

# Discovering Roles and Anomalies in Graphs: Theory and Applications

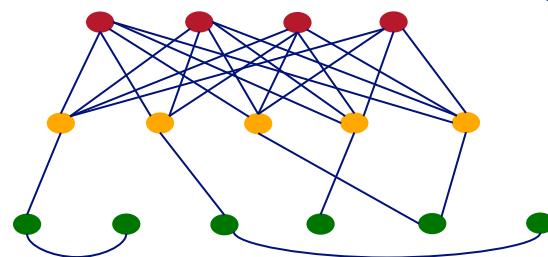
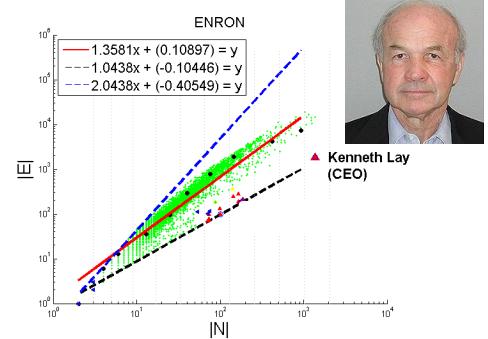
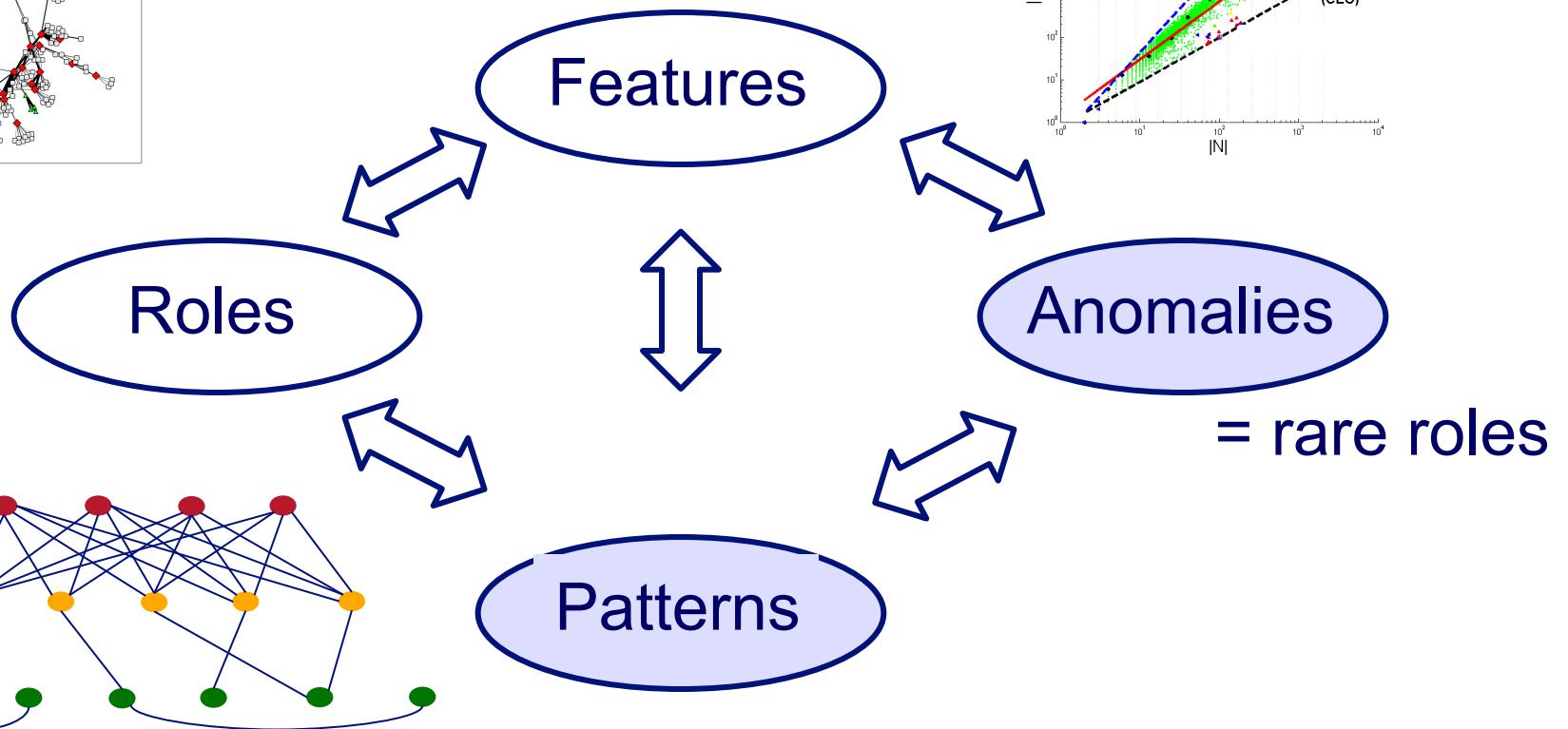
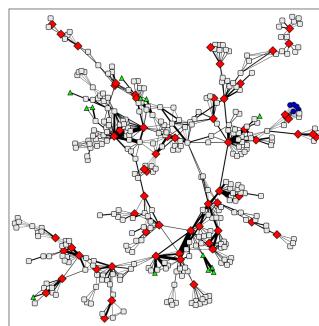
Part 2: Patterns and Anomalies

*Tina Eliassi-Rad* (Rutgers)

*Christos Faloutsos* (CMU)

ECML PKDD 2013 Tutorial

# OVERVIEW - high level:



## Resource:

Open source system for mining huge graphs:

PEGASUS project (PEta GrAph mining System)

- [www.cs.cmu.edu/~pegasus](http://www.cs.cmu.edu/~pegasus)
- code and papers

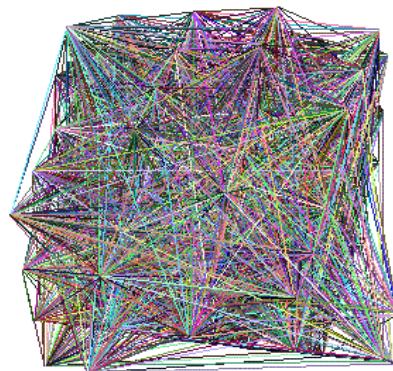


# Roadmap

- ➡ • Patterns in graphs
  - Overview
  - Static graphs
  - Weighted graphs
  - Time-evolving graphs
- Anomaly Detection
- Application: ebay fraud
- Conclusions

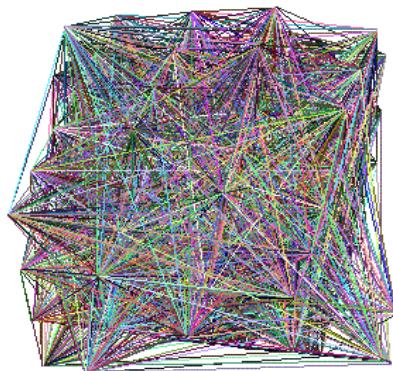


# Problem #1 - network and graph mining



- What does the Internet look like?
- What does FaceBook look like?
  
- What is ‘normal’ / ‘abnormal’ ?
- which patterns/laws hold?

# Problem #1 - network and graph mining

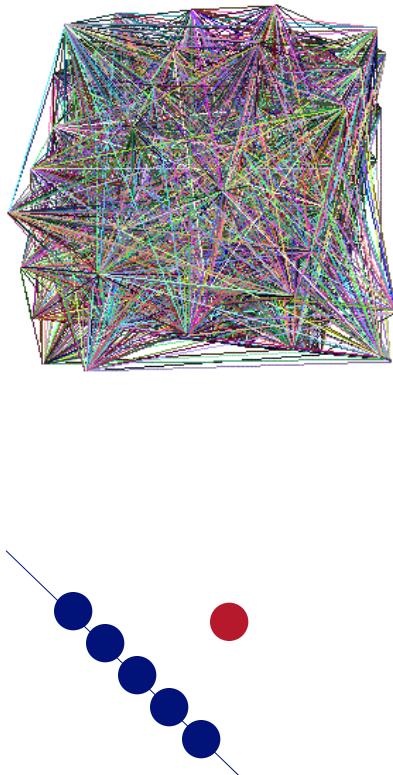


- 
- 
- 

- What does the Internet look like?
- What does FaceBook look like?
  
- What is ‘normal’ / ‘abnormal’ ?
- Which patterns/laws hold?
  - To spot **anomalies** (rarities), we have to discover **patterns**

# Problem #1 - network and graph mining

- What does the Internet look like?
- What does FaceBook look like?
- What is ‘normal’ / ‘abnormal’ ?
- Which patterns/laws hold?
  - To spot **anomalies** (rarities), we have to discover **patterns**
  - **Large** datasets reveal patterns/anomalies that may be invisible otherwise...



# Graph mining

- Are real graphs random?

# Laws and patterns

- Are real graphs random?
- A: NO!!
  - Diameter
  - In- and out- degree distributions
  - Other (surprising) patterns
- So, let's look at the data



# Real Graph Patterns

## Unweighted

**Static**

- P01.** Power-law degree distribution [Faloutsos et. al. '99, Kleinberg et. al. '99, Chakrabarti et. al. '04, Newman '04]
- P02.** Triangle Power Law [Tsourakakis '08]
- P03.** Eigenvalue Power Law [Siganos et. al. '03]
- P04.** Community structure [Flake et. al. '02, Girvan and Newman '02]
- P05.** Clique Power Laws [Du et. al. '09]

**Dynamic**

- P06.** Densification Power Law [Leskovec et. al. '05]
- P07.** Small and shrinking diameter [Albert and Barabási '99, Leskovec et. al. '05, McGlohon et. al. '08]
- P08.** Gelling point [McGlohon et. al. '08]
- P09.** Constant size 2<sup>nd</sup> and 3<sup>rd</sup> connected components [McGlohon et. al. '08]
- P10.** Principal Eigenvalue Power Law [Akoglu et. al. '08]
- P11.** Bursty/self-similar edge/weight additions [Gomez and Santonja '98, Gribble et. al. '98, Crovella and Bestavros '99, McGlohon et. al. '08]

## Weighted

- P12.** Snapshot Power Law [McGlohon et. al. '08]

- P13.** Weight Power Law [McGlohon et. al. '08]
- P14.** Skewed call duration distributions [Vaz de Melo et. al. '10]

[RTG: A Recursive Realistic Graph Generator using Random Typing](#)  
 Leman Akoglu and Christos Faloutsos. *ECML PKDD'09*.

# Roadmap

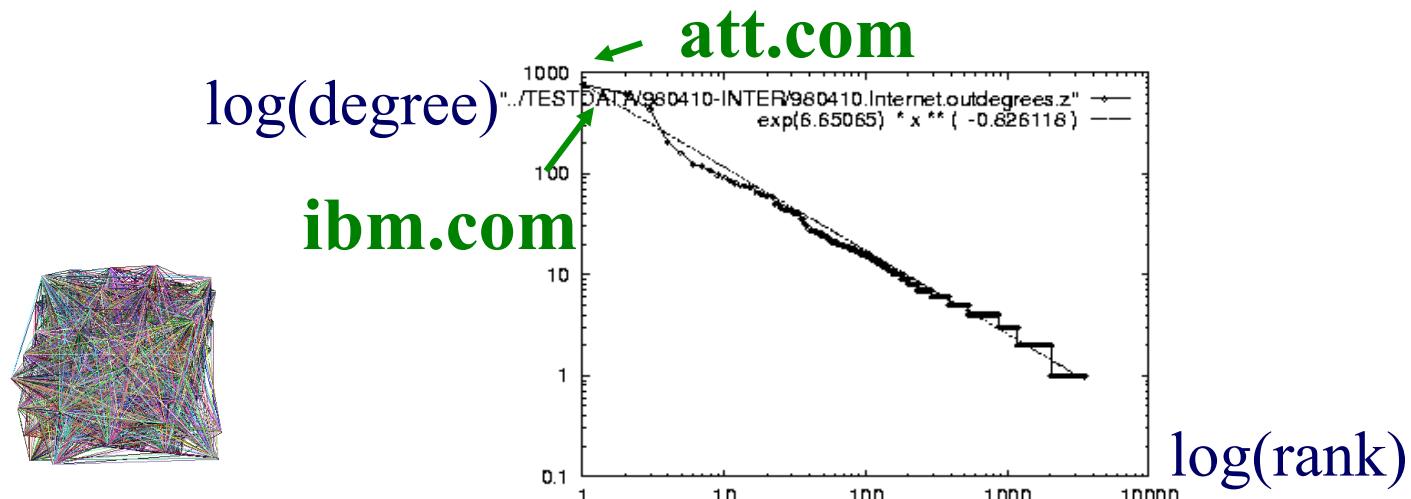
- Patterns in graphs
  - Overview
  - Static graphs
  - Weighted graphs
  - Time-evolving graphs
- Anomaly Detection
- Application: ebay fraud
- Conclusions



# Solution# S.1

- Power law in the degree distribution  
[SIGCOMM99]

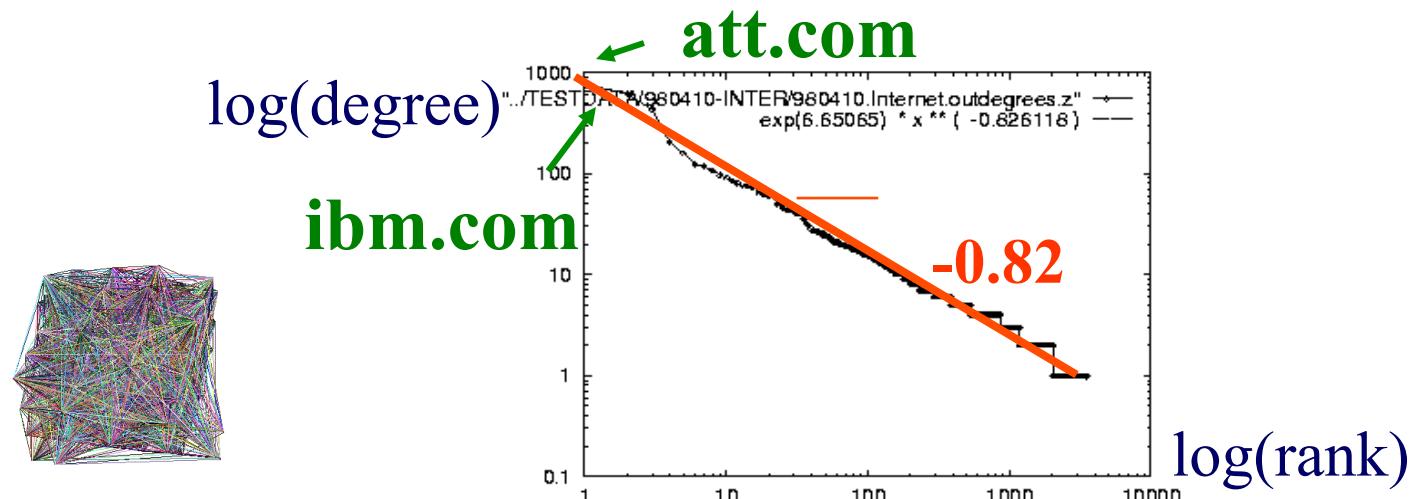
Internet Domains



# Solution# S.1

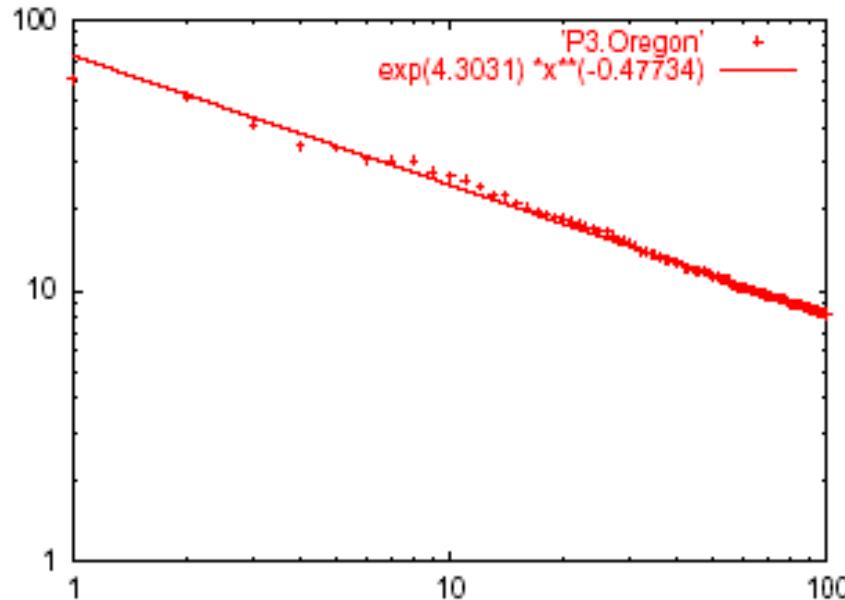
- Power law in the degree distribution [SIGCOMM99]

Internet Domains



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

Eigenvalue



Exponent = slope

$$E = -0.48$$

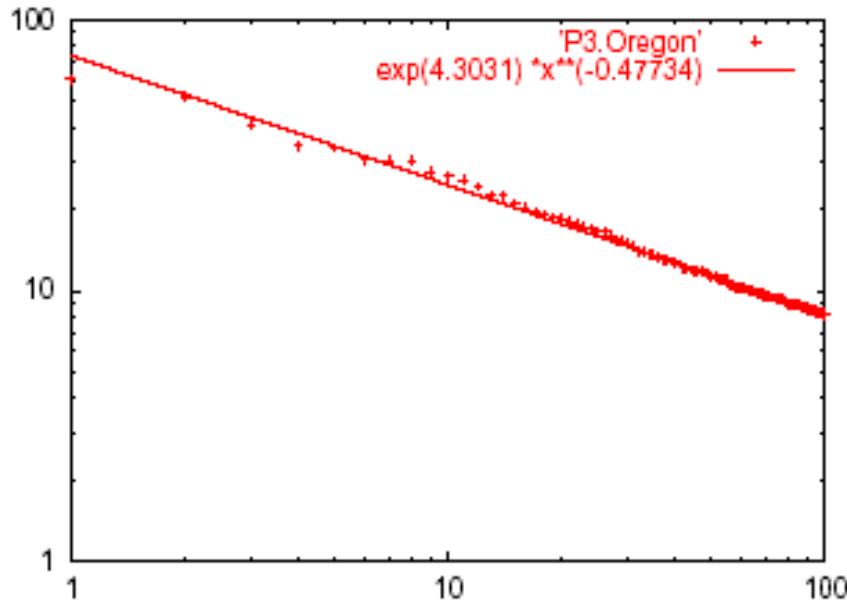
May 2001

Rank of decreasing eigenvalue

- A2: power law in the eigenvalues of the adjacency matrix

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

Eigenvalue



Exponent = slope

$$E = -0.48$$

May 2001

Rank of decreasing eigenvalue

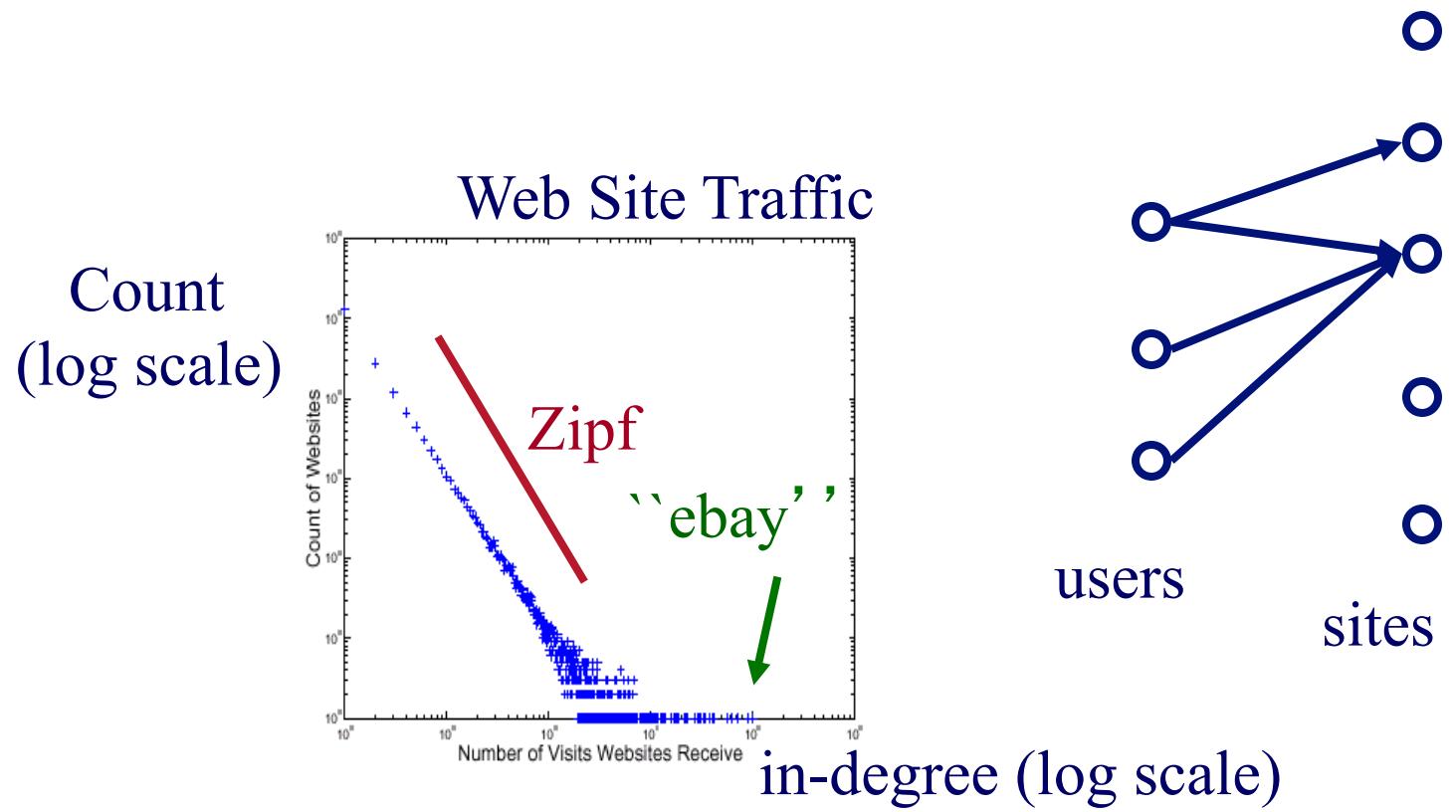
- [Mihail, Papadimitriou '02]: slope is  $\frac{1}{2}$  of rank exponent

