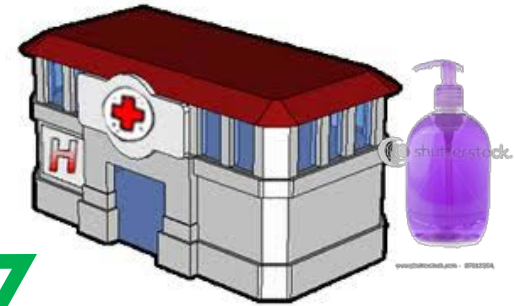
**Problem:**  
*Given  $k$  units of disinfectant, distribute them to maximize hospitals saved*



Hospital



15%



Another Hospital



# Fractional Asymmetric Immunization

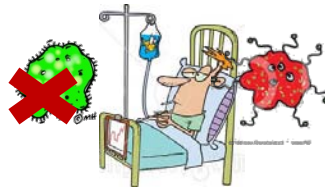


## Problem:

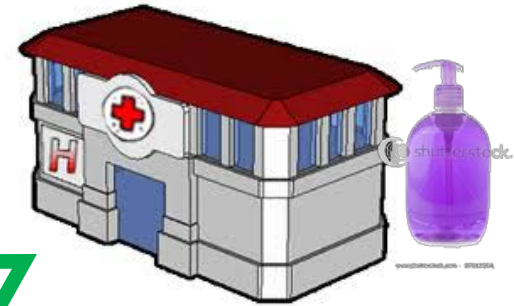
*Given  $k$  units of disinfectant, distribute them to maximize hospitals saved @ 365 days*



Hospital



15%



Another Hospital



# Running Time

Wall-Clock  
Time

> 1 week

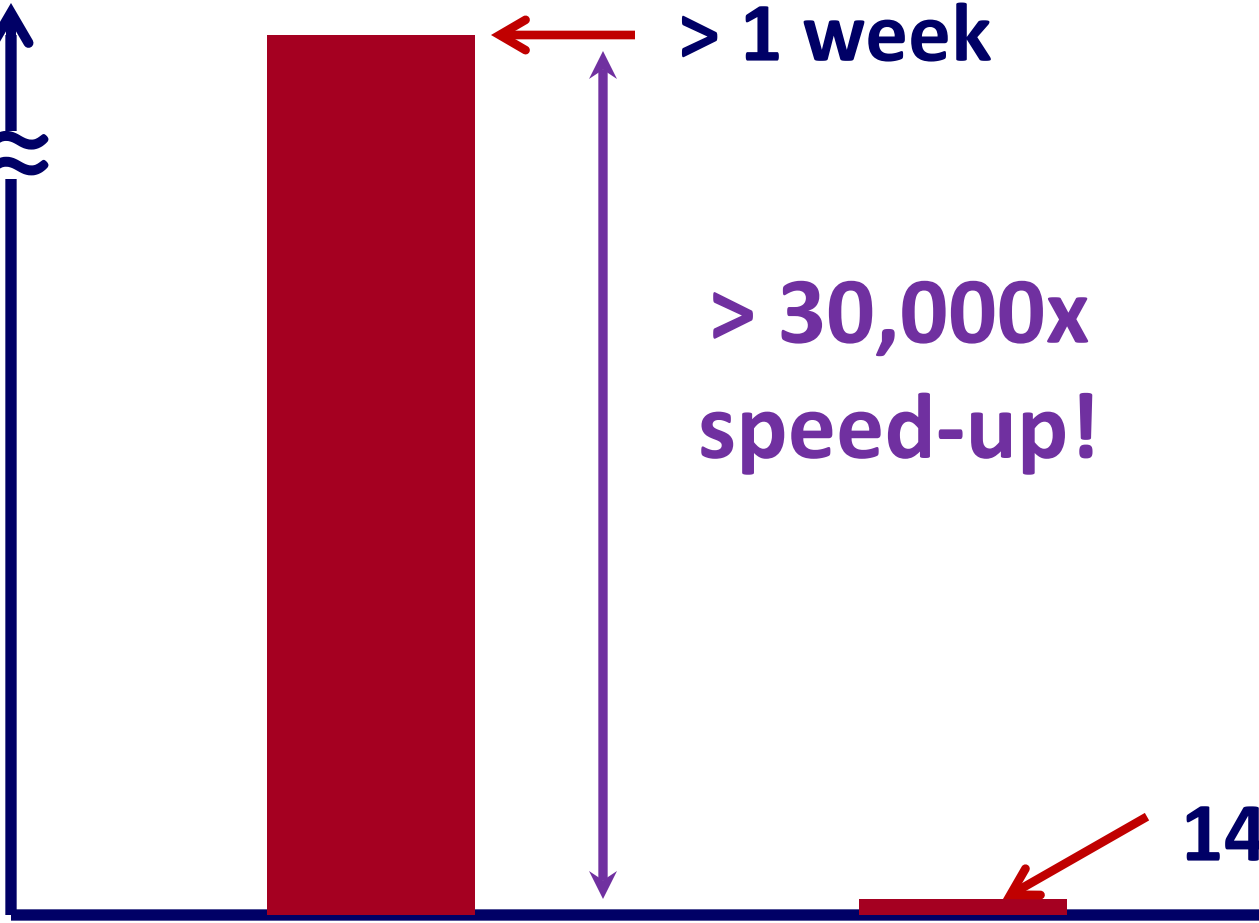
> 30,000x  
speed-up!