# But:

## How about graphs from other domains?

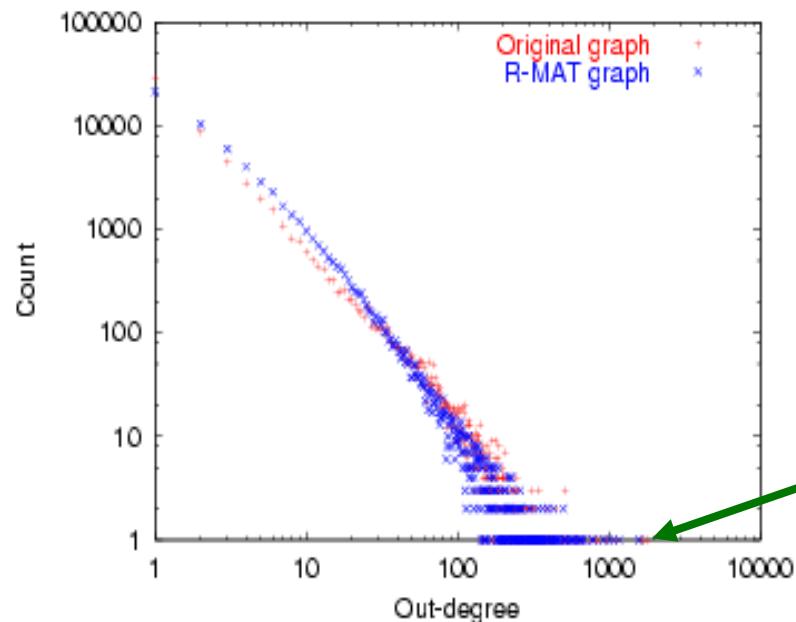
# More power laws:

- web hit counts [w/ A. Montgomery]



# epinions.com

count



- who-trusts-whom  
[Richardson +  
Domingos, KDD  
2001]

trusts-2000-people user

# And numerous more

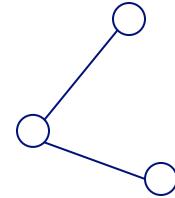
- # of sexual contacts
- Income [Pareto] – ‘80-20 distribution’
- Duration of downloads [Bestavros+]
- Duration of UNIX jobs ('mice and elephants')
- Size of files of a user
- ...
- ‘Black swans’

# Roadmap

- Patterns in graphs
  - overview
  - Static graphs
    - S1: Degree, S2: Eigenvalues
    - S3-4: Triangles, S5: Cliques
    - Radius plot
    - Other observations ('EigenSpokes')
  - Weighted graphs
  - Time-evolving graphs

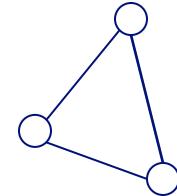


# Solution# S.3: Triangle ‘Laws’



- Real social networks have a lot of triangles

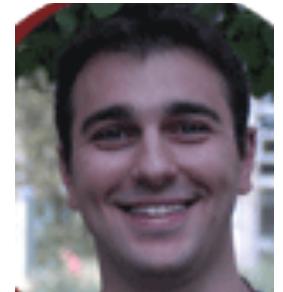
# Solution# S.3: Triangle ‘Laws’



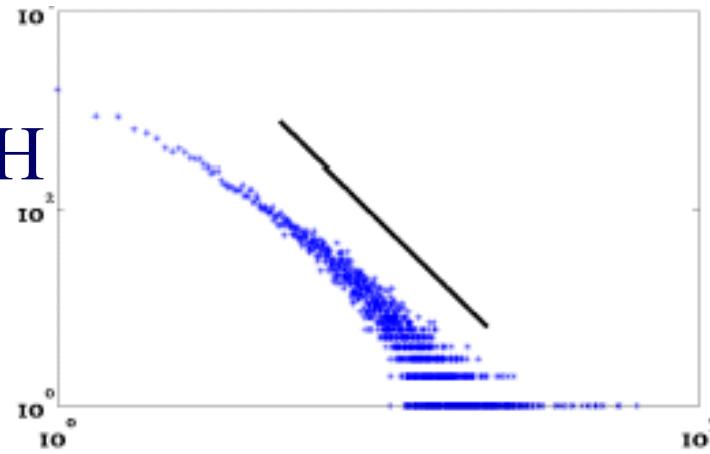
- Real social networks have a lot of triangles
  - Friends of friends are friends
- Any patterns?

# Triangle Law: #S.3

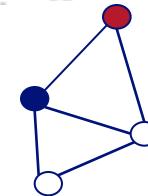
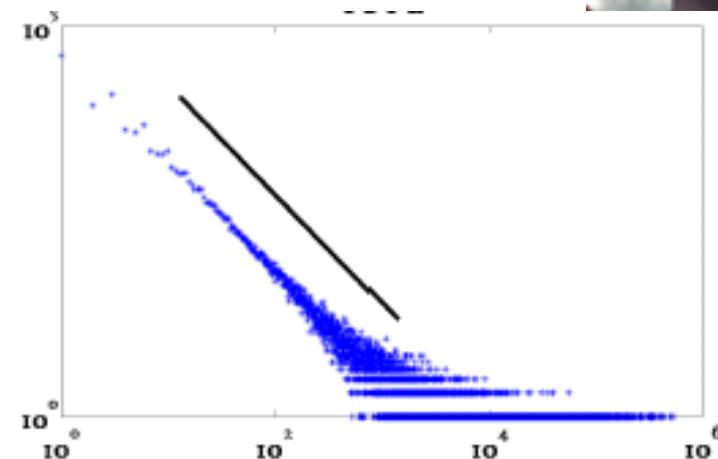
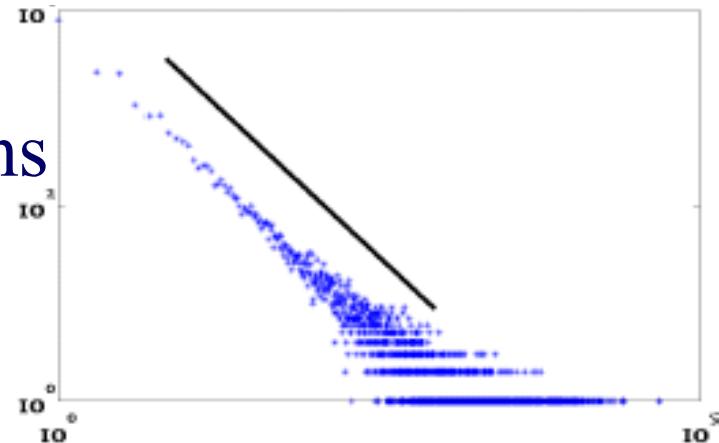
[Tsourakakis ICDM 2008]



HEP-TH



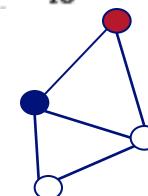
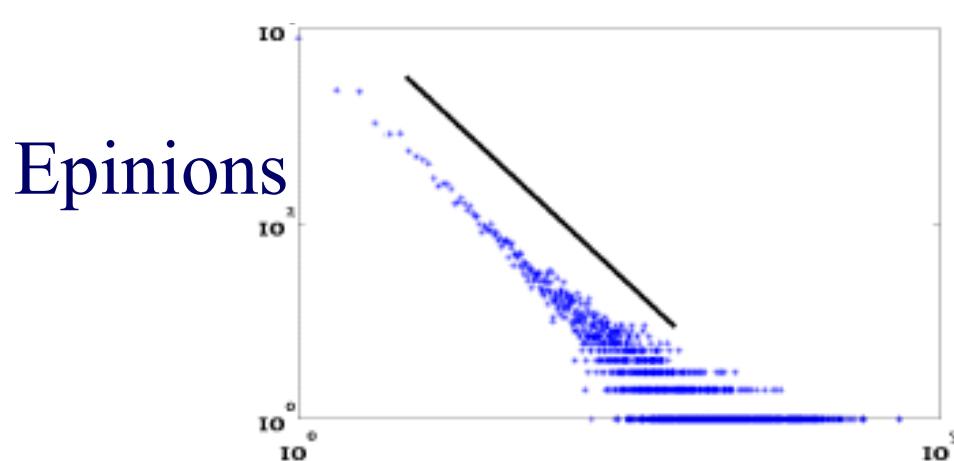
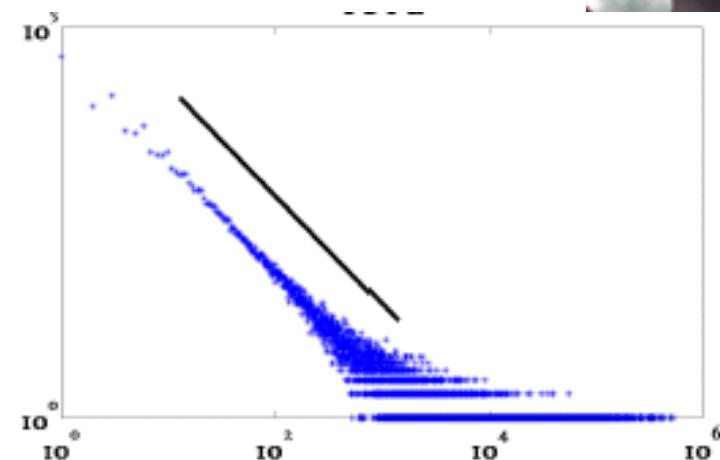
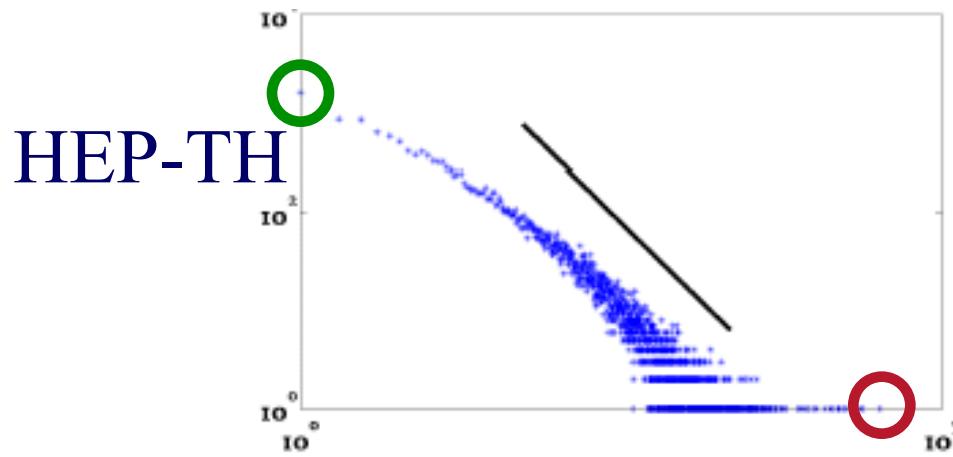
Epinions



X-axis: # of participating triangles  
Y: count ( $\sim$  pdf)

# Triangle Law: #S.3

## [Tsourakakis ICDM 2008]

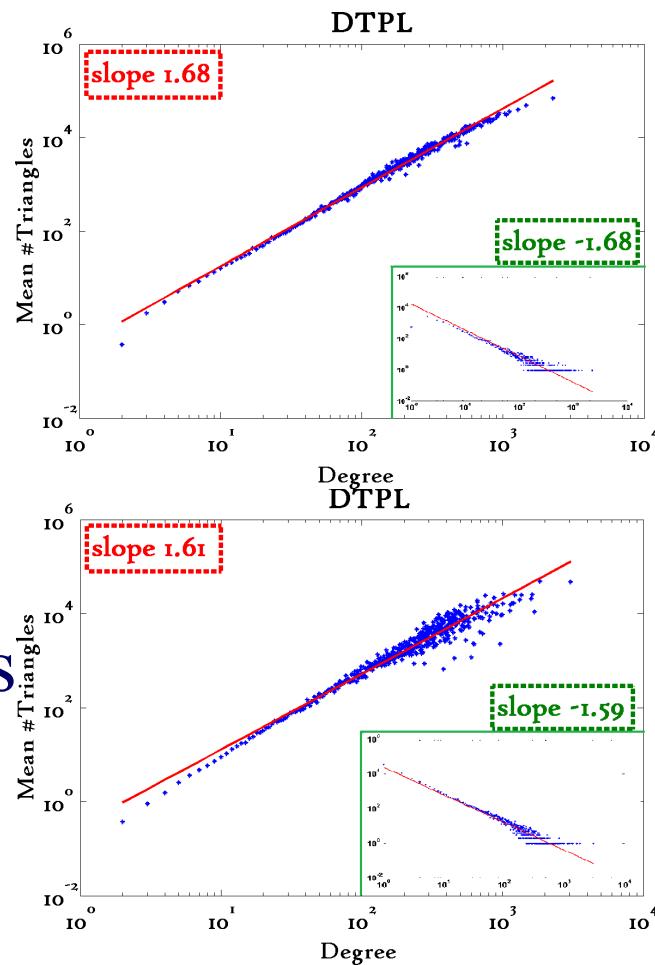


X-axis: # of participating triangles  
Y: count ( $\sim$  pdf)

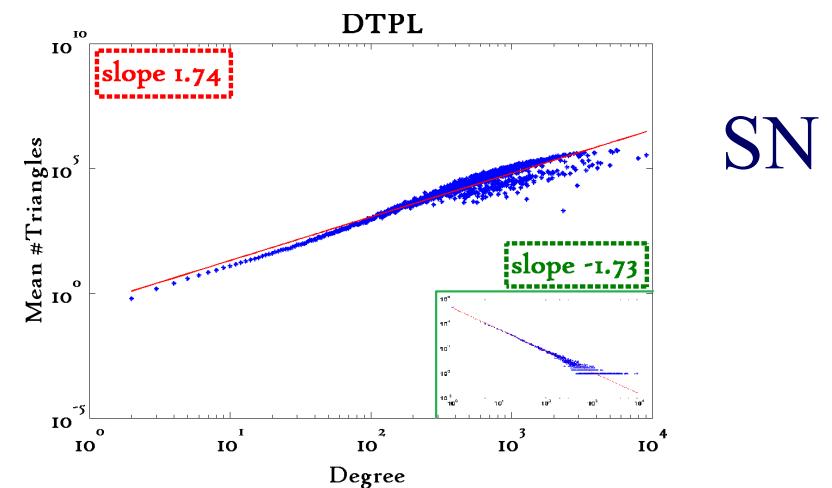
# Triangle Law: #S.4

## [Tsourakakis ICDM 2008]

Reuters



Epinions



SN

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

# Triangle Law: Computations

## [Tsourakakis ICDM 2008]

Triangles are expensive to compute  
(3-way join; several approximation algorithms)

Q: Can we do this computation quickly?

# Triangle Law: Computations

## [Tsourakakis ICDM 2008]

Triangles are expensive to compute  
(3-way join; several approx. algos)

Q: Can we do this computation quickly?

A: Yes!

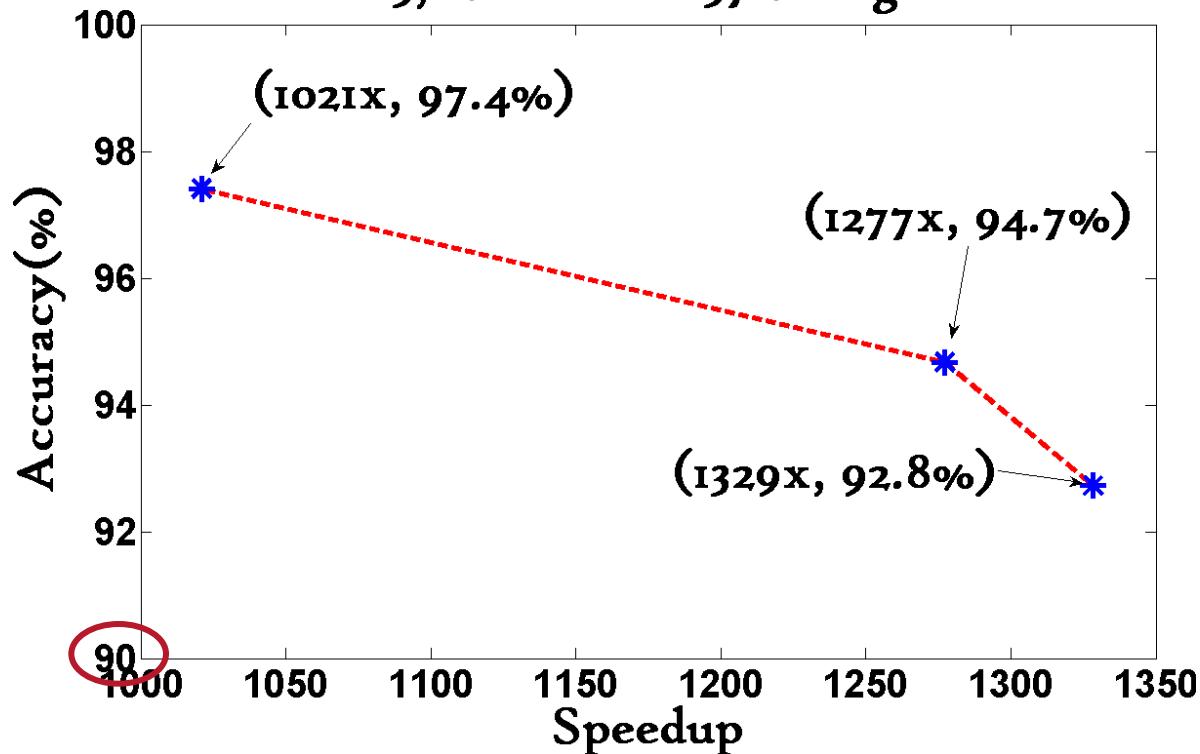
**#triangles = 1/6 Sum (  $\lambda_i^3$  )**  
(and, because of skewness (S2) ,  
we only need the top few eigenvalues!)