↓ *better*

14 secs

Simulations

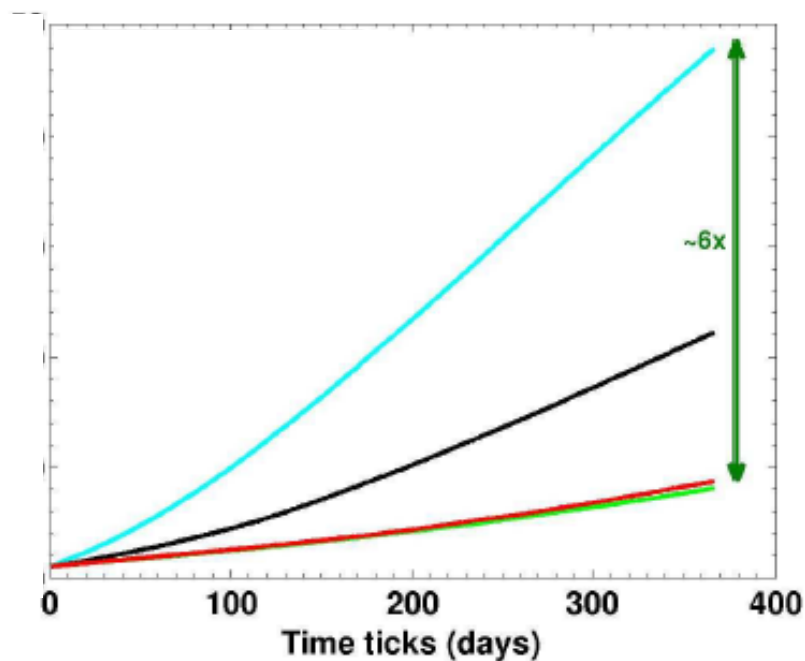
SMART-ALLOC



# Experiments



# infected



uniform

↓ *better*

SMART-ALLOC

$K = 120$

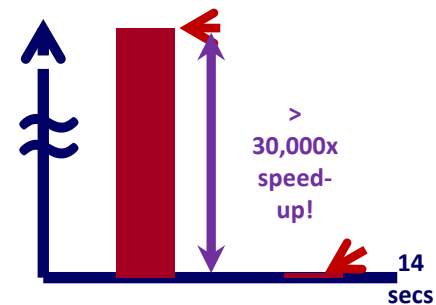
# epochs

# What is the ‘silver bullet’?

A: Try to decrease connectivity of graph

Q: how to measure connectivity?

- Avg degree? Max degree?
- Std degree / avg degree ?
- Diameter?
- Modularity?
- ‘Conductance’ (~min cut size)?
- Some combination of above?



# What is the 'silver bullet'?

A: Try to decrease connectivity of graph

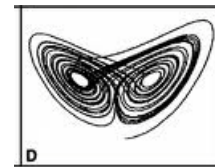
Q: how to measure connectivity?

A: first **eigenvalue** of adjacency matrix

Q1: why??

(Q2: defn & intuition of eigenvalue ? )

Avg degree  
Max degree  
Diameter  
Modularity  
'Conductance'



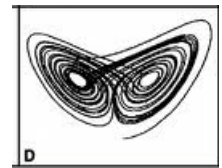
# Why eigenvalue?

A1: ‘G2’ theorem and ‘eigen-drop’:

- For (almost) **any** type of virus
- For **any** network
- -> no epidemic, if small-enough first eigenvalue ( $\lambda_1$ ) of *adjacency* matrix

*Threshold Conditions for Arbitrary Cascade Models on Arbitrary Networks*, B. Aditya Prakash, Deepayan Chakrabarti, Michalis Faloutsos, Nicholas Valler, Christos Faloutsos, ICDM 2011, Vancouver, Canada





# Why eigenvalue?

A1: ‘G2’ theorem and ‘eigen-drop’:

- For (almost) **any** type of virus
- For **any** network
- -> no epidemic, if small-enough first eigenvalue ( $\lambda_1$ ) of *adjacency* matrix
- Heuristic: for immunization, try to min  $\lambda_1$
- The smaller  $\lambda_1$ , the closer to extinction.

# G2 theorem



## *Threshold Conditions for Arbitrary Cascade Models on Arbitrary Networks*



B. Aditya Prakash, Deepayan Chakrabarti,  
Michalis Faloutsos, Nicholas Valler,  
Christos Faloutsos



IEEE ICDM 2011, Vancouver

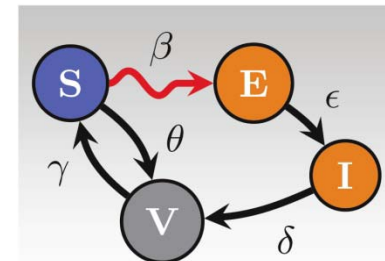
extended version, in arxiv

<http://arxiv.org/abs/1004.0060>

~10 pages proof

# Our thresholds for some models

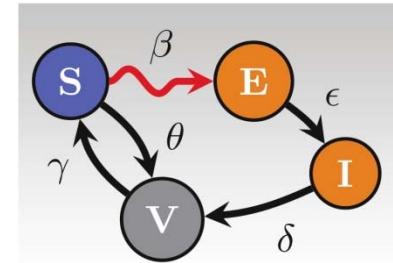
- $s = \text{effective strength}$
- $s < 1$  : *below threshold*



Models	Effective Strength (s)	Threshold (tipping point)
SIS, SIR, SIRS, SEIR	$s = \lambda \left( \frac{\beta}{\delta} \right)$	$s = 1$
SIV, SEIV	$s = \lambda \cdot \left( \frac{\beta\gamma}{\delta(\gamma + \theta)} \right)$	
<u>S</u> <u>(H.I.V.)</u> <sub>2</sub> <u>I</u> <sub>1</sub> <u>V</u> <sub>1</sub> <u>V</u> <sub>2</sub>	$s = \lambda \cdot \left( \frac{\beta_1 v_2 + \beta_2 \epsilon}{v_2 (\epsilon + v_1)} \right)$	

# Our thresholds for some models

- $s = \text{effective strength}$
- $s < 1$  : *below threshold*

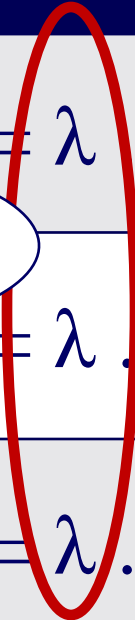


No immunity

Temp. immunity

	Effective Strength	Threshold (tipping point)
SIS, SIR, SIRS, SEIR	$s = \lambda \left( \frac{\beta}{\delta} \right)$	$s = 1$
SIV, SEIV	$s = \lambda \cdot \left( \frac{\beta\gamma}{\delta(\gamma + \theta)} \right)$	
<u>S</u> , <u>I</u> , <u>I</u> , <u>V</u> , <u>V</u> <u>(H.I.V.)</u> <sub>2</sub>	$s = \lambda \cdot \left( \frac{\beta_1 v_2 + \beta_2 \epsilon}{v_2(\epsilon + v_1)} \right)$	

w/ incubation



# Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: Cascade analysis
  - (Fractional) Immunization
  - intuition behind  $\lambda_1$
- Conclusions

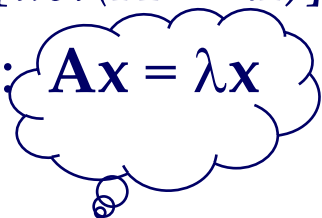


# Intuition for $\lambda$

## “Official” definitions:

- Let  $A$  be the adjacency matrix. Then  $\lambda$  is the root with the largest magnitude of the characteristic polynomial of  $A$  [ $\det(A - \lambda I)$ ].

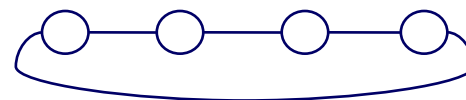
- Also:  $\mathbf{Ax} = \lambda \mathbf{x}$



Neither gives much intuition!

## “Un-official” Intuition

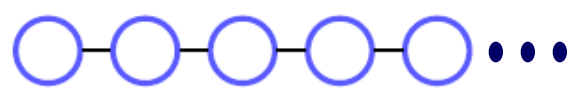
- For ‘homogeneous’ graphs,  $\lambda \approx \text{degree}$



- $\lambda \sim \text{avg degree}$ 
  - done right, for skewed degree distributions

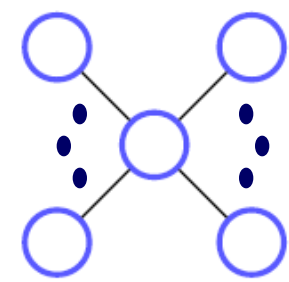
# Largest Eigenvalue ( $\lambda$ )

better connectivity  $\longrightarrow$  higher  $\lambda$



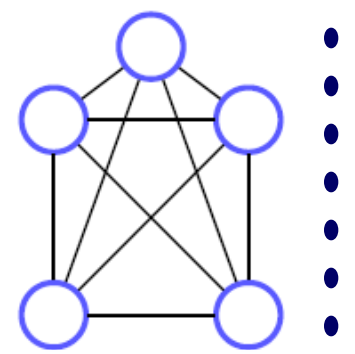
$$\lambda \approx 2$$

(a) Chain



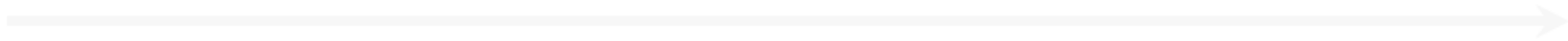
$$\lambda = \sqrt{N}$$

(b) Star



$$\lambda = N-1$$

(c) Clique



$$\lambda \approx 2$$

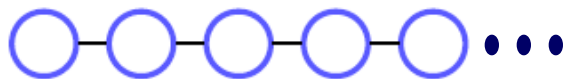
$$\lambda = 31.67$$

$$\lambda = 999$$

$N = 1000$  nodes

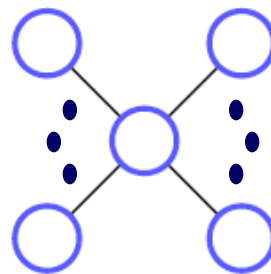
# Largest Eigenvalue ( $\lambda$ )