# Triangle Law: Computations

## [Tsourakakis ICDM 2008]

Wikipedia graph 2006-Nov-04

$\approx 3.1\text{M}$  nodes  $\approx 37\text{M}$  edges

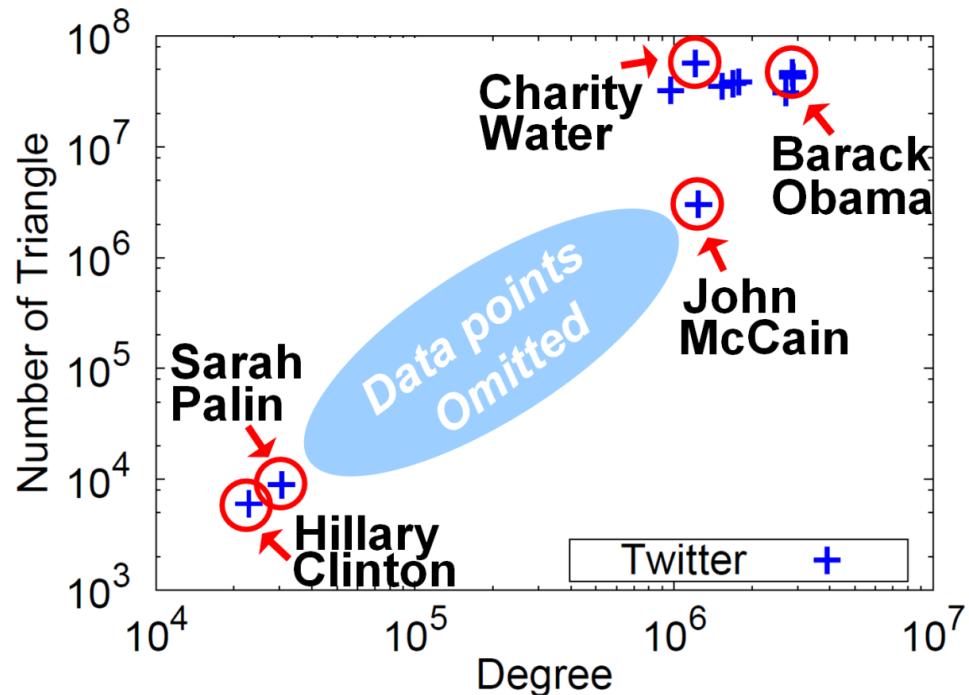


1000x+ speed-up, >90% accuracy

# Triangle Counting on Big Graphs

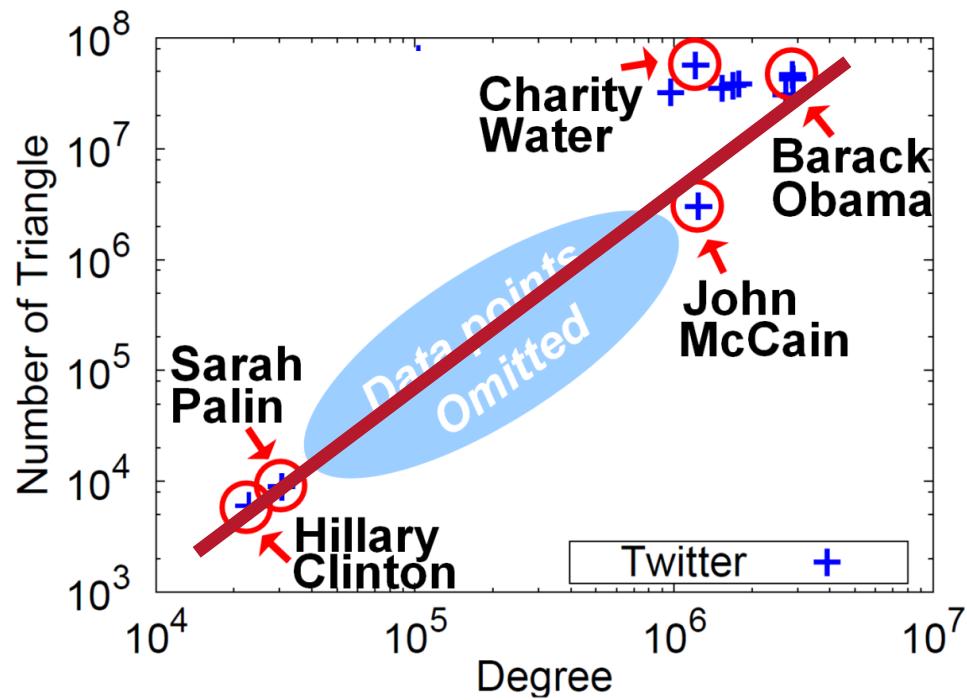
Anomalous nodes in Twitter ( $\sim 3$  billion edges)  
[U Kang, Brendan Meeder, +, PAKDD'11]

# Triangle Counting on Big Graphs



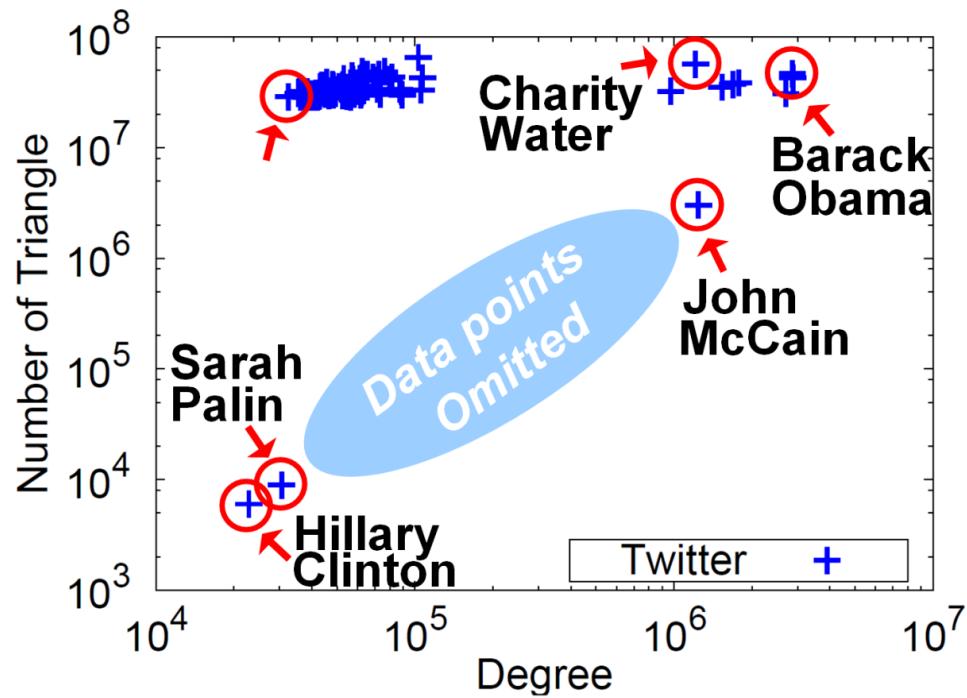
Anomalous nodes in Twitter ( $\sim 3$  billion edges)  
 [U Kang, Brendan Meeder, +, PAKDD'11]

# Triangle Counting on Big Graphs



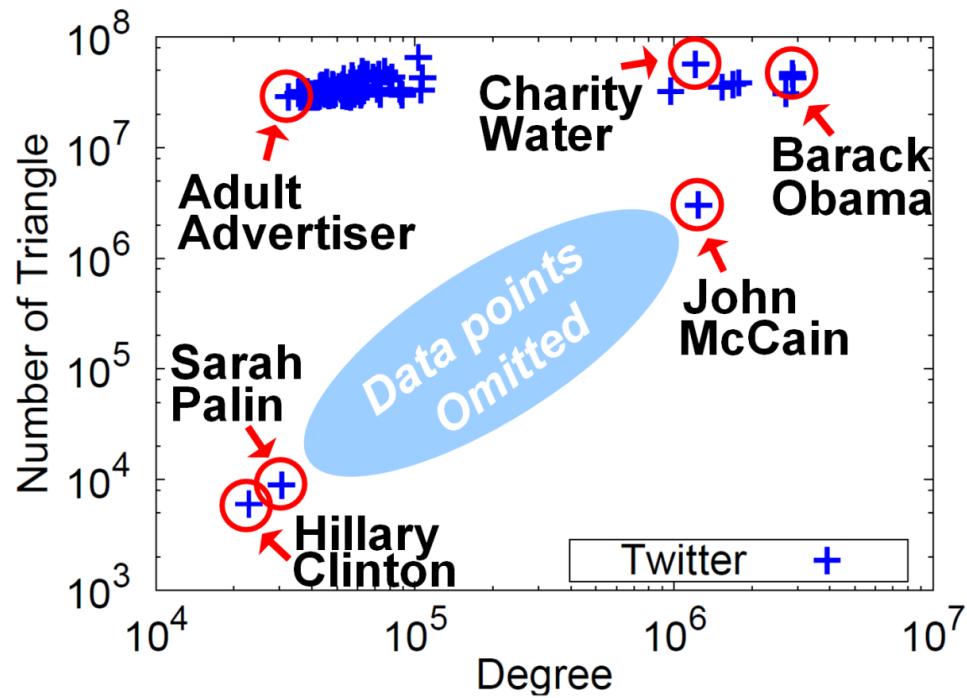
Anomalous nodes in Twitter (~ 3 billion edges)  
 [U Kang, Brendan Meeder, +, PAKDD'11]

# Triangle Counting on Big Graphs



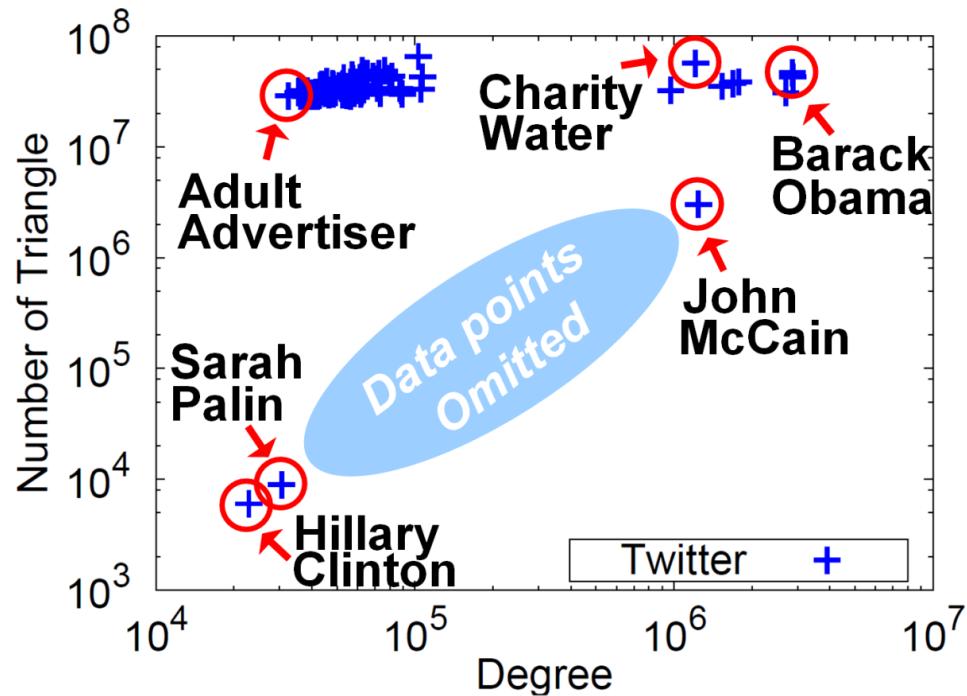
Anomalous nodes in Twitter ( $\sim 3$  billion edges)  
[U Kang, Brendan Meeder, +, PAKDD'11]

# Triangle Counting on Big Graphs



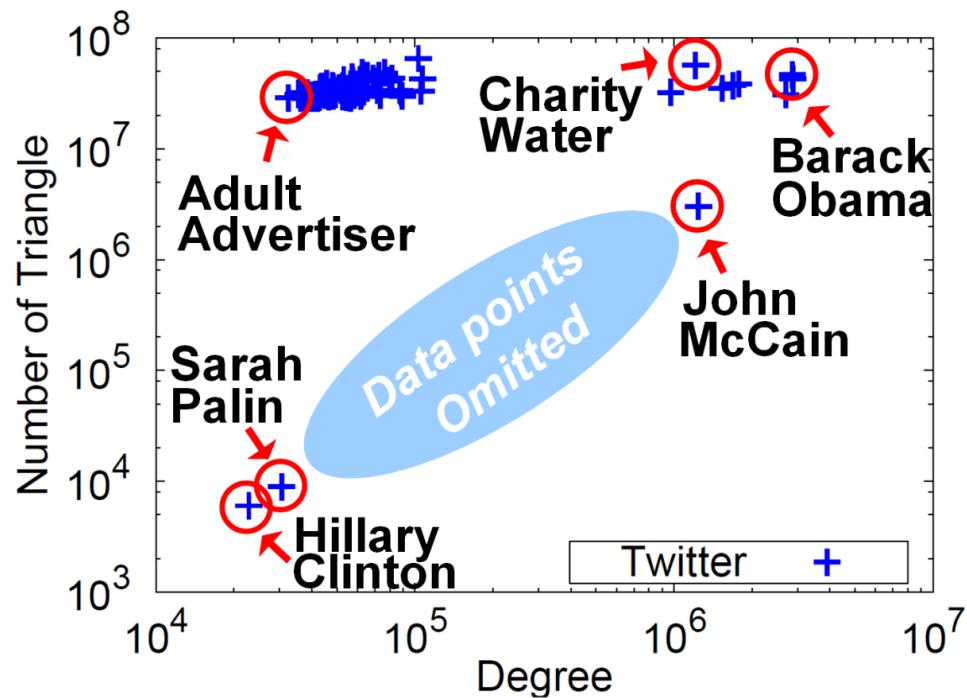
Anomalous nodes in Twitter ( $\sim 3$  billion edges)  
 [U Kang, Brendan Meeder, +, PAKDD'11]

# Triangle Counting on Big Graphs



Q: How to compute # triangles in B-node graph? ( $O(d_{\max}^2)$ )?

# Triangle Counting on Big Graphs



Q: How to compute # triangles in B-node graph? A: **cubes of eigvals**

# Roadmap

- Patterns in graphs
  - overview
  - Static graphs
    - S1: Degree, S2: Eigenvalues
    - S3-4: Triangles, S5: Cliques
    - Radius plot
    - Other observations ('EigenSpokes')
  - Weighted graphs
  - Time-evolving graphs



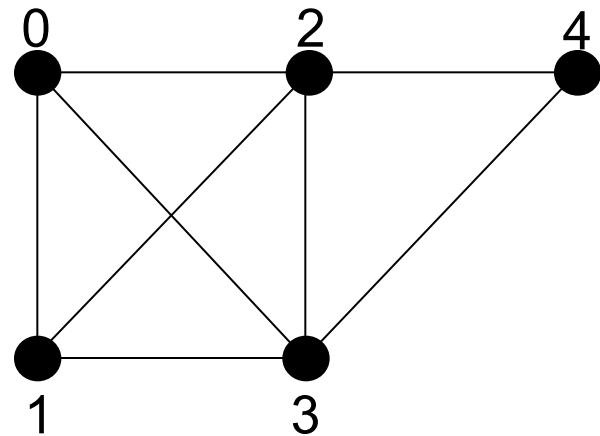
# How about cliques?

- *Large Human Communication Networks Patterns and a Utility-Driven Generator.*  
**Nan Du**, Christos Faloutsos, Bai Wang,  
Leman Akoglu, KDD 2009



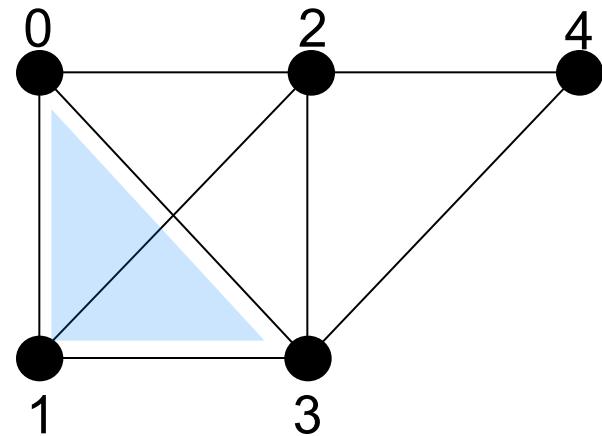
# Cliques

- Clique is a complete subgraph.
- If a clique can not be contained by any larger clique, it is called the maximal clique.



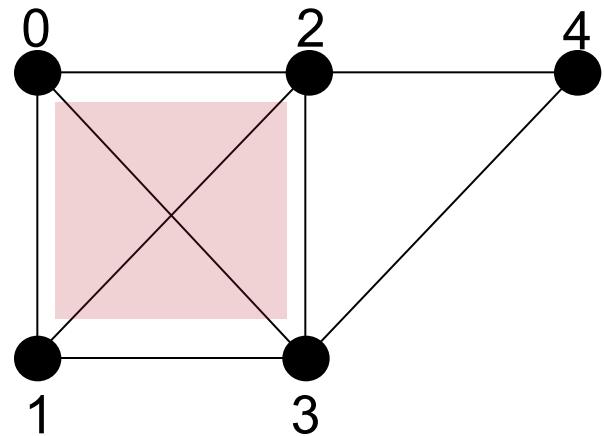
# Clique

- Clique is a complete subgraph.
- If a clique can not be contained by any larger clique, it is called the maximal clique.



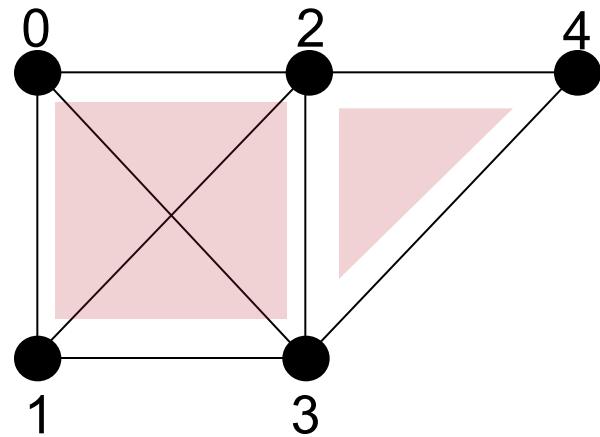
# Clique

- Clique is a complete subgraph.
- If a clique can not be contained by any larger clique, it is called the maximal clique.



# Clique

- Clique is a complete subgraph.
- If a clique can not be contained by any larger clique, it is called the maximal clique.
- $\{0,1,2\}$ ,  $\{0,1,3\}$ ,  $\{1,2,3\}$   $\{2,3,4\}$ ,  $\{0,1,2,3\}$  are cliques;
- $\{\mathbf{0,1,2,3}\}$  and  $\{\mathbf{2,3,4}\}$  are the maximal cliques.



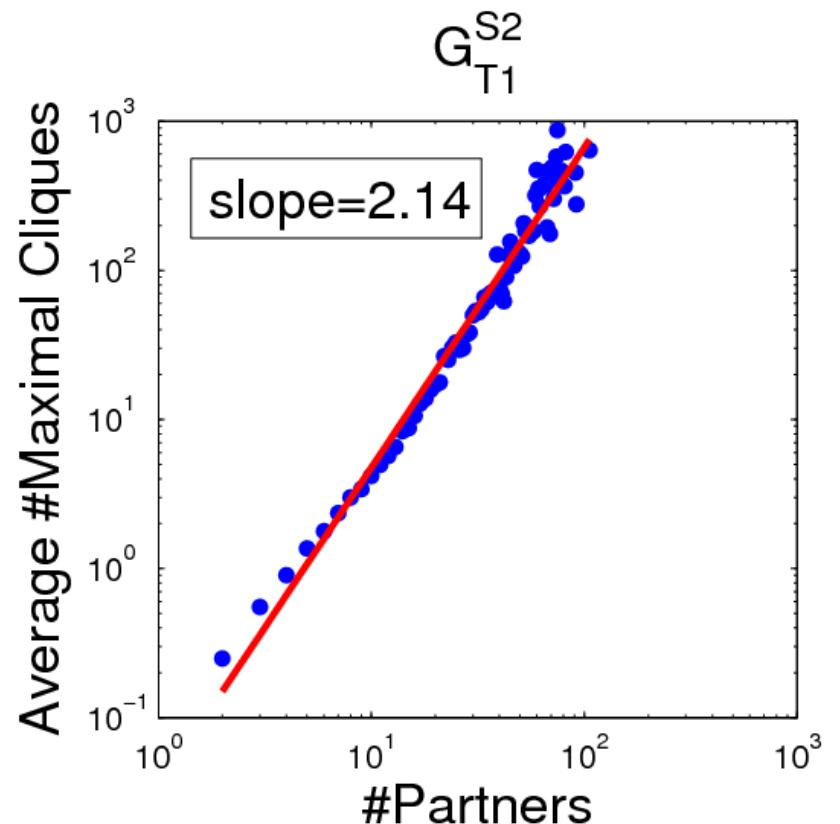
## S5: Clique-Degree Power-Law

- Power law:

$$C_{avg}^{d_i} \propto d_i^\alpha$$

# maximal cliques of node i      degree of node i

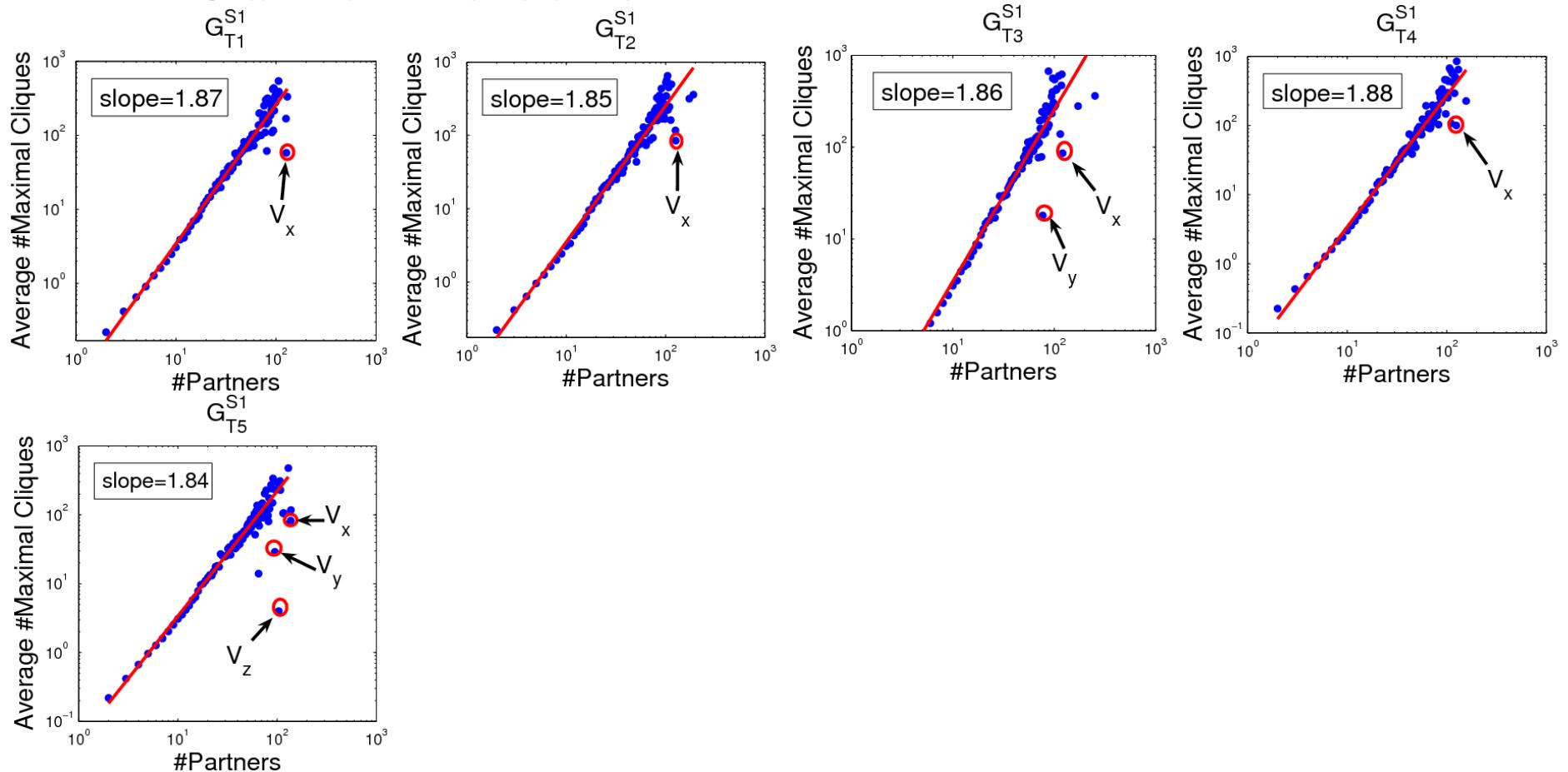
$\alpha$  is the power law exponent  
 $\alpha \in [1.8, 2.2]$  for S1~S3



*More friends, even more social circles !*

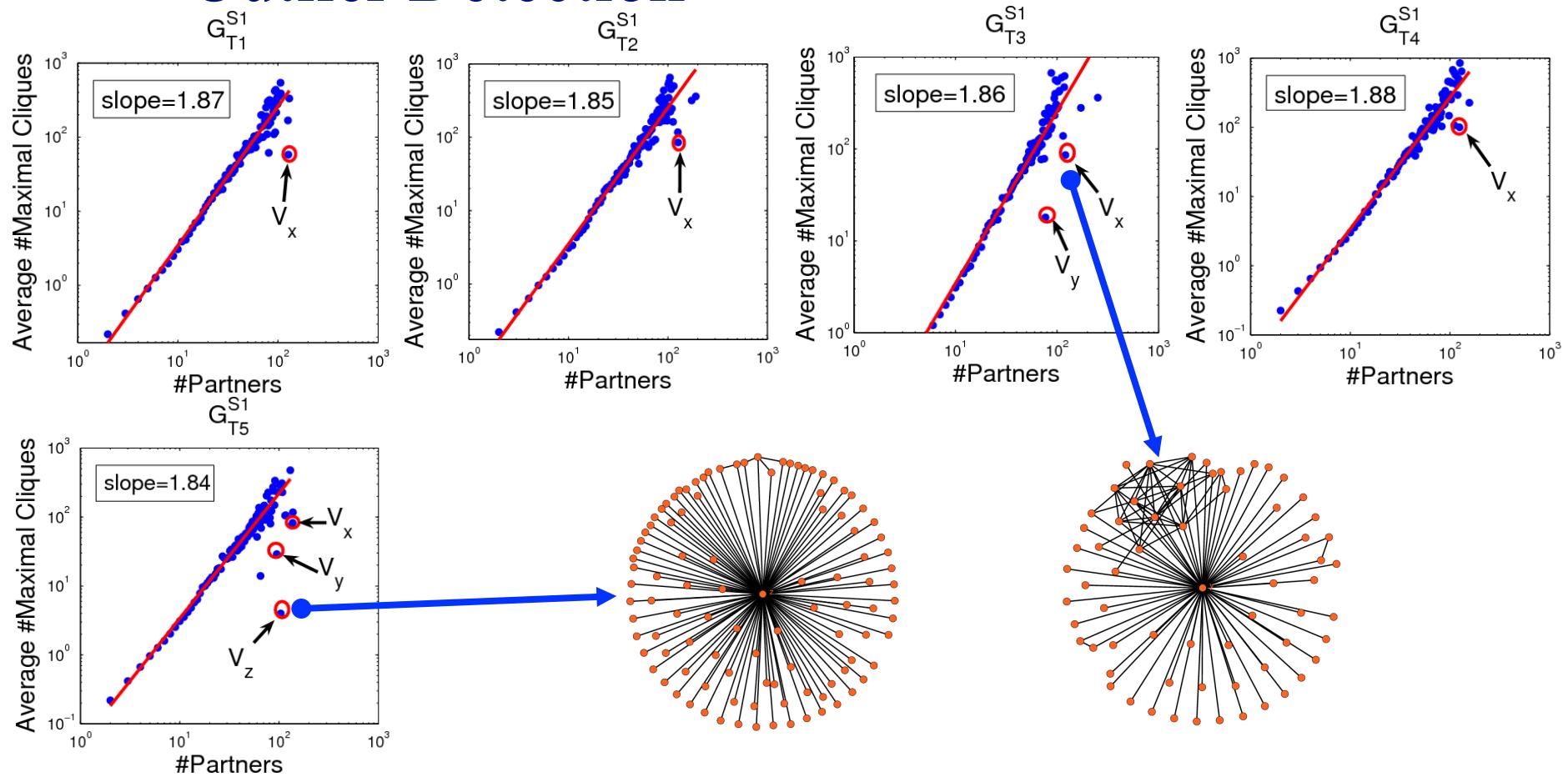
## S5: Clique-Degree Power-Law

- Outlier Detection



## S5: Clique-Degree Power-Law

- Outlier Detection



# Roadmap

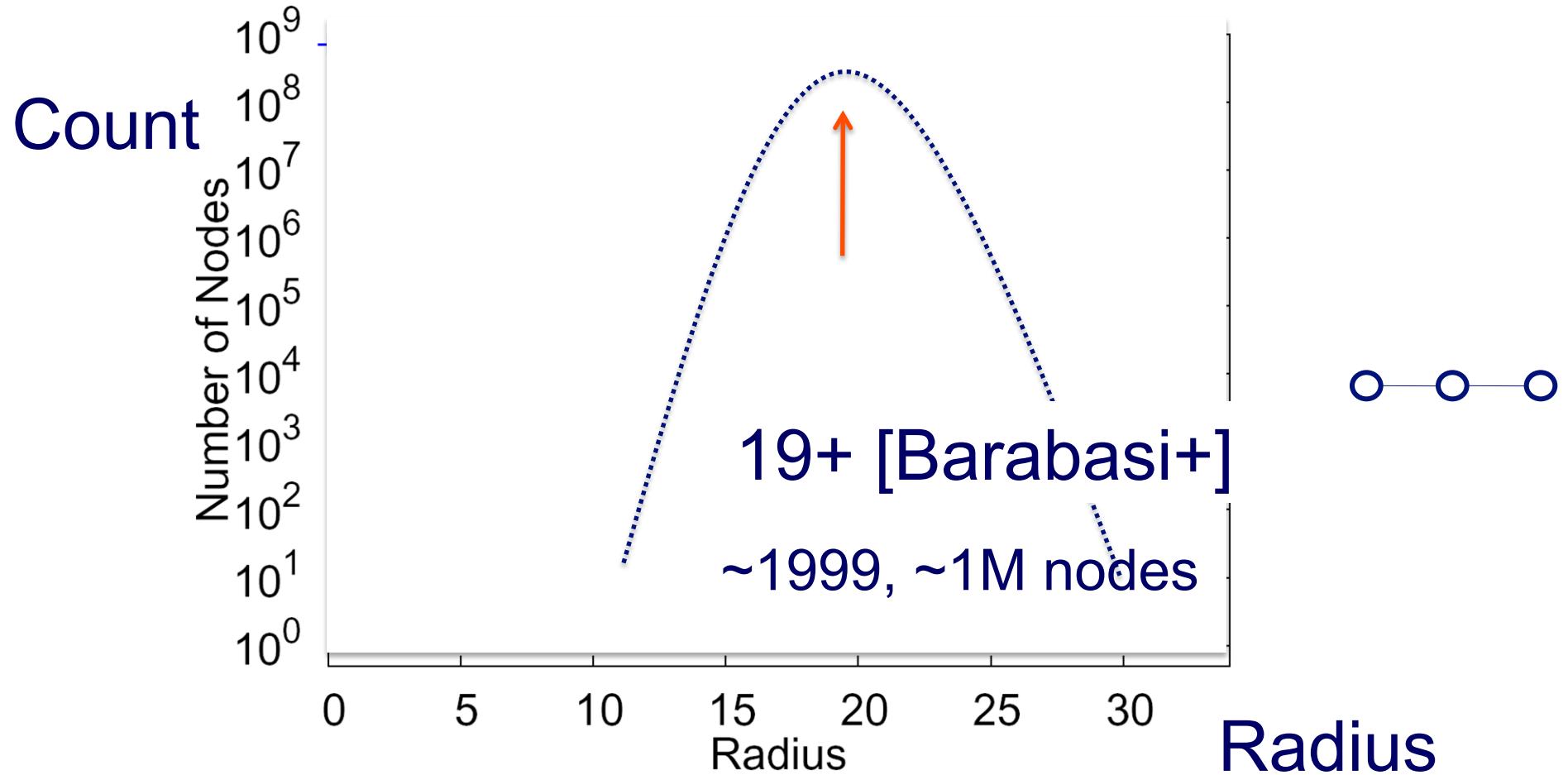
- Patterns in graphs
  - overview
  - Static graphs
    - S1: Degree, S2: Eigenvalues
    - S3-4: Triangles, S5: Cliques
    - Radius plot
    - Other observations ('EigenSpokes')
  - Weighted graphs
  - Time-evolving graphs

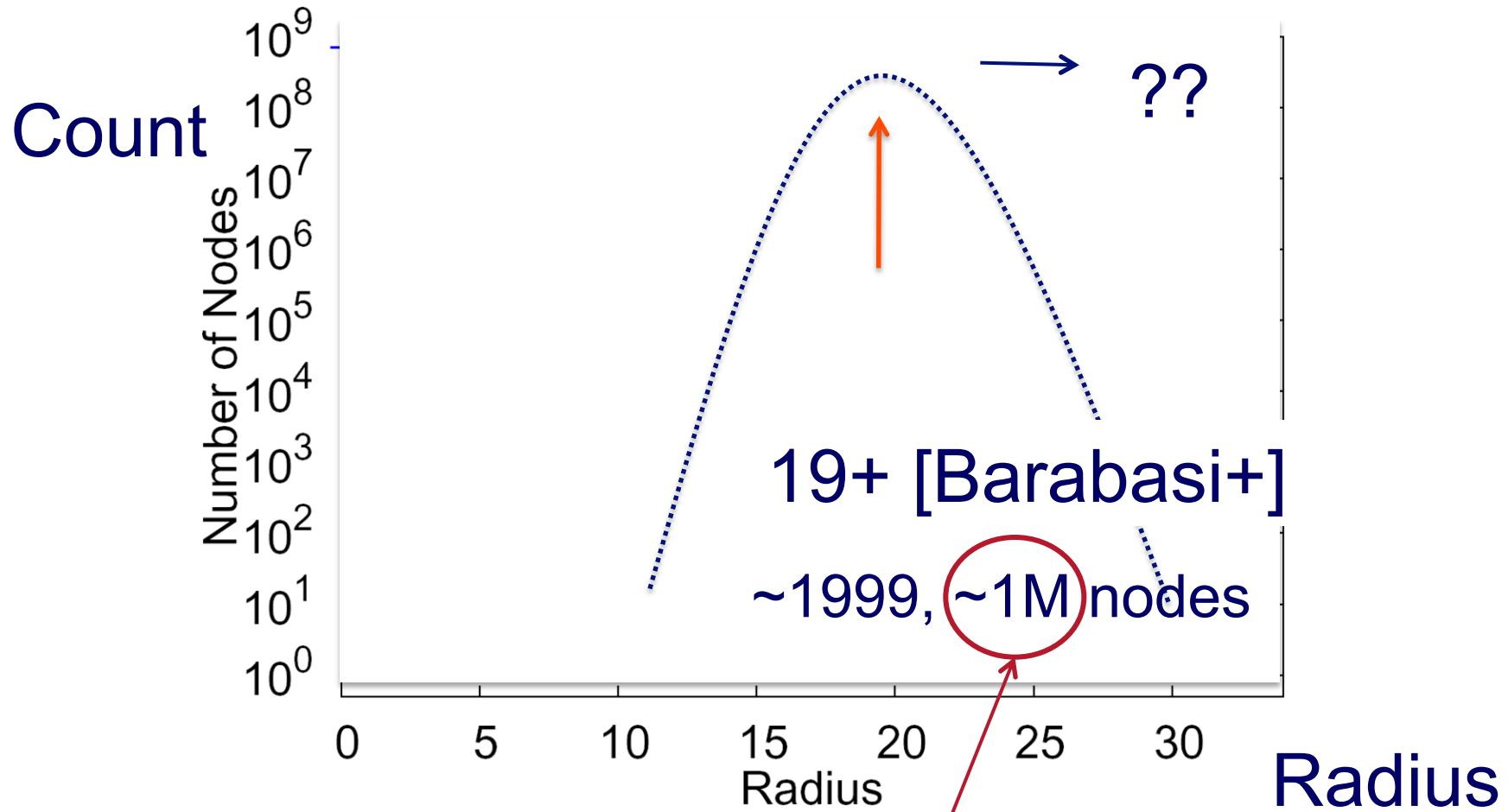




# HADI for Diameter Estimation

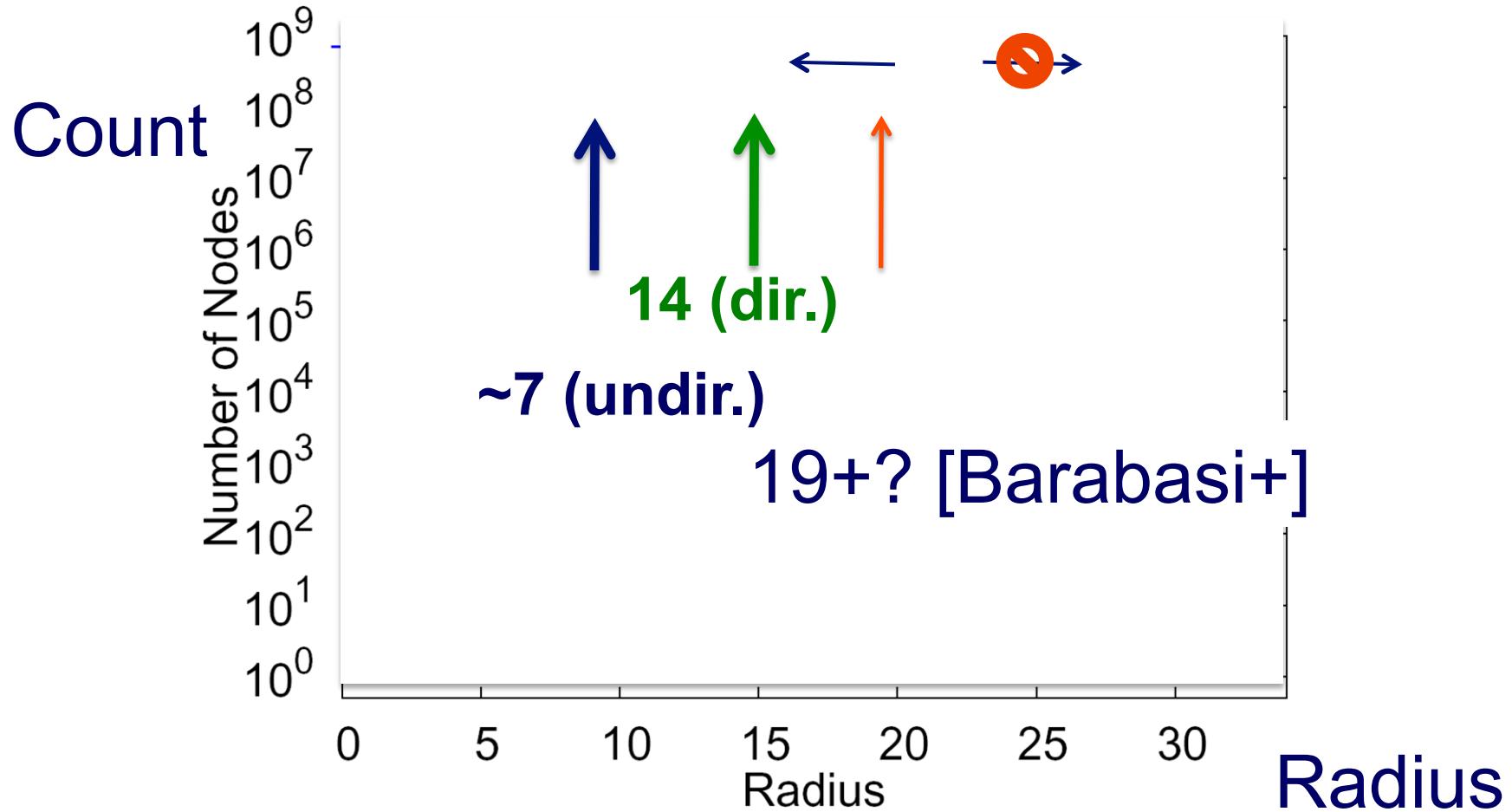
- *Radius Plots for Mining Tera-byte Scale Graphs*  
**U Kang**, Charalampos Tsourakakis, Ana Paula Appel, Christos Faloutsos, Jure Leskovec,  
SDM'10
- Naively: diameter needs  $O(N^2)$  space and up to  $O(N^3)$  time – **prohibitive** ( $N \sim 1B$ )
- Our HADI: linear on  $E$  ( $\sim 10B$ )
  - Near-linear scalability w.r.t. # machines
  - Several optimizations → 5x faster





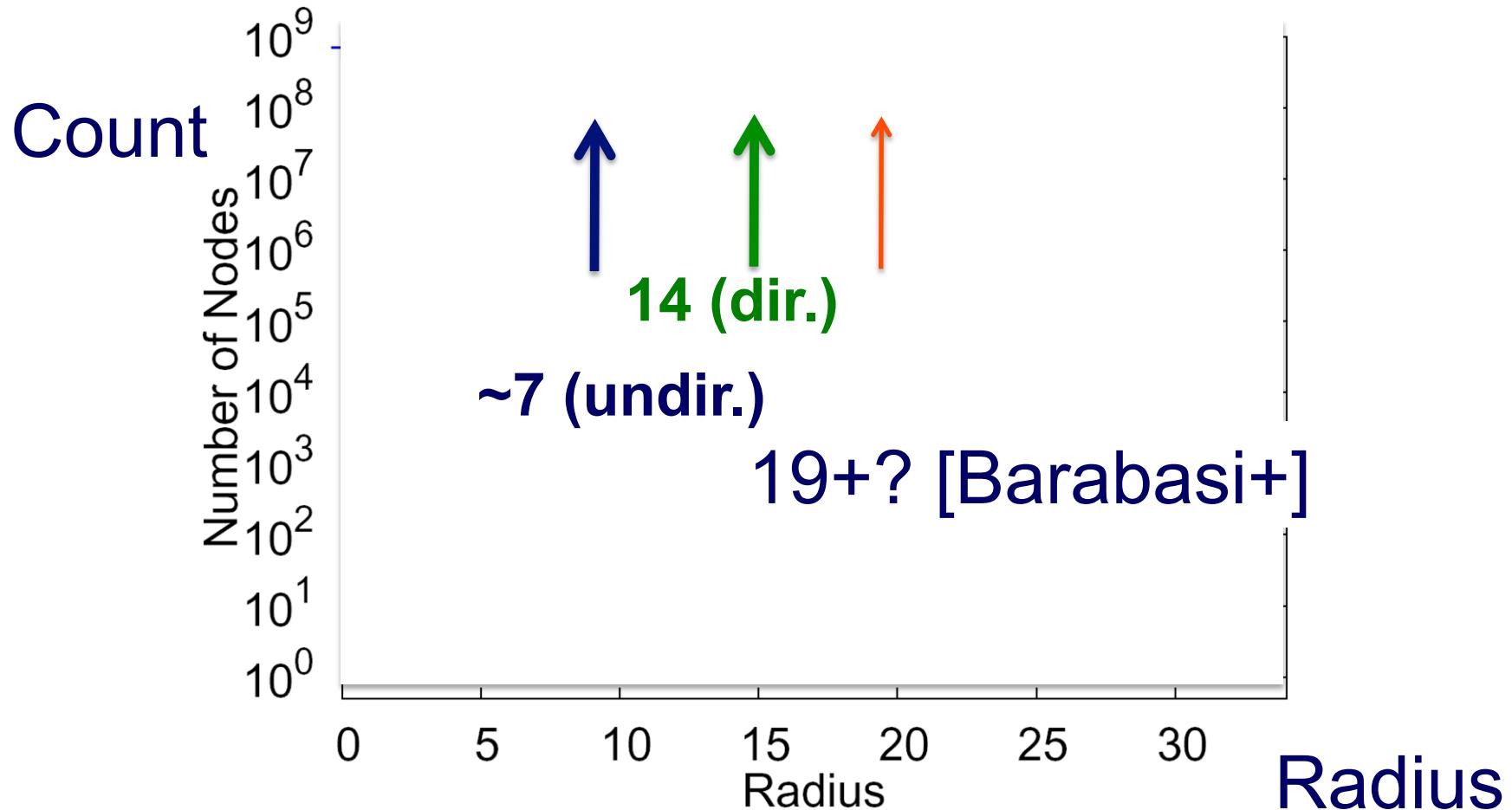
YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

- Largest publicly available graph ever studied.