better connectivity  $\longrightarrow$  higher  $\lambda$



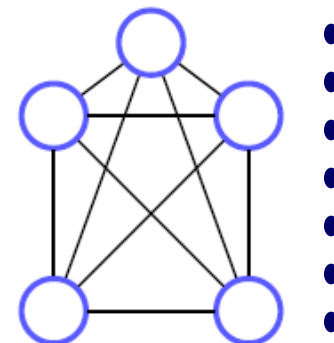
$$\lambda \approx 2$$

(a) Chain



$$\lambda = \sqrt{N}$$

(b) Star



$$\lambda = N-1$$

(c) Clique

$\lambda \approx 2$   
 $N = 1000$  nodes

WWW, Seoul

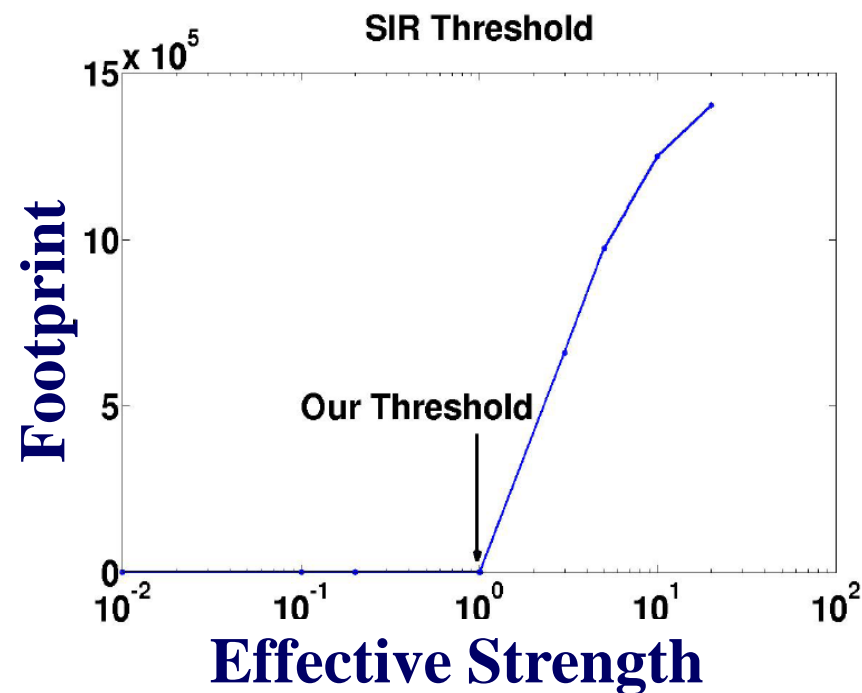
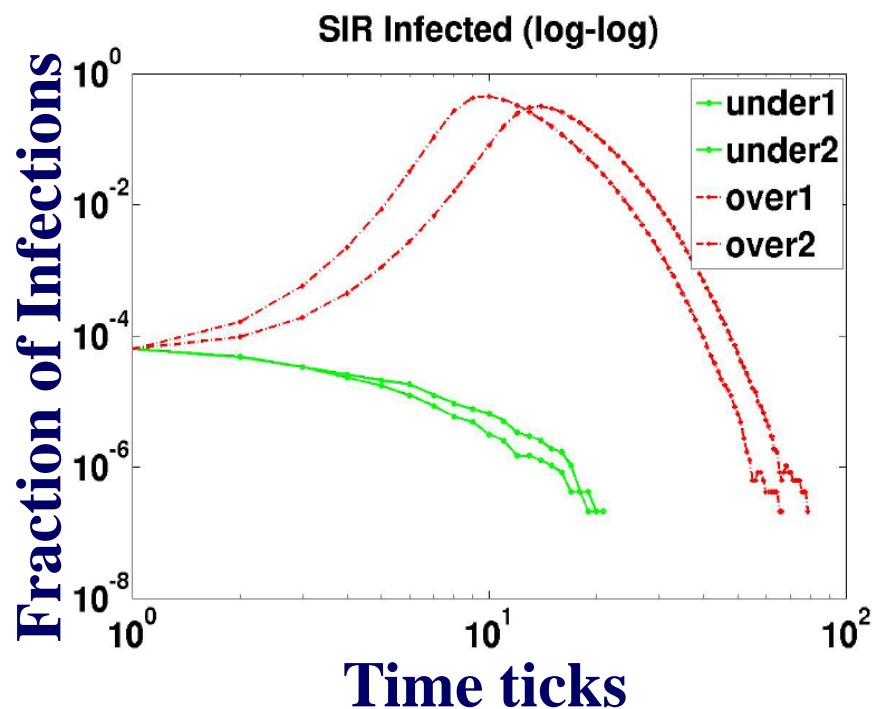
$\lambda = 31.67$