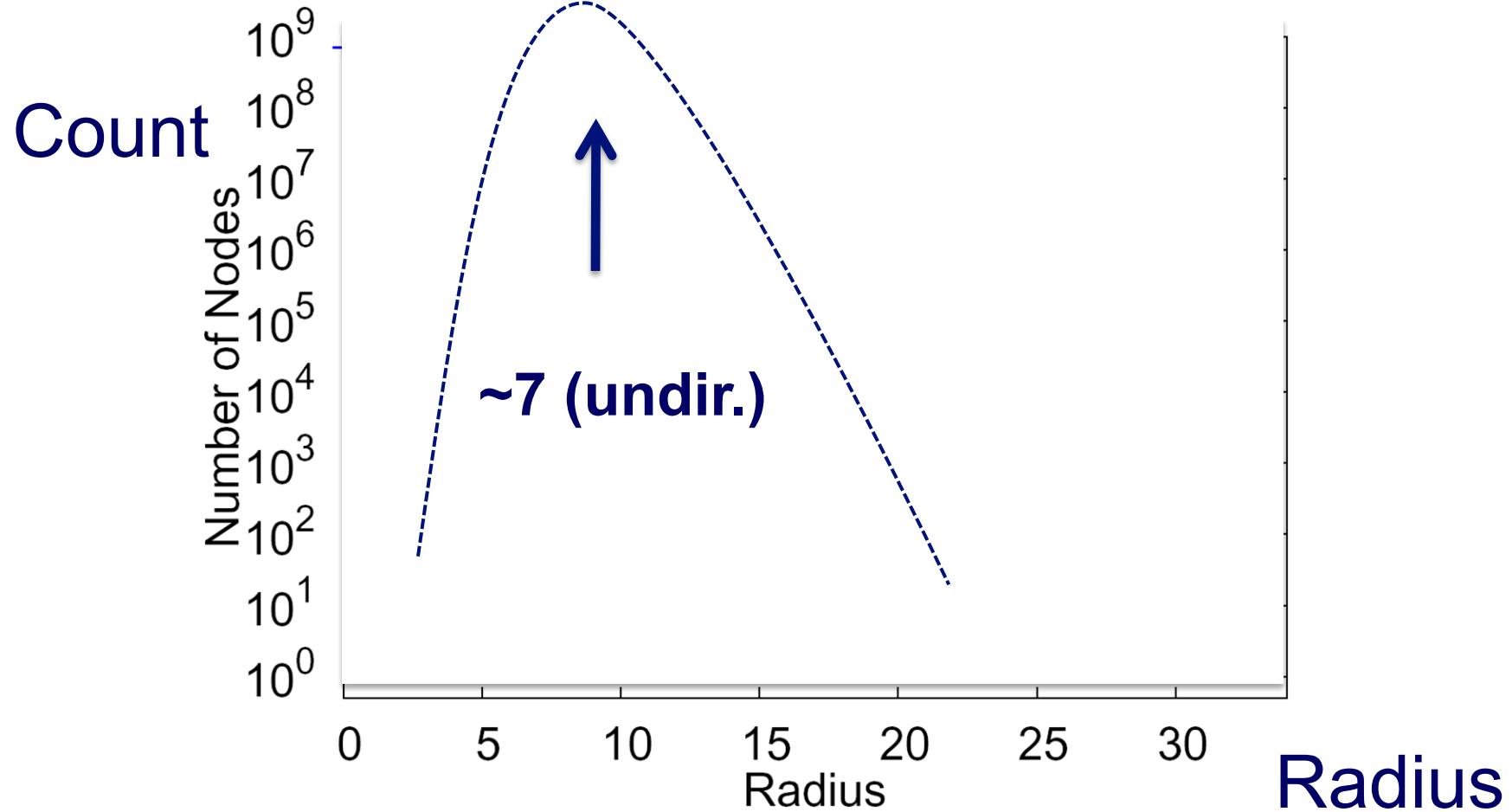
YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

- Largest publicly available graph ever studied.

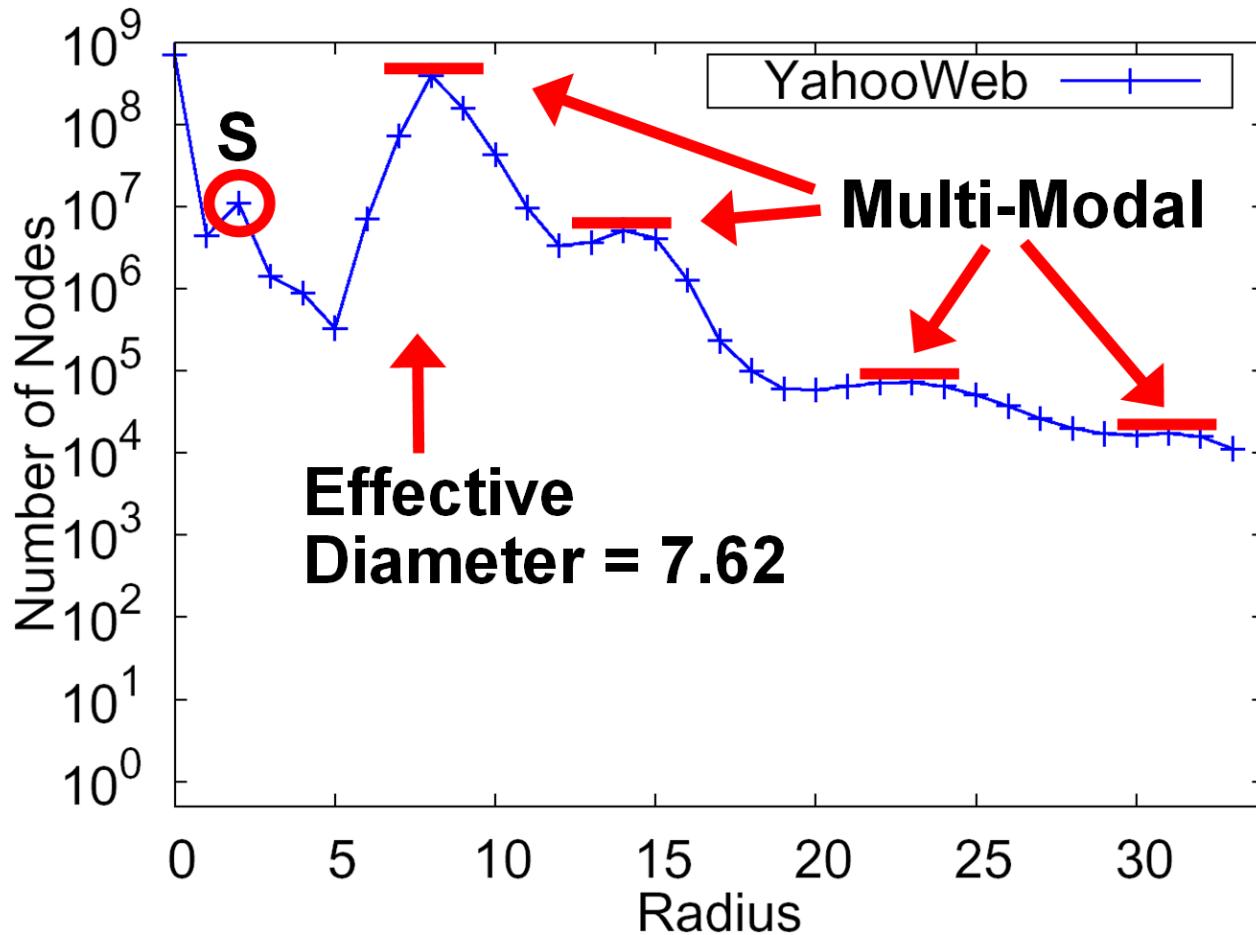


YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

- 7 degrees of separation (!)
- Diameter: shrunk

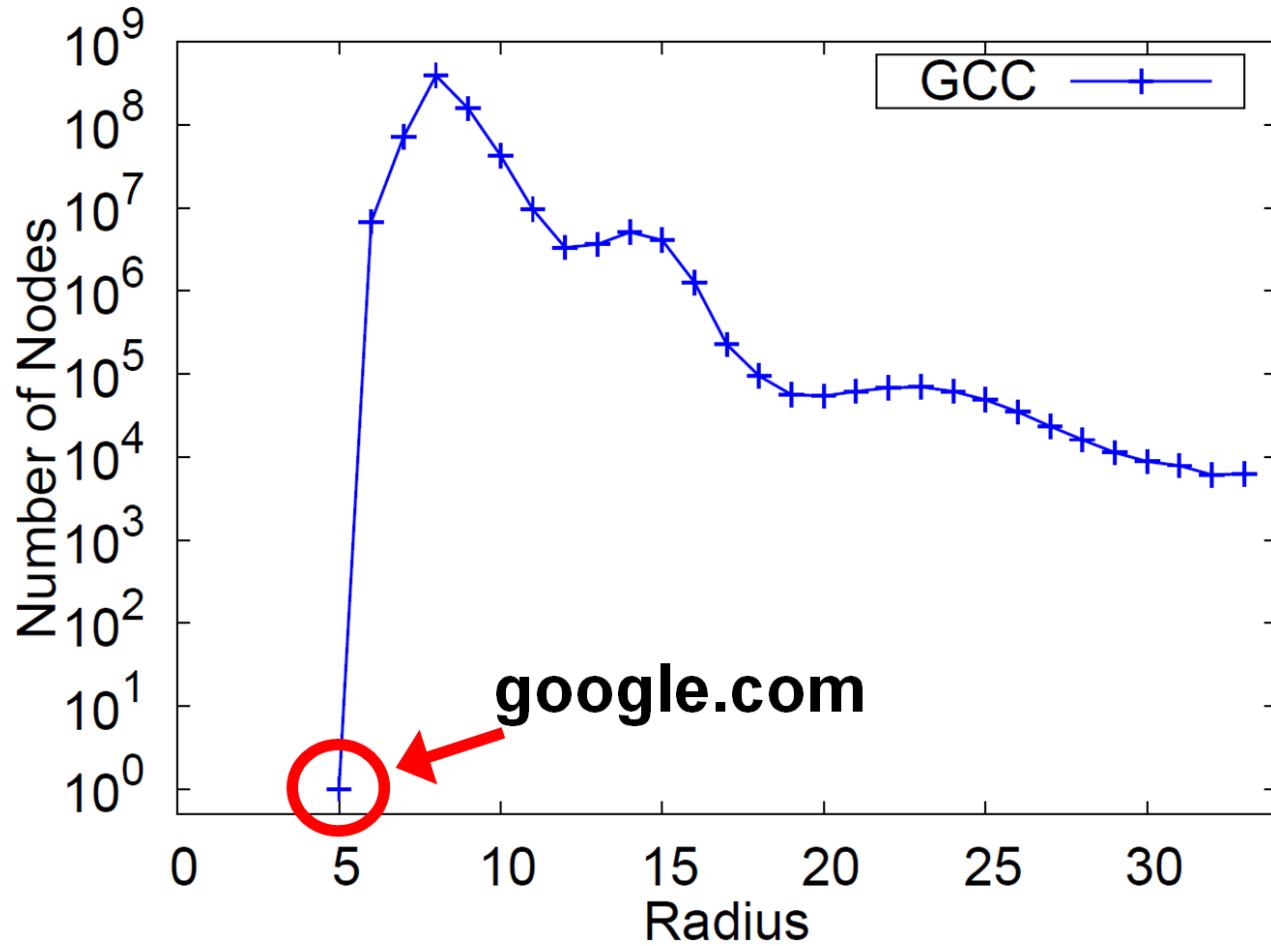


YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)  
Q: Shape?

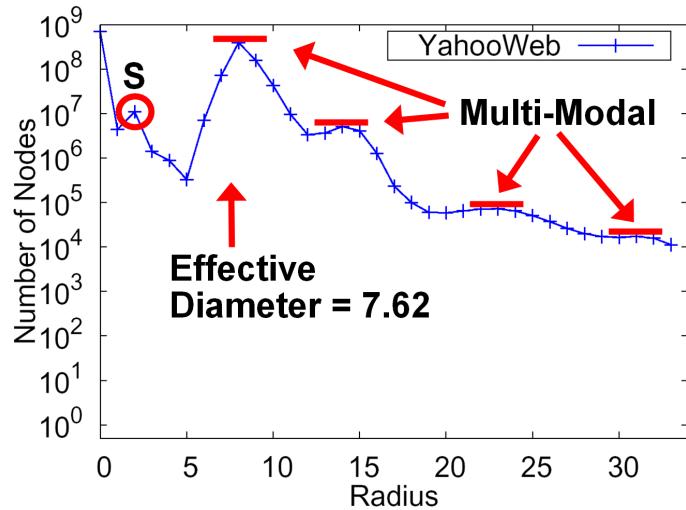


YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

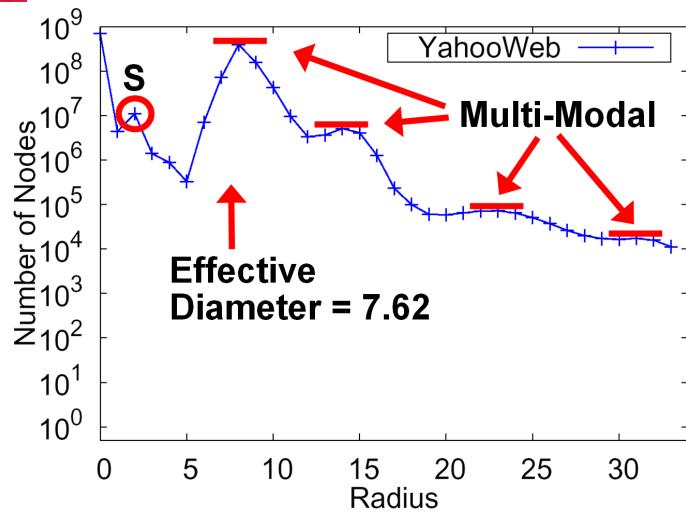
- Effective diameter: surprisingly small
- Multi-modality (?!)



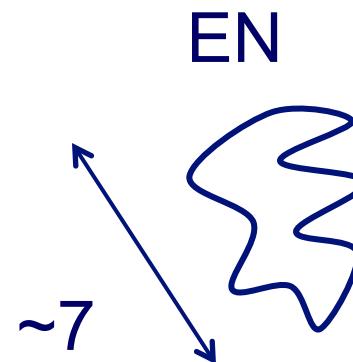
## Radius Plot of **GCC** of YahooWeb Graph



- YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)
- Effective diameter: surprisingly small
  - Multi-modality: probably mixture of cores

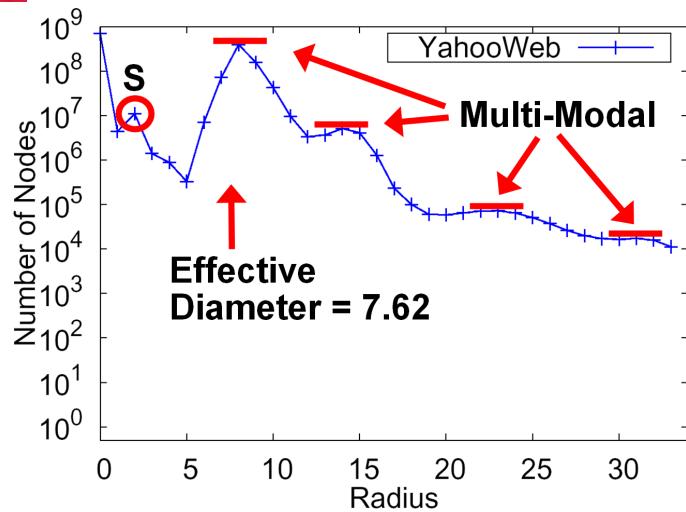


Conjecture:

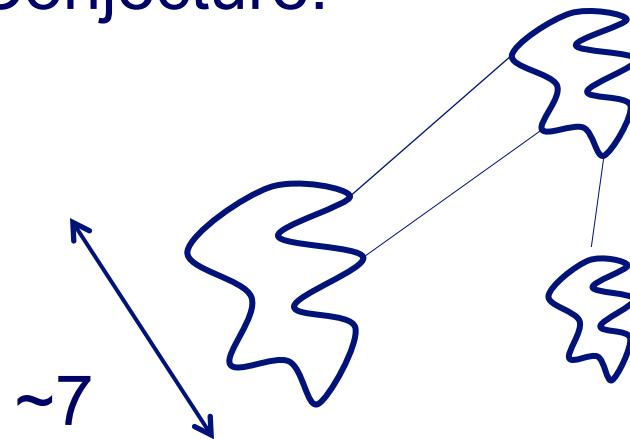


YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

- Effective diameter: surprisingly small
- Multi-modality: probably mixture of cores



Conjecture:



YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)

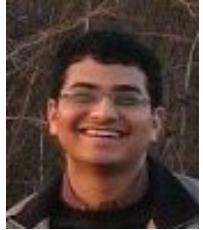
- effective diameter: surprisingly small
- Multi-modality: probably mixture of cores

# Roadmap

- Patterns in graphs
  - overview
  - Static graphs
    - S1: Degree, S2: Eigenvalues
    - S3-4: Triangles, S5: Cliques
    - Radius plot
    - Other observations ('EigenSpokes')
  - Weighted graphs
  - Time-evolving graphs



# S6: EigenSpokes



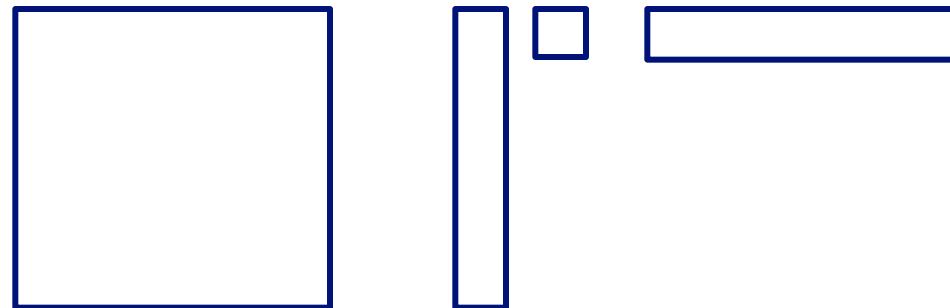
*EigenSpokes: Surprising Patterns and Scalable Community Chipping in Large Graphs.*