(c) 2014, C. Faloutsos

$\lambda = 999$



# Examples: Simulations – SIR (mumps)

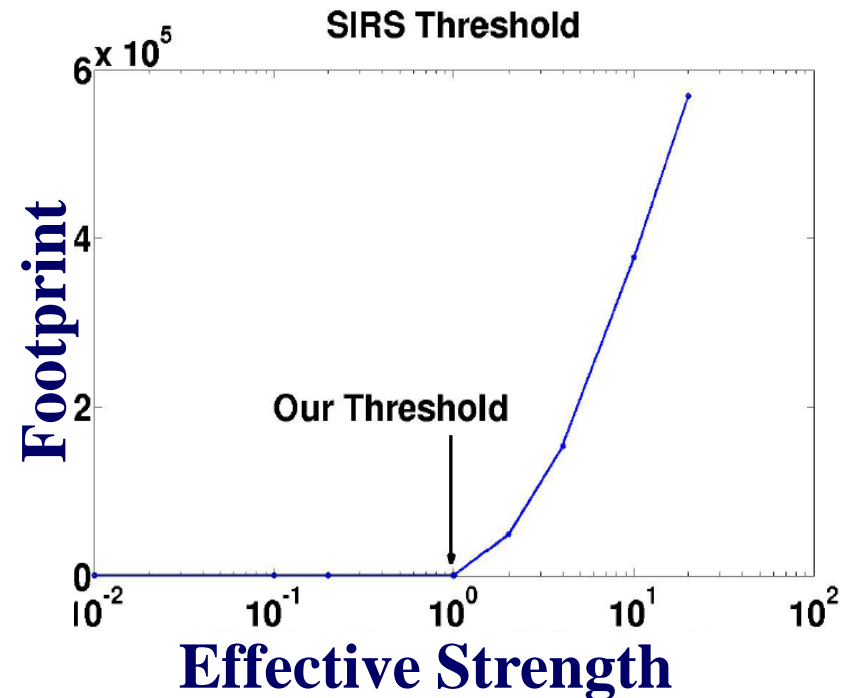
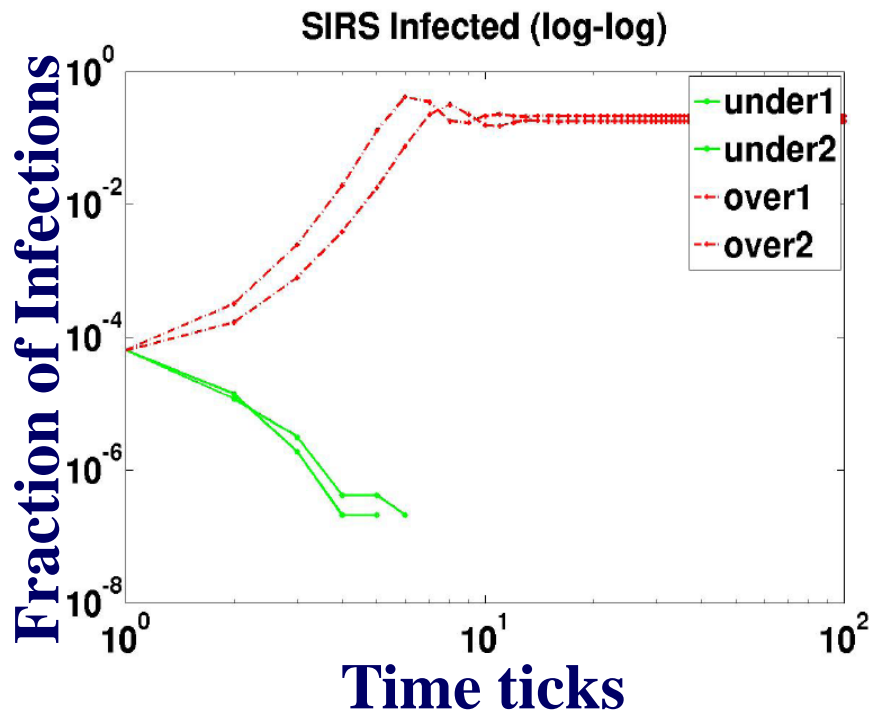


(a) Infection profile

(b) “Take-off” plot

PORTLAND graph: *synthetic population,*  
*31 million links, 6 million nodes*

# Examples: Simulations – SIRS (pertusis)



(a) Infection profile

(b) "Take-off" plot

PORTLAND graph: *synthetic population,*  
*31 million links, 6 million nodes*

## Part3: Immunization - conclusion

In (**almost any**) immunization setting,

- Allocate resources, such that to
- **Minimize  $\lambda_1$**
- (*regardless* of virus specifics)
  
- Conversely, in a market penetration setting
  - Allocate resources to
  - Maximize  $\lambda_1$

# Roadmap

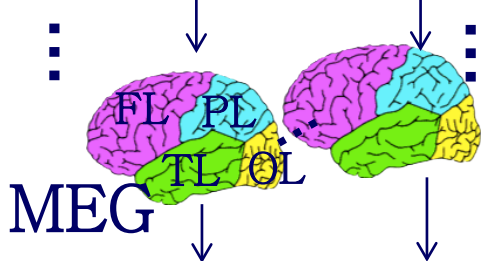


- Introduction – Motivation
  - Why study (big) graphs?
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
- Part#3: Cascades and immunization
- ➔ • Future directions
- Conclusions