**B. Aditya Prakash**, Mukund Seshadri,  
Ashwin Sridharan, Sridhar Machiraju, and  
Christos Faloutsos: PAKDD 2010.

# EigenSpokes

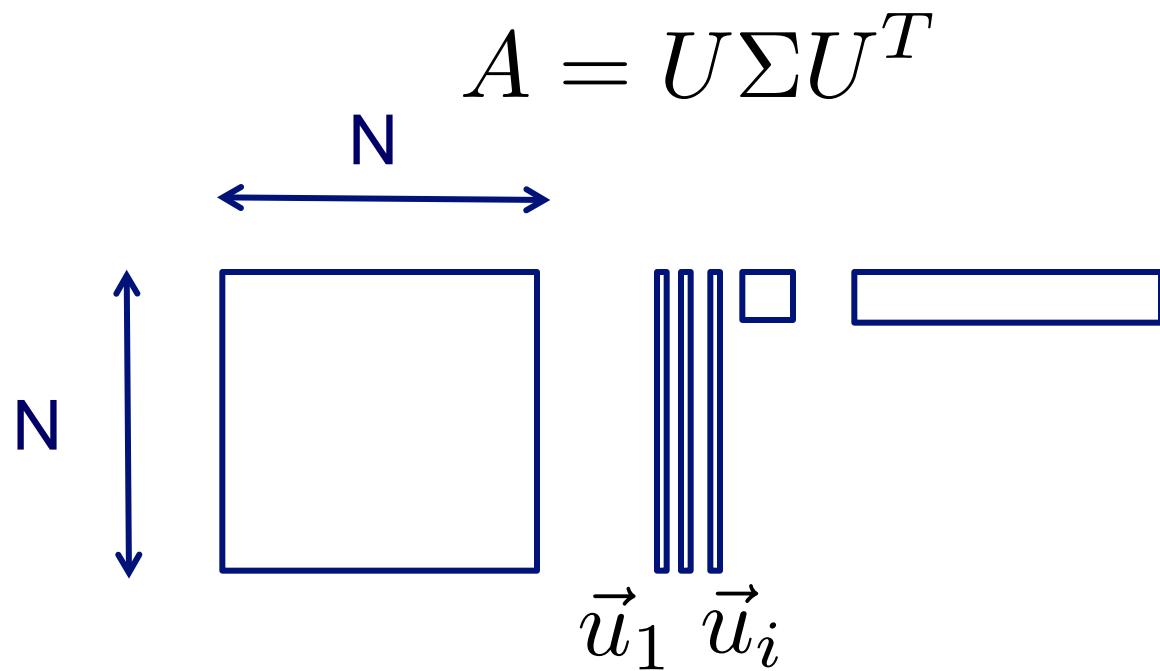
- Eigenvectors of adjacency matrix
  - equivalent to singular vectors (symmetric, undirected graph)

$$A = U\Sigma U^T$$



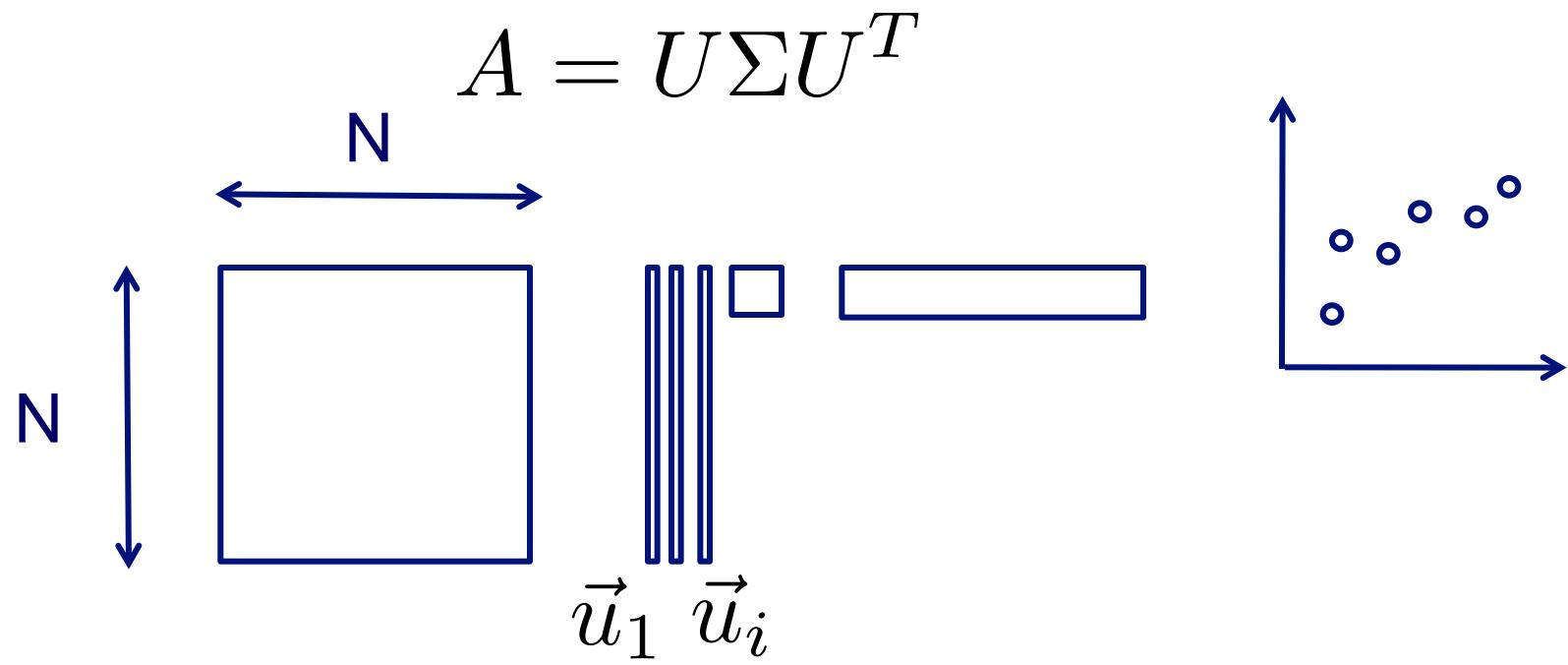
# EigenSpokes

- Eigenvectors of adjacency matrix
  - equivalent to singular vectors (symmetric, undirected graph)



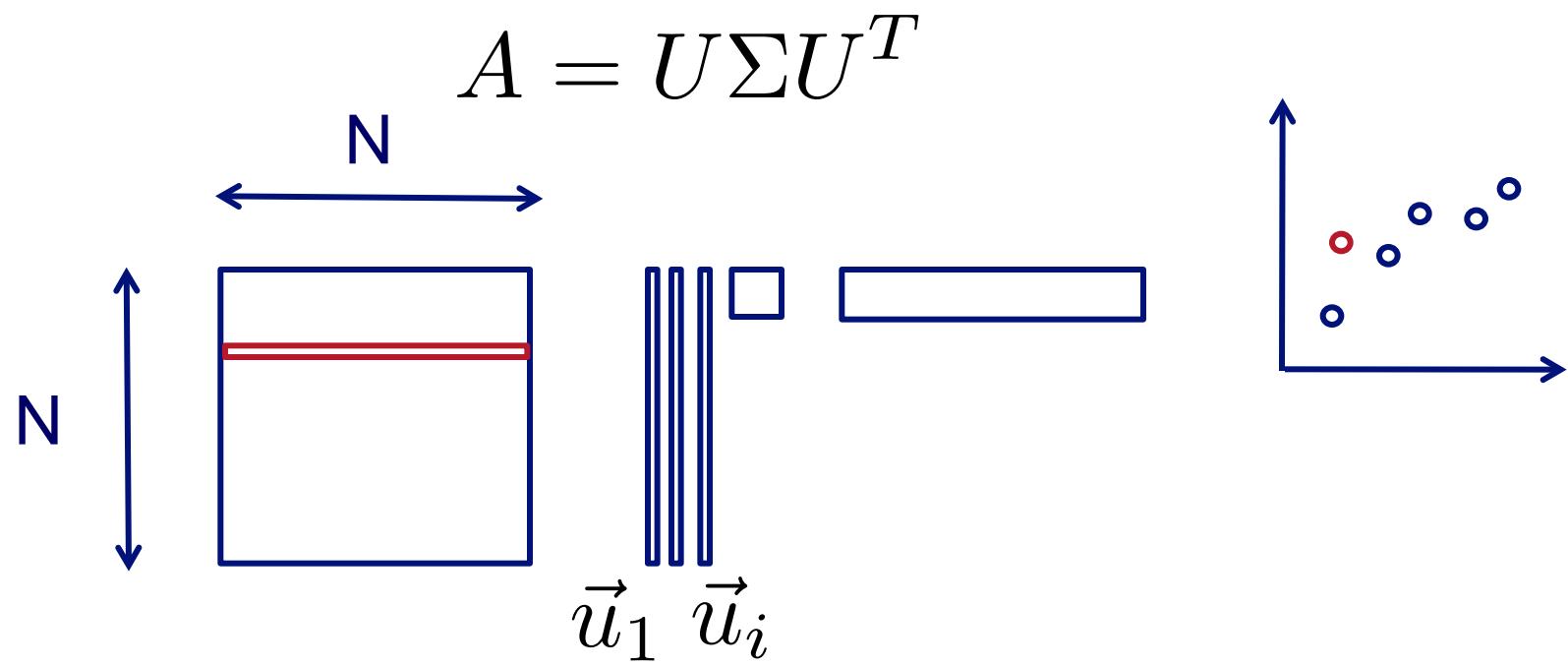
# EigenSpokes

- Eigenvectors of adjacency matrix
  - equivalent to singular vectors  
(symmetric, undirected graph)



# EigenSpokes

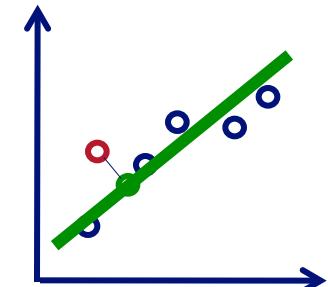
- Eigenvectors of adjacency matrix
  - equivalent to singular vectors (symmetric, undirected graph)



# EigenSpokes

- Eigenvectors of adjacency matrix
  - equivalent to singular vectors  
(symmetric, undirected graph)

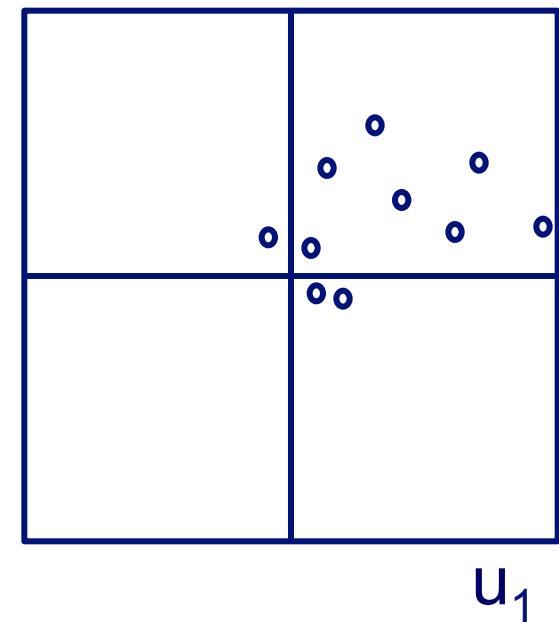
$$A = U\Sigma U^T$$



# EigenSpokes

- EigenEigen (EE) Plot
  - Scatter plot of scores of  $u_1$  vs.  $u_2$
- One would expect
  - Many points @ origin
  - A few scattered ~randomly

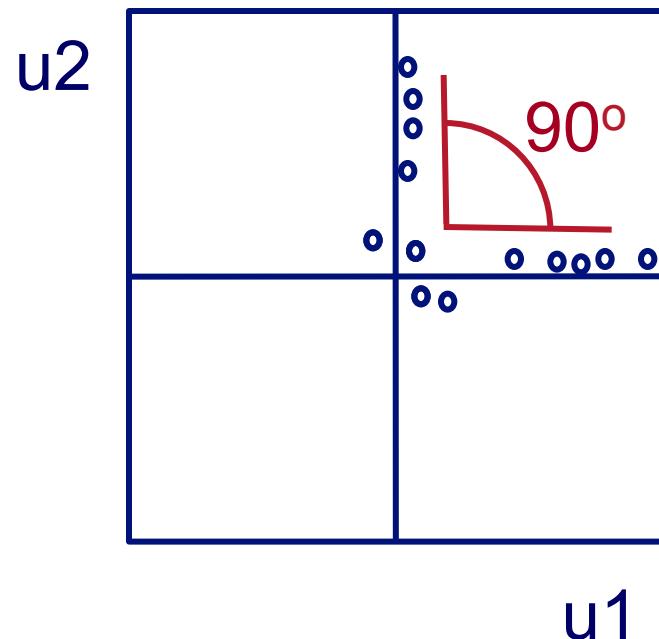
2<sup>nd</sup> Principal component  
 $u_2$



1<sup>st</sup> Principal component

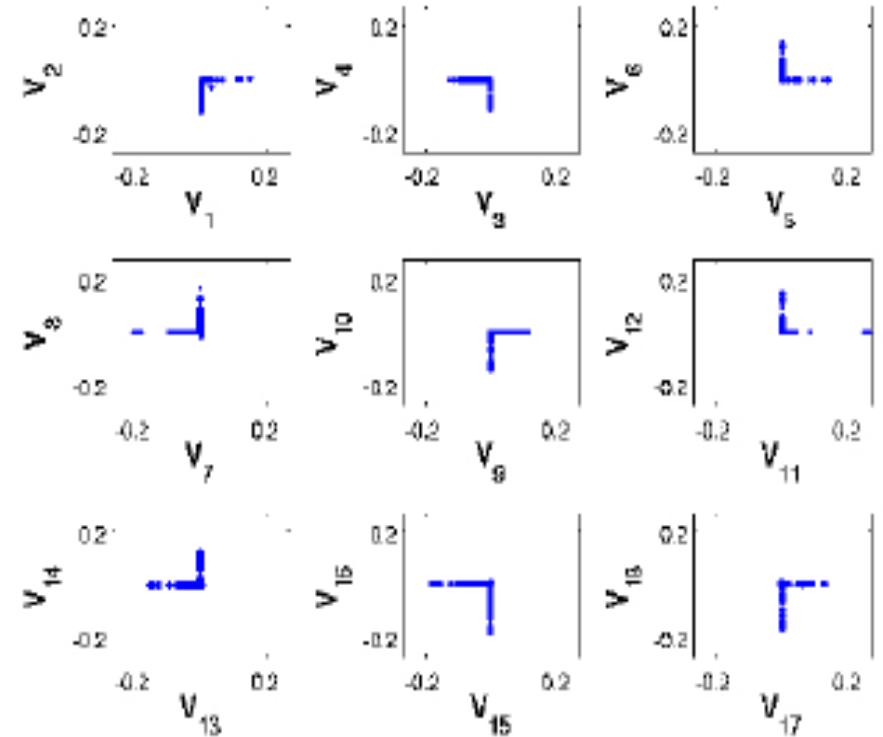
# EigenSpokes

- EigenEigen (EE) Plot
  - Scatter plot of scores of  $u_1$  vs  $u_2$
- One would expect
  - Many points @ origin
  - A few scattered ~randomly



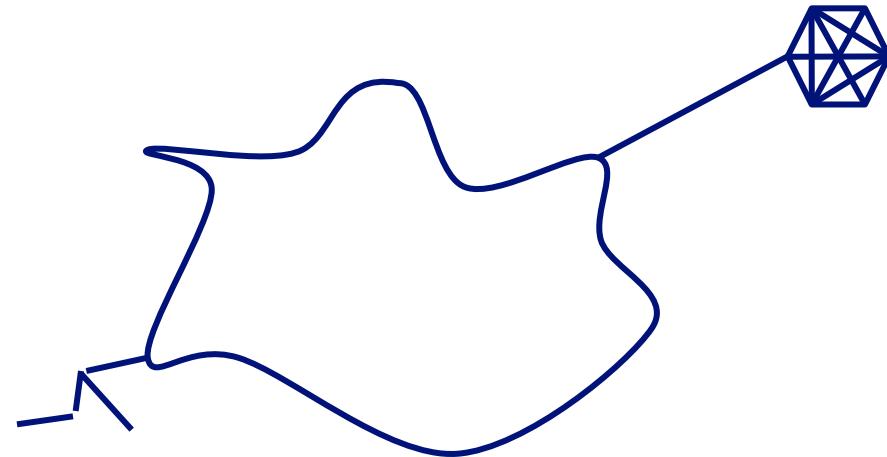
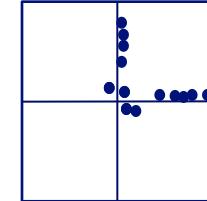
# EigenSpokes - pervasiveness

- Present in mobile social graph
  - Across time and space
- Patent citation graph



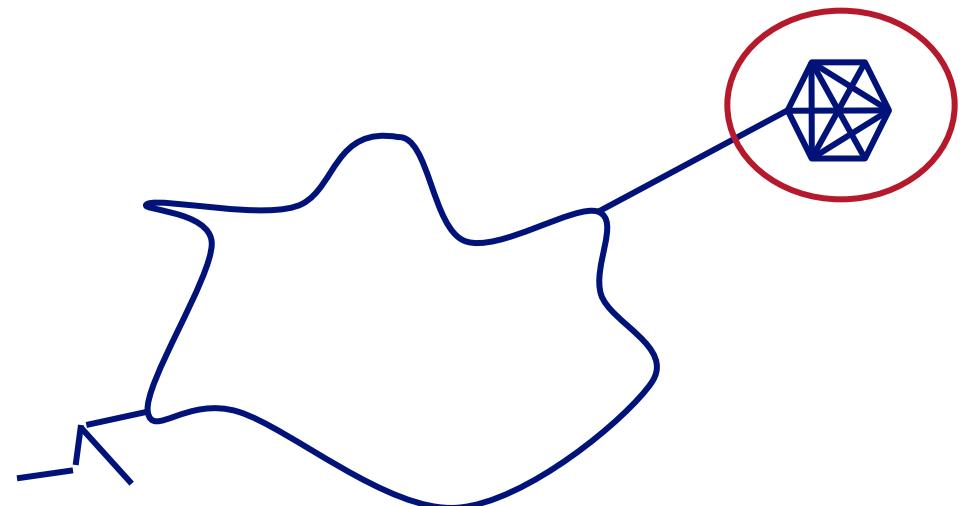
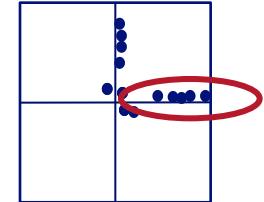
# EigenSpokes - explanation

Near-cliques, or near-bipartite-cores, loosely connected



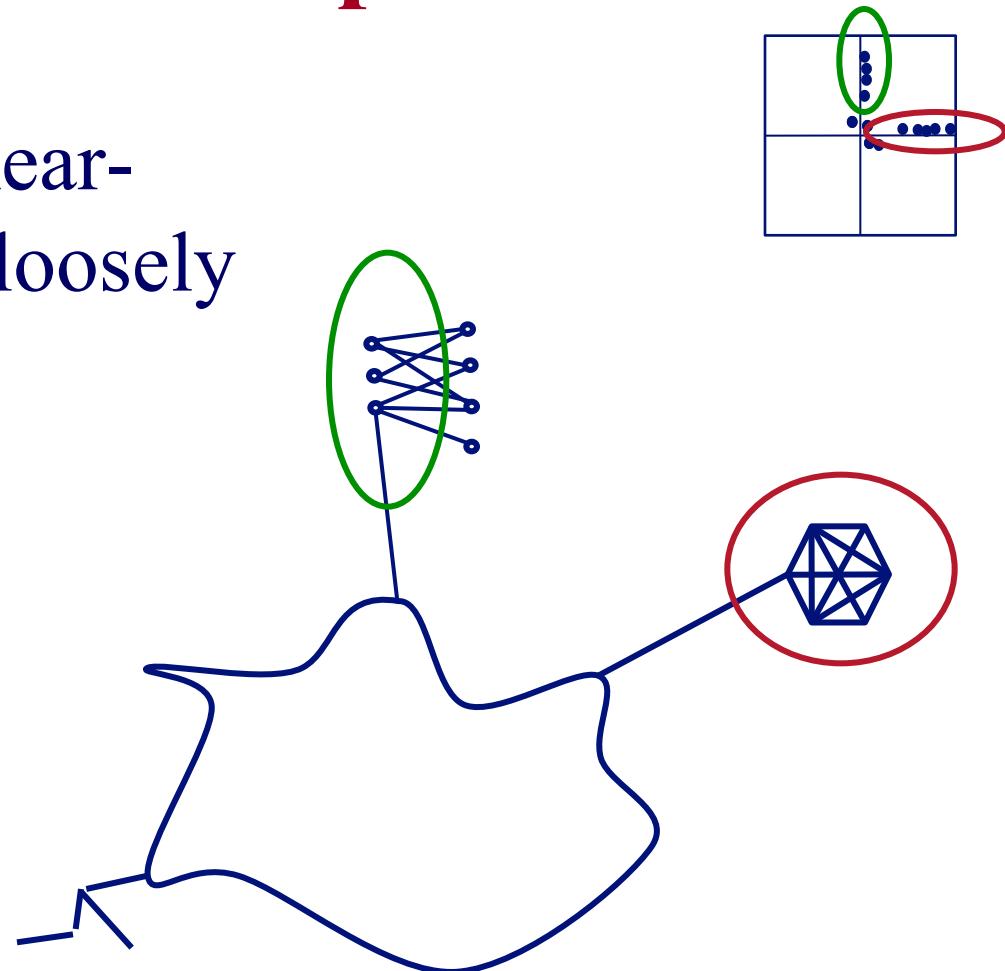
# EigenSpokes - explanation

Near-cliques, or near-bipartite-cores, loosely connected



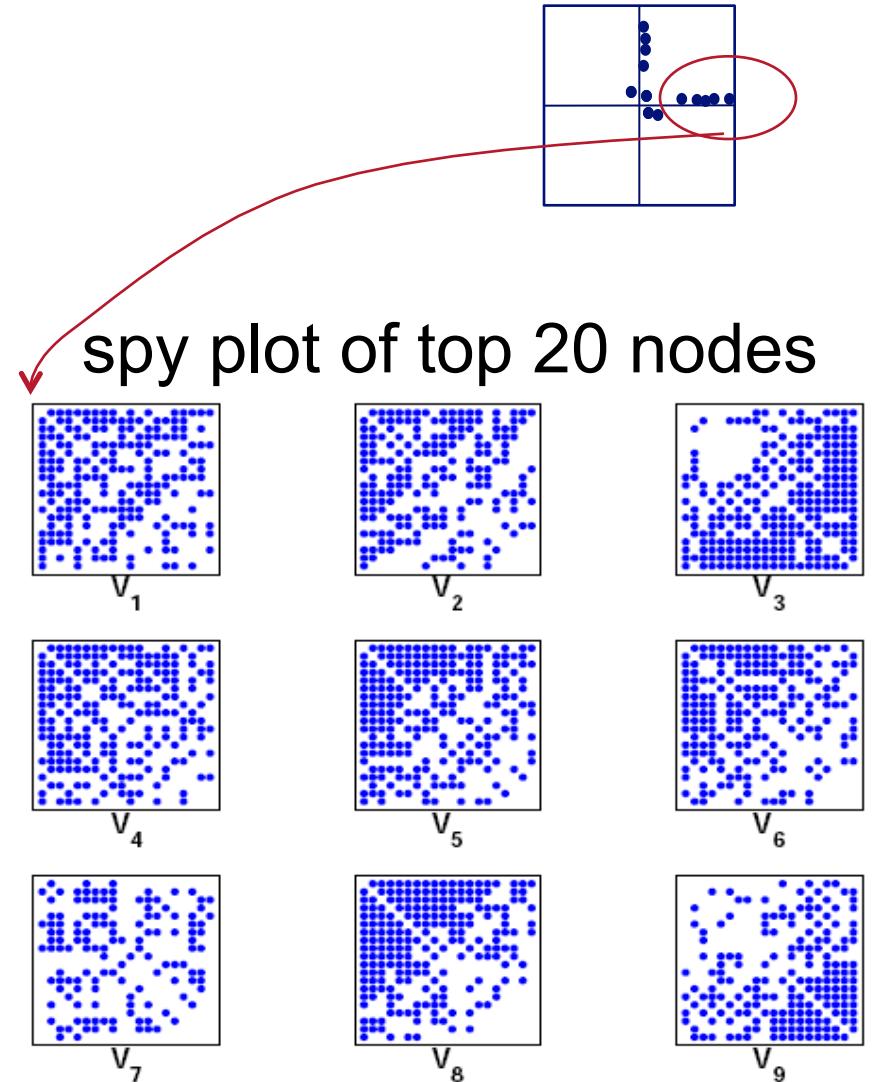
# EigenSpokes - explanation

Near-cliques, or near-bipartite-cores, loosely connected



# EigenSpokes - explanation

Near-cliques, or near-bipartite-cores, loosely connected



So what?

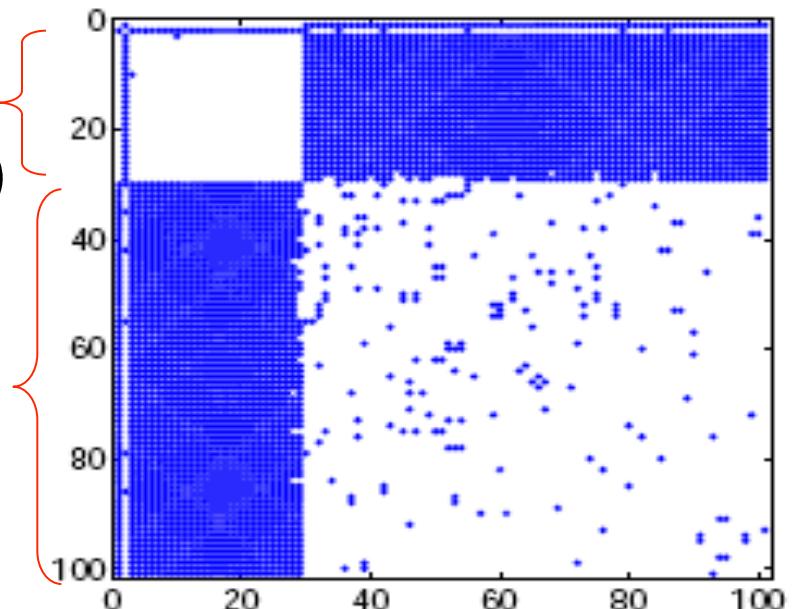
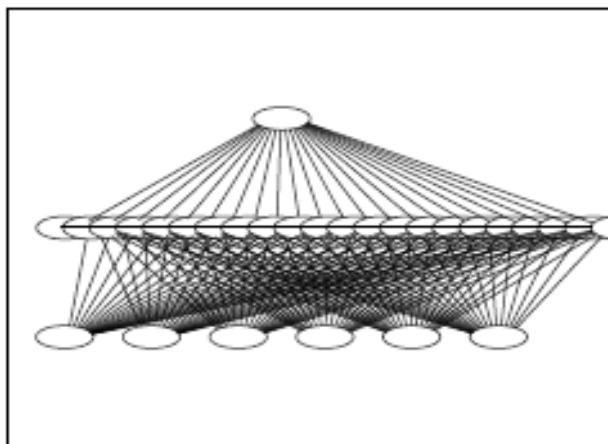
- Extract nodes with high *scores*
- High connectivity
- Good “communities”

# Bipartite Communities!

patents from  
same inventor(s)

‘cut-and-paste’  
bibliography!

magnified bipartite community



# Roadmap

- Patterns in graphs
  - Overview
  - Static graphs
  - Weighted graphs
  - Time-evolving graphs
- Anomaly Detection
- Application: ebay fraud
- Conclusions



# Observations on Weighted Graphs

- A: yes - even more ‘laws’ !



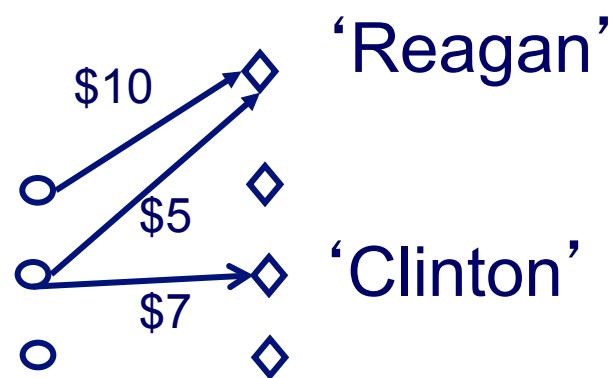
M. McGlohon, L. Akoglu, and C. Faloutsos  
*Weighted Graphs and Disconnected Components:  
Patterns and a Generator*. KDD 2008

# Observation W.1: Fortification

*Q: How do the weights  
of nodes relate to degree?*

# Observation W.1: Fortification

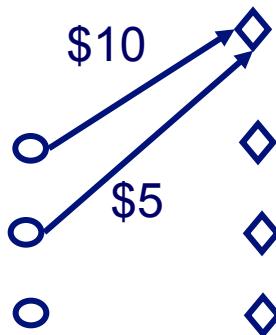
More donors,  
more \$ ?



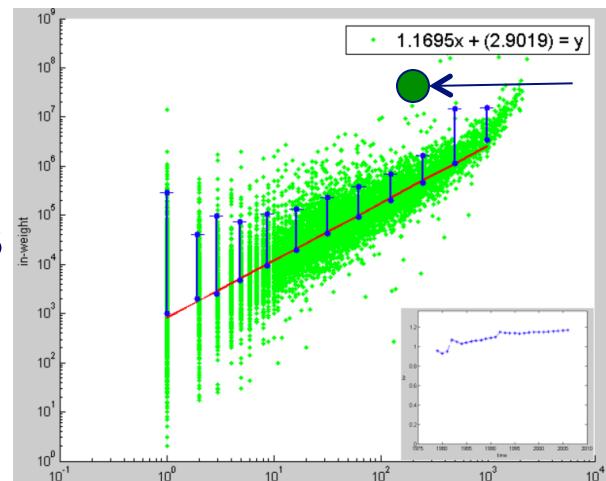
# Observation W.1: fortification: Snapshot Power Law

- Weight: super-linear on in-degree
- exponent ‘iw’ :  $1.01 < iw < 1.26$

More donors,  
even more \$



In-weights  
(\$)



Edges (# donors)

e.g. John Kerry,  
\$10M received,  
from 1K donors

# Roadmap

- Patterns in graphs
  - Overview
  - Static graphs
  - Weighted graphs
  - Time-evolving graphs
- 
- Anomaly Detection
- Application: ebay fraud
- Conclusions



# Problem: Time evolution

- with Jure Leskovec  
(CMU → Stanford)

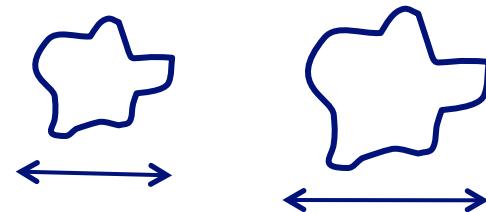


- and Jon Kleinberg  
(Cornell – sabb. @ CMU)



## T.1 Evolution of the Diameter

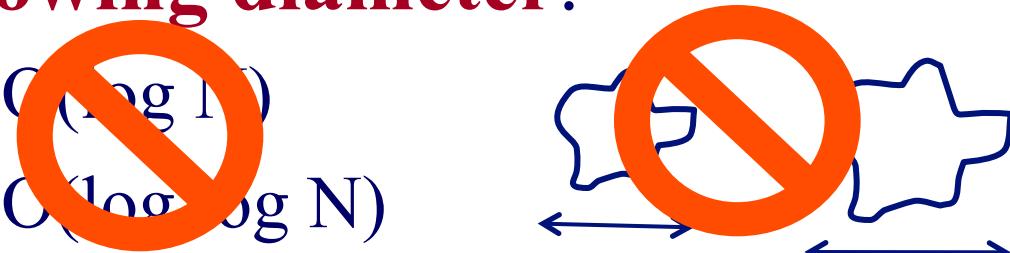
- Prior work on Power Law graphs hinted at **slowly growing diameter**:
  - diameter  $\sim O(\log N)$
  - diameter  $\sim O(\log \log N)$
- What is happening in real data?



## T.1 Evolution of the Diameter

- Prior work on Power Law graphs hints at **slowly growing diameter**:

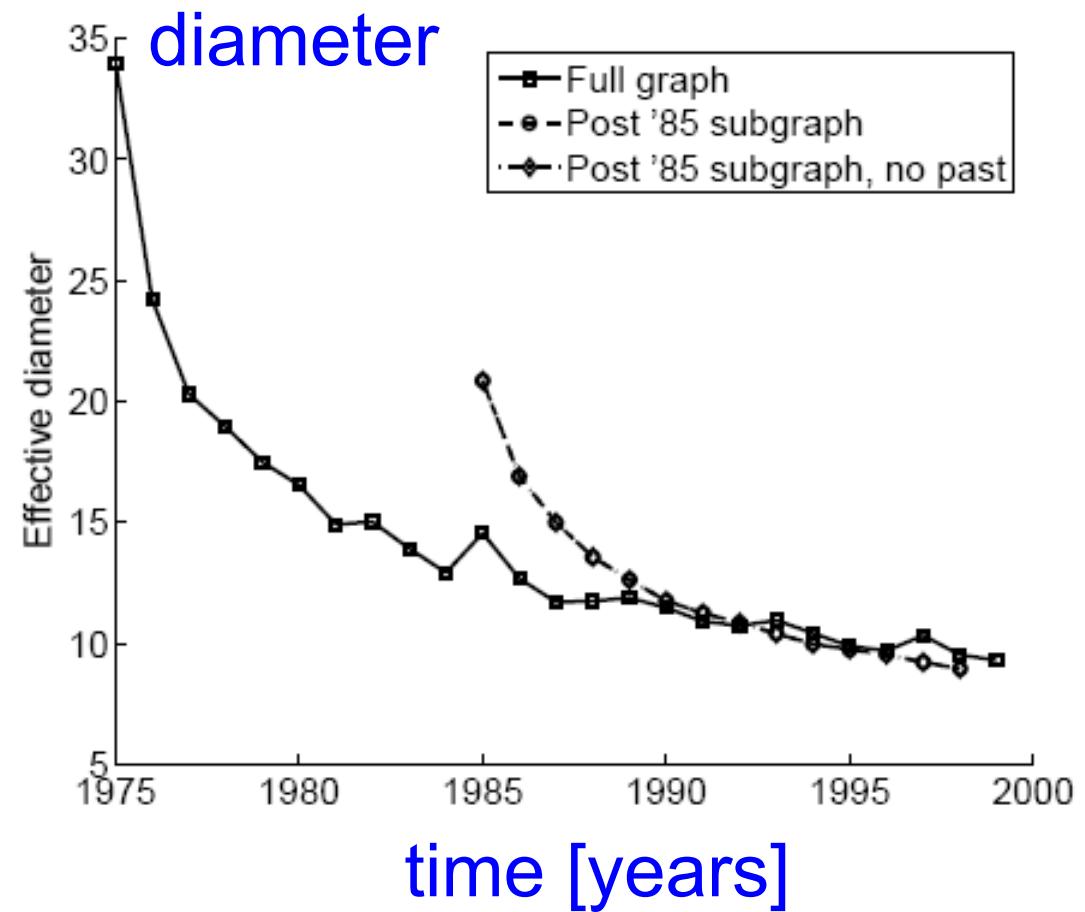
- diameter  $\sim O(\log N)$
  - diameter  $\sim O(10^{\sigma} \log N)$



- What is happening in real data?
- Diameter **shrinks** over time

# T.1 Diameter – “Patents”

- Patent citation network
- 25 years of data
- @1999
  - 2.9 M nodes
  - 16.5 M edges



## T.2 Temporal Evolution of the Graphs

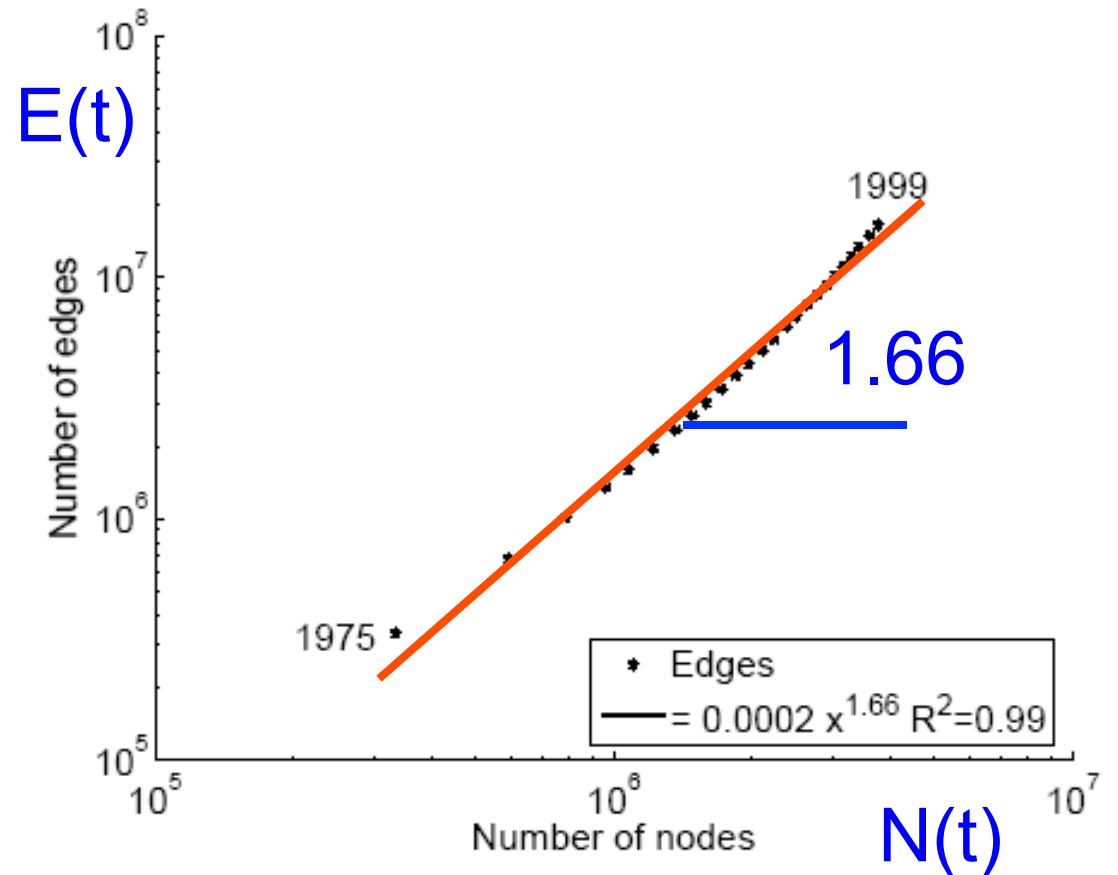
- $N(t)$  ... nodes at time  $t$
- $E(t)$  ... edges at time  $t$
- Suppose that
$$N(t+1) = 2 * N(t)$$
- Q: what is your guess for
$$E(t+1) =? 2 * E(t)$$

## T.2 Temporal Evolution of the Graphs

- $N(t)$  ... nodes at time t
- $E(t)$  ... edges at time t
- Suppose that
$$N(t+1) = 2 * N(t)$$
- Q: what is your guess for
$$E(t+1) = ? \cdot 2 * E(t)$$
- A: over-doubled!
  - But obeying the “Densification Power Law”

## T.2 Densification – Patent Citations

- Citations among patents granted
- @1999
  - 2.9 M nodes
  - 16.5 M edges
- Each year is a datapoint



# Roadmap

- Patterns in graphs
  - ...
  - Time-evolving graphs
    - T1: Shrinking diameter
    - T2: Densification
    - T3: Connected components
    - T4: Popularity over time
    - T5: Phone-call patterns
  - ...



# More on Time-evolving graphs

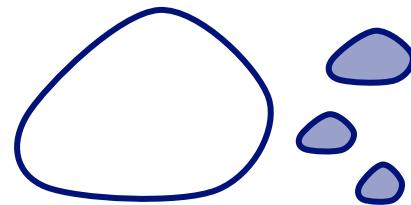
M. McGlohon, L. Akoglu, and C. Faloutsos.  
*Weighted Graphs and Disconnected Components:  
Patterns and a Generator*. KDD 2008

# Observation T.3: NLCC behavior

*Q: How do the NLCCs emerge and join with the GCC?*

(“NLCC” = non-largest conn. components)

- Do they continue to grow in size?
  - or do they shrink?
  - or stabilize?

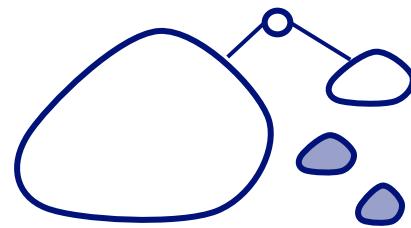


# Observation T.3: NLCC behavior

*Q: How do NLCCs emerge and join with the GCC?*

(“NLCC” = non-largest conn. components)

- Do they continue to grow in size?
  - or do they shrink?
  - or stabilize?



# Observation T.3: NLCC behavior

*Q: How do NLCCs emerge and join with the GCC?*

(“NLCC” = non-largest conn. components)

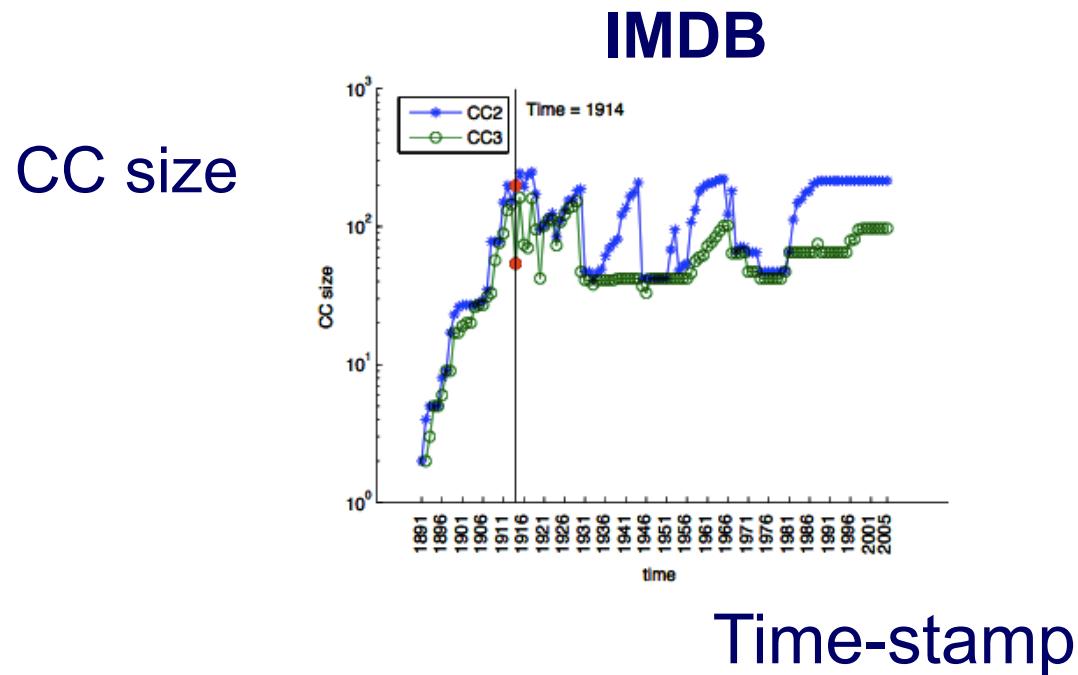
**YES** – Do they continue to grow in size?