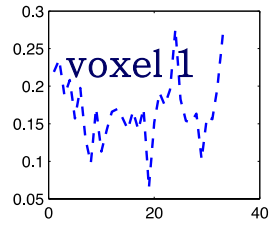
# Brain connectivity

“apple” “Is it edible?” (y/n)

“knife” “Can it hurt you?” (y/n)



MEG



voxel 2

⋮

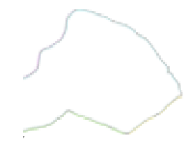
voxel 306

Frontal lobe  
(attention)



Temporal lobe  
(language)

Parietal lobe  
(movement)

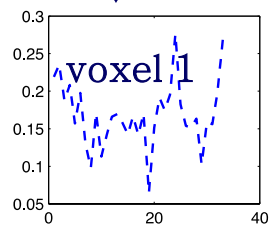
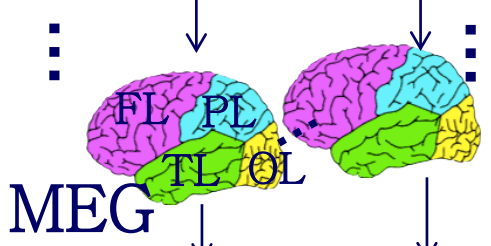


Occipital lobe  
(vision)



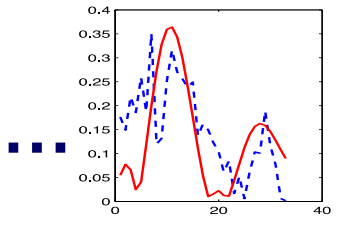
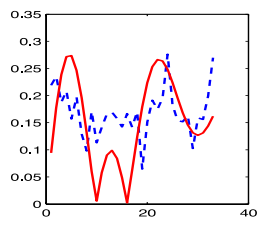
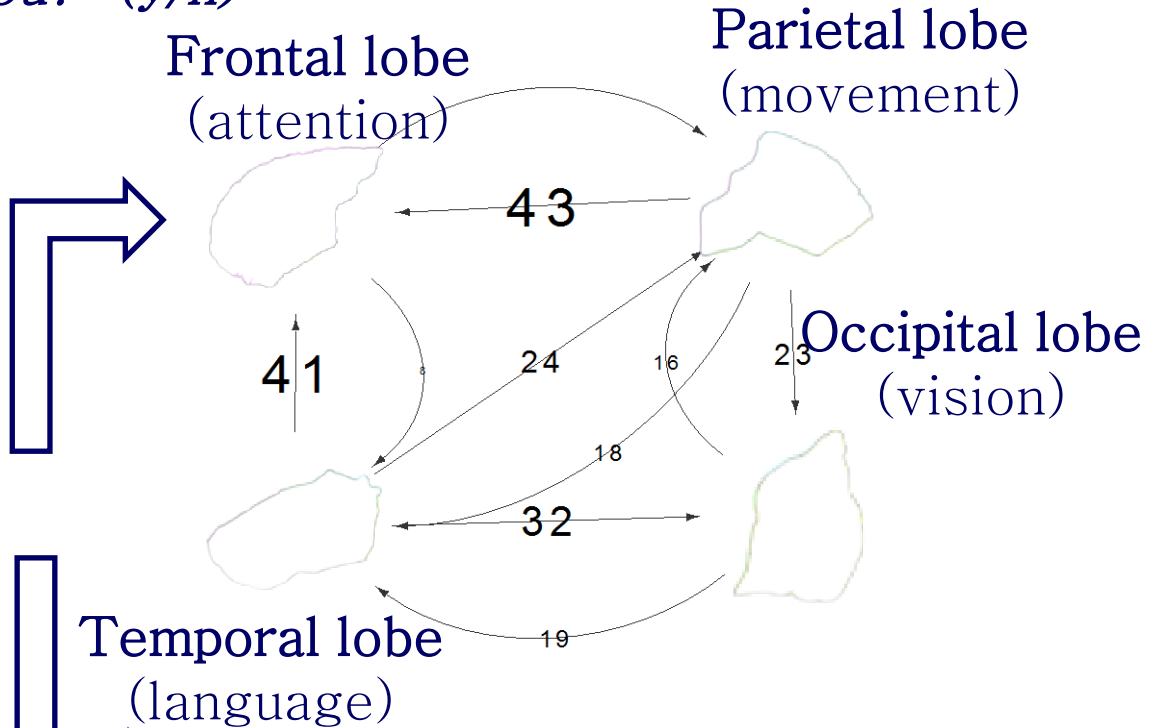
# Brain connectivity

“apple” “Is it edible?” (y/n)  
“knife” “Can it hurt you?” (y/n)



voxel 2

⋮  
voxel 306



# Roadmap



- Introduction – Motivation
  - Why study (big) graphs?
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
- Part#3: Cascades and immunization
- Future directions
- ➔ • Acknowledgements and Conclusions

# Thanks



© 2008 Yahoo!

*Disclaimer: All opinions are mine; not necessarily reflecting the opinions of the funding agencies*

Thanks to: NSF IIS-0705359, IIS-0534205,  
CTA-INARC; Yahoo (M45), LLNL, IBM, SPRINT,  
Google, INTEL, HP, iLab



# Project info: PEGASUS



[www.cs.cmu.edu/~pegasus](http://www.cs.cmu.edu/~pegasus)

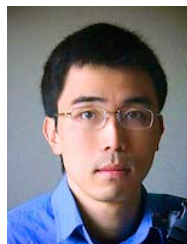
Results on large graphs: with Pegasus +  
hadoop + M45

Apache license

Code, papers, manual, video



Prof. U Kang



Prof. Polo Chau

# Cast



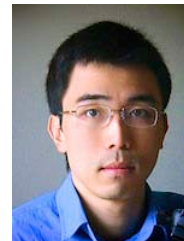
Akoglu,  
Leman



Araujo,  
Miguel



Beutel,  
Alex



Chau,  
Polo



Kang, U



Koutra,  
Danai



Lee,  
Jay Yoon



Prakash,  
Aditya



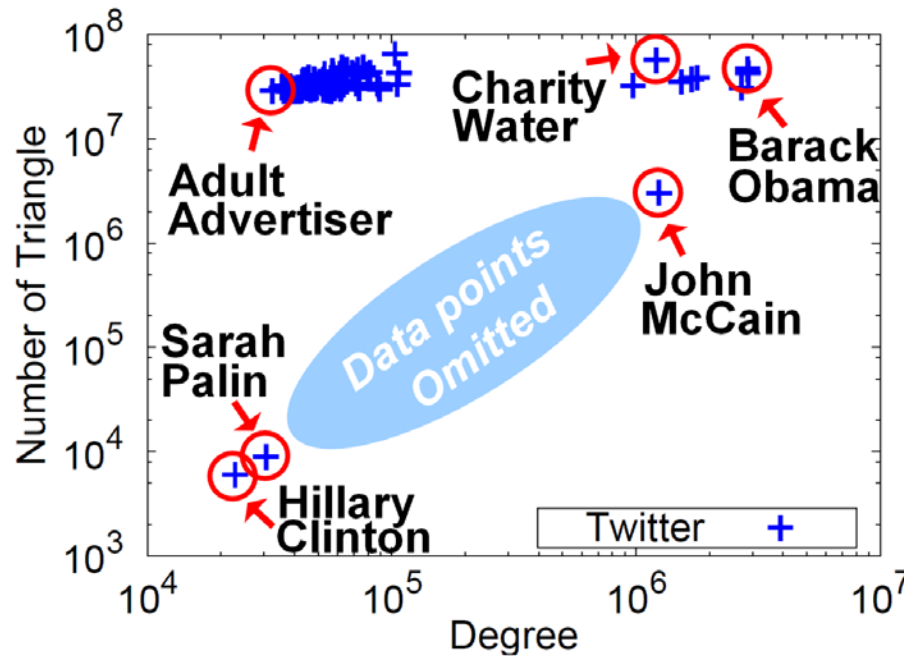
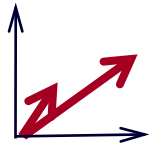
Papalexakis,  
Vagelis



Shah,  
Neil

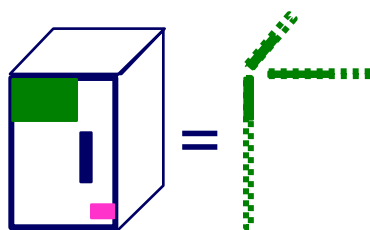
# CONCLUSION#1 – Big data

- Large datasets reveal patterns/outliers that are invisible otherwise



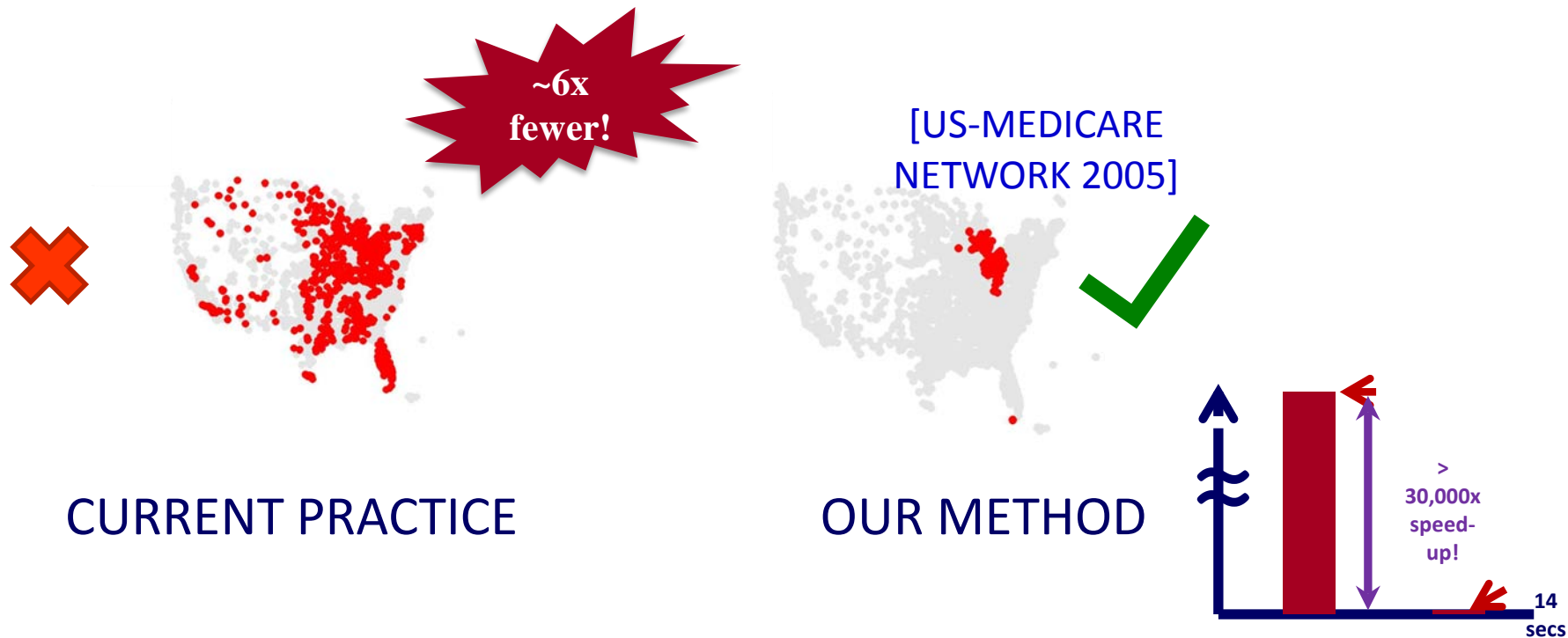
# CONCLUSION#2 – tensors

- powerful tool



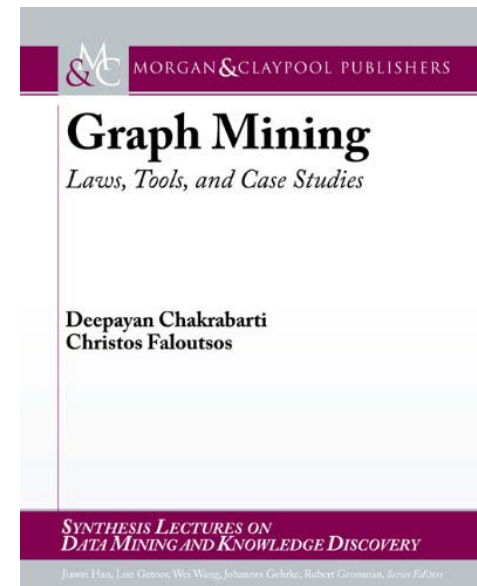
# CONCLUSION#3 – eigen-drop

- Cascades & immunization: G2 theorem & eigenvalue



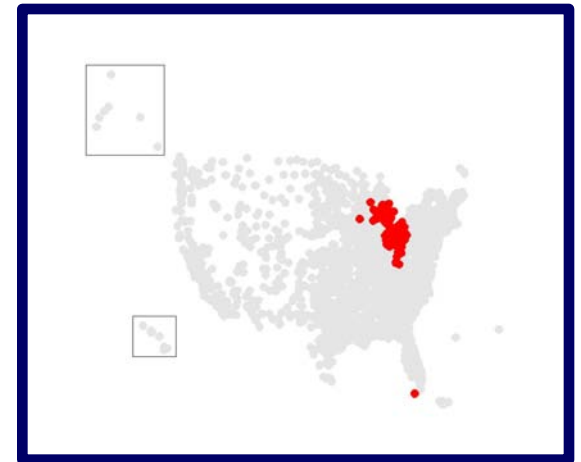
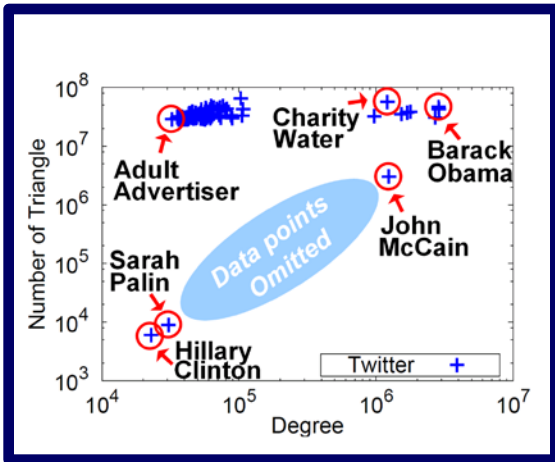
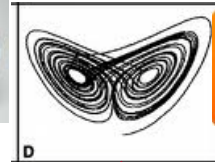
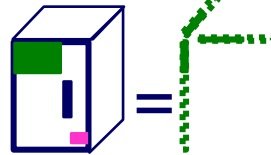
# References

- 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! Questions?

## Cross-disciplinarity