**YES** – or do they shrink?

**YES** – or stabilize?

# Observation T.3: NLCC behavior

- After the gelling point, the GCC takes off, but the NLCCs remain ~constant (actually, oscillate).



# Roadmap

- Patterns in graphs
  - ...
  - Time-evolving graphs
    - T1: Shrinking diameter
    - T2: Densification
    - T3: Connected components
    - T4: Popularity over time
    - T5: Phone-call patterns
  - ...

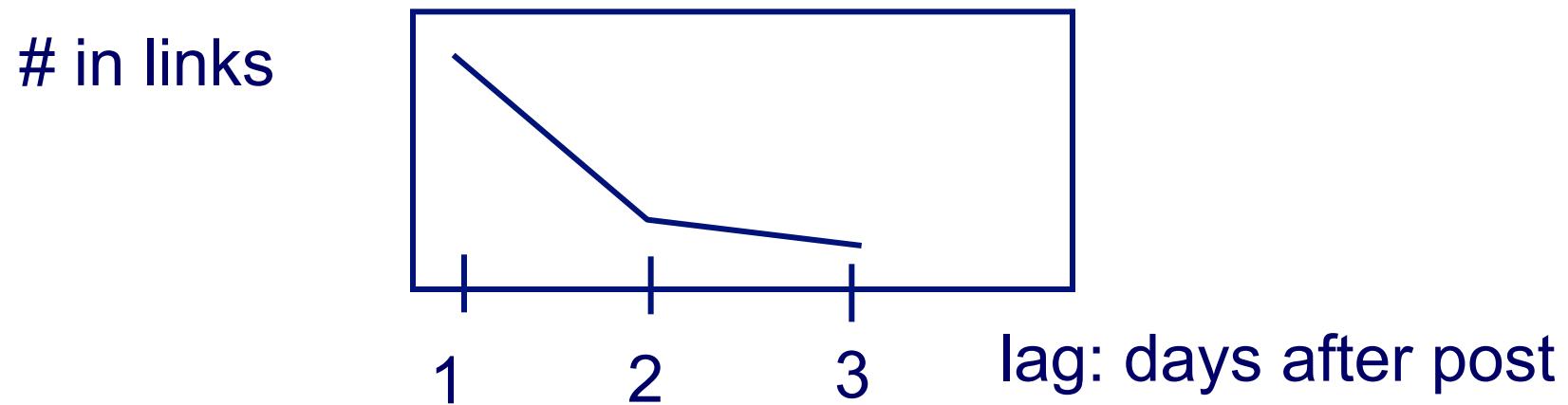


# Timing for Blogs

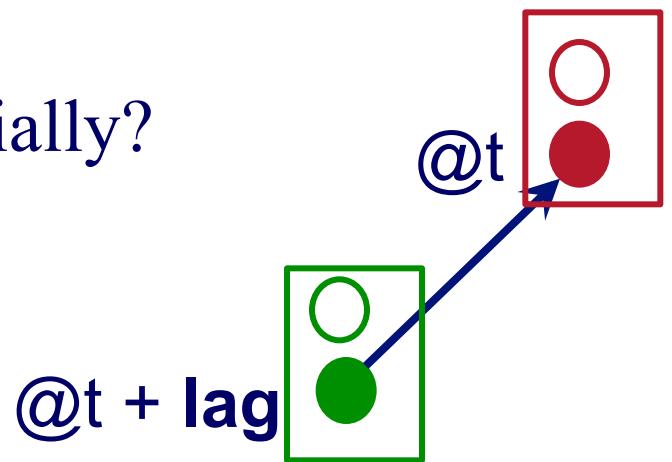
- with Mary McGlohon (CMU→Google)
- Jure Leskovec (CMU→Stanford)
- Natalie Glance (now at Google)
- Matt Hurst (now at MSR)

[SDM' 07]

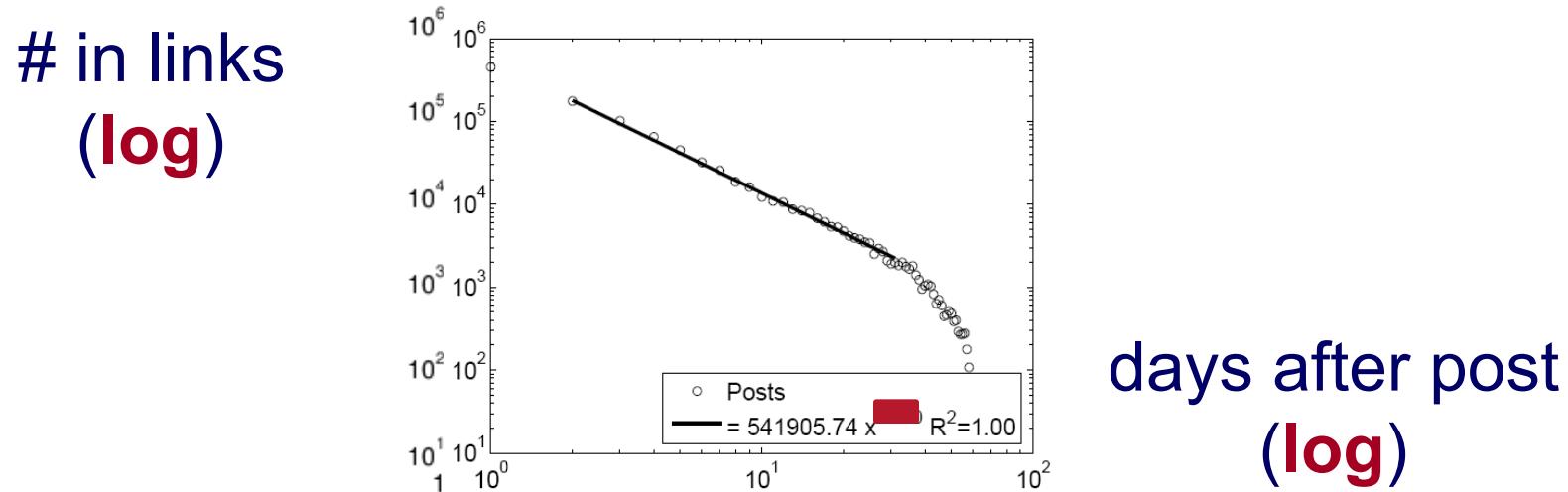
## T.4 : popularity over time



Post popularity drops-off – exponentially?

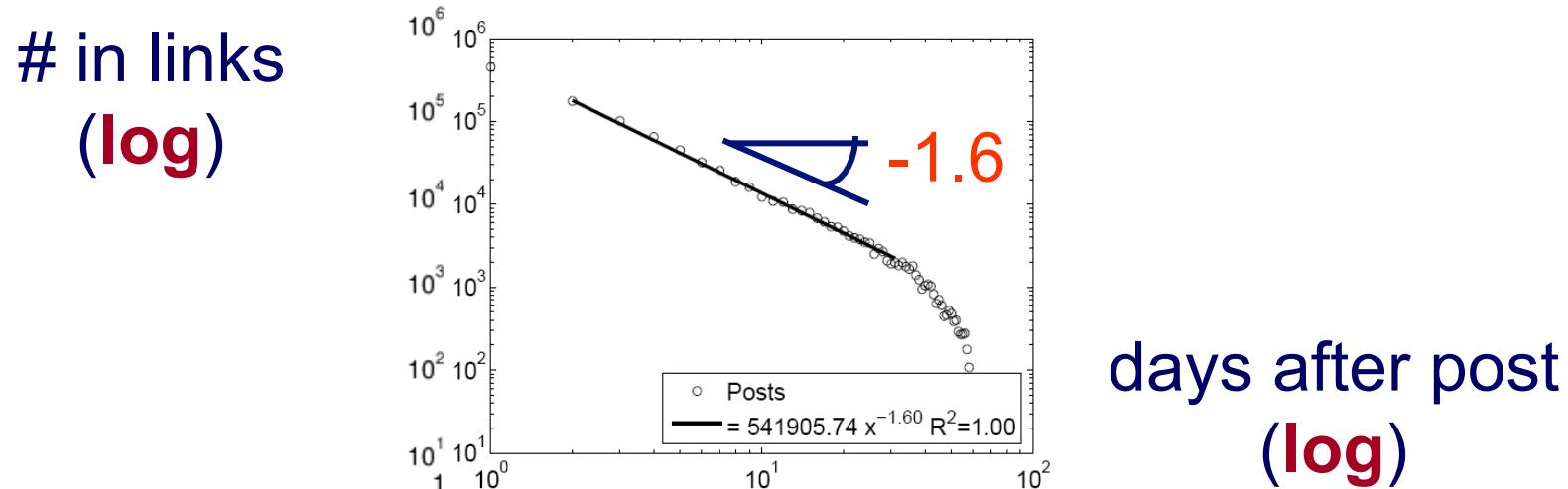


## T.4 : popularity over time



Post popularity drops-off – exponentially?  
POWER LAW!  
Exponent?

## T.4 : popularity over time

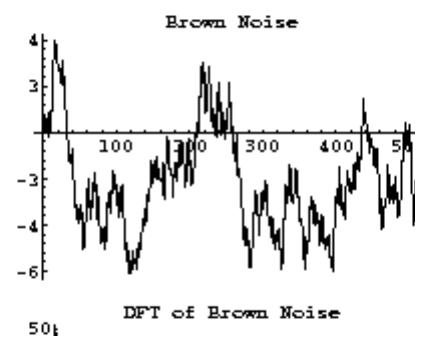


Post popularity drops-off – exponentially?

POWER LAW!

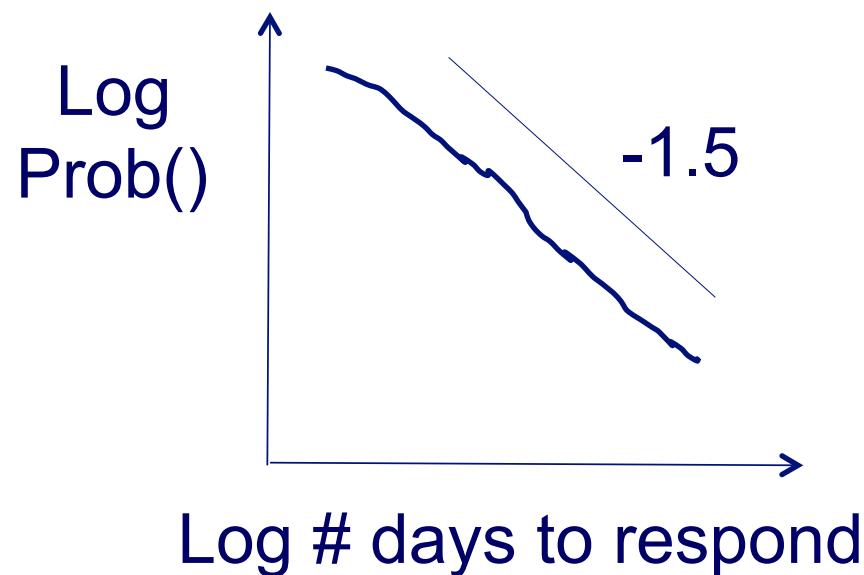
Exponent? -1.6

- close to -1.5: Barabasi's stack model
- and like the zero-crossings of a random walk



# -1.5 slope

J. G. Oliveira & A.-L. Barabási. Human Dynamics: The Correspondence Patterns of Darwin and Einstein. *Nature* **437**, 1251 (2005) . [[PDF](#)]



# Roadmap

- Patterns in graphs
  - ...
  - Time-evolving graphs
    - T1: Shrinking diameter
    - T2: Densification
    - T3: Connected components
    - T4: Popularity over time
    - T5: Phone-call patterns
  - ...



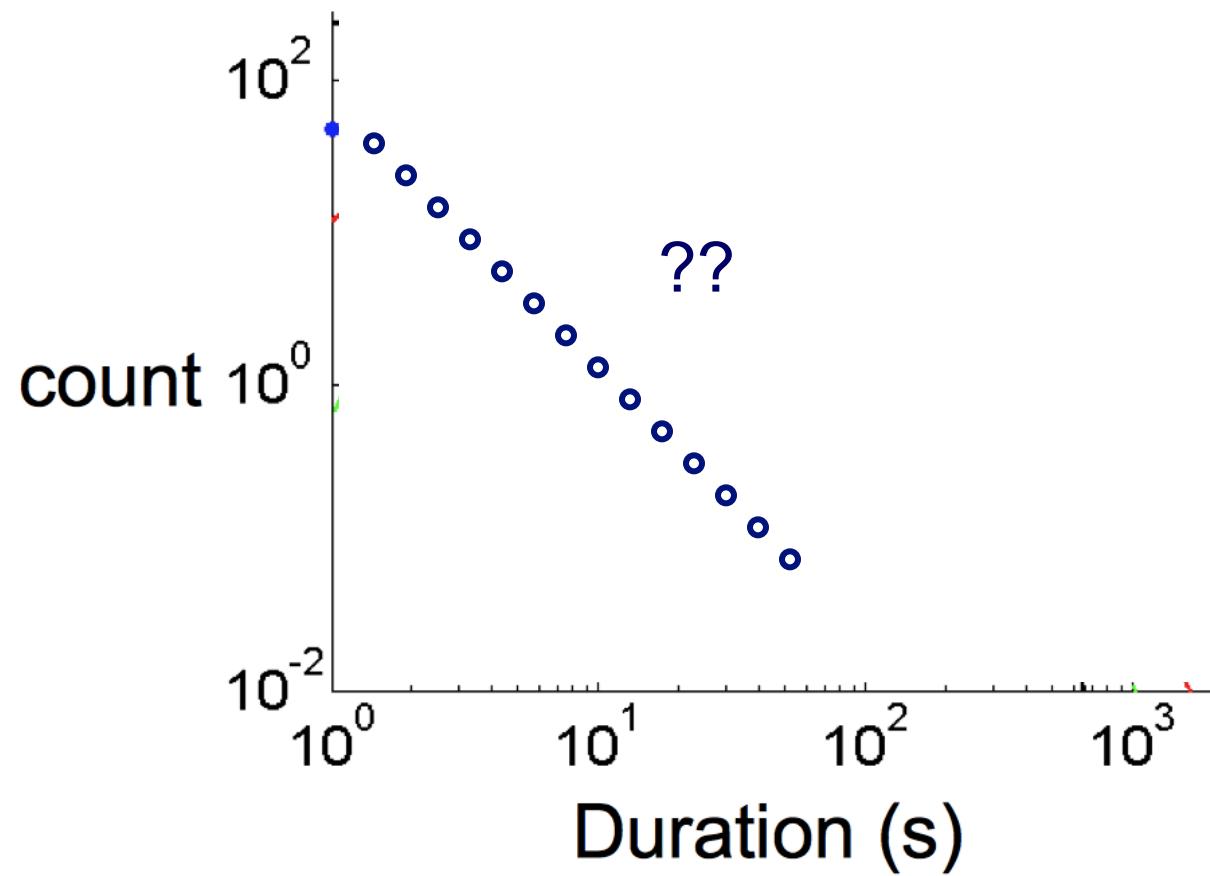
# T.5: duration of phonecalls

*Surprising Patterns for the Call Duration  
Distribution of Mobile Phone Users.*

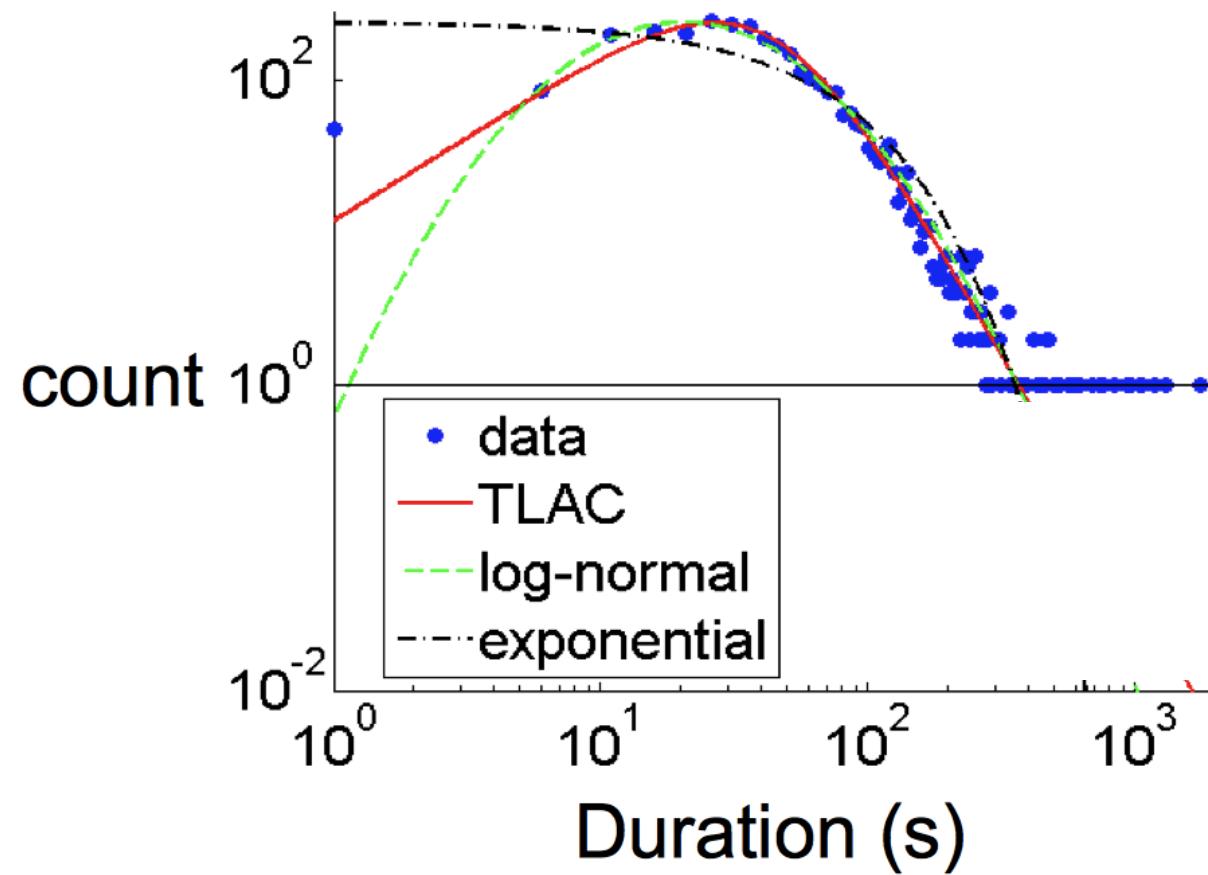


Pedro O. S. Vaz de Melo, Leman Akoglu,  
Christos Faloutsos, Antonio A. F. Loureiro.  
PKDD 2010

# Probably, power law (?)

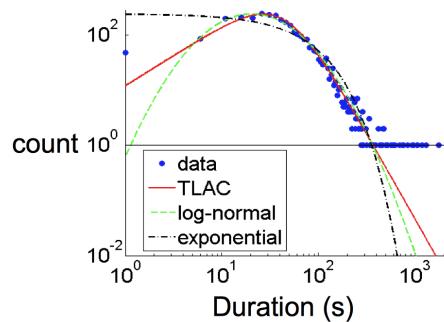


# No Power Law!



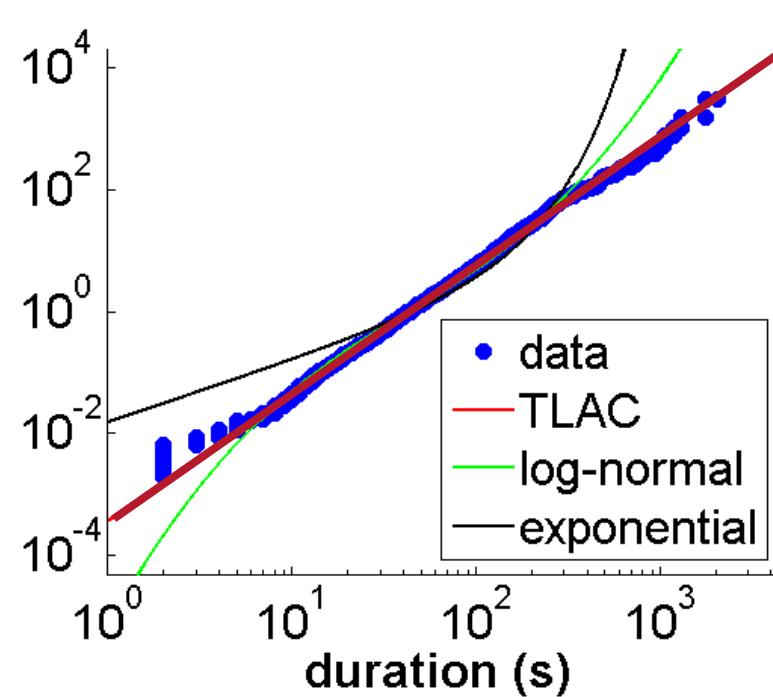
# 'TLaC: Lazy Contractor'

- The longer a task (phone-call) has taken, the even longer it will take



Odds ratio =  
*Casualties(<x)*:  
*Survivors(>=x)*

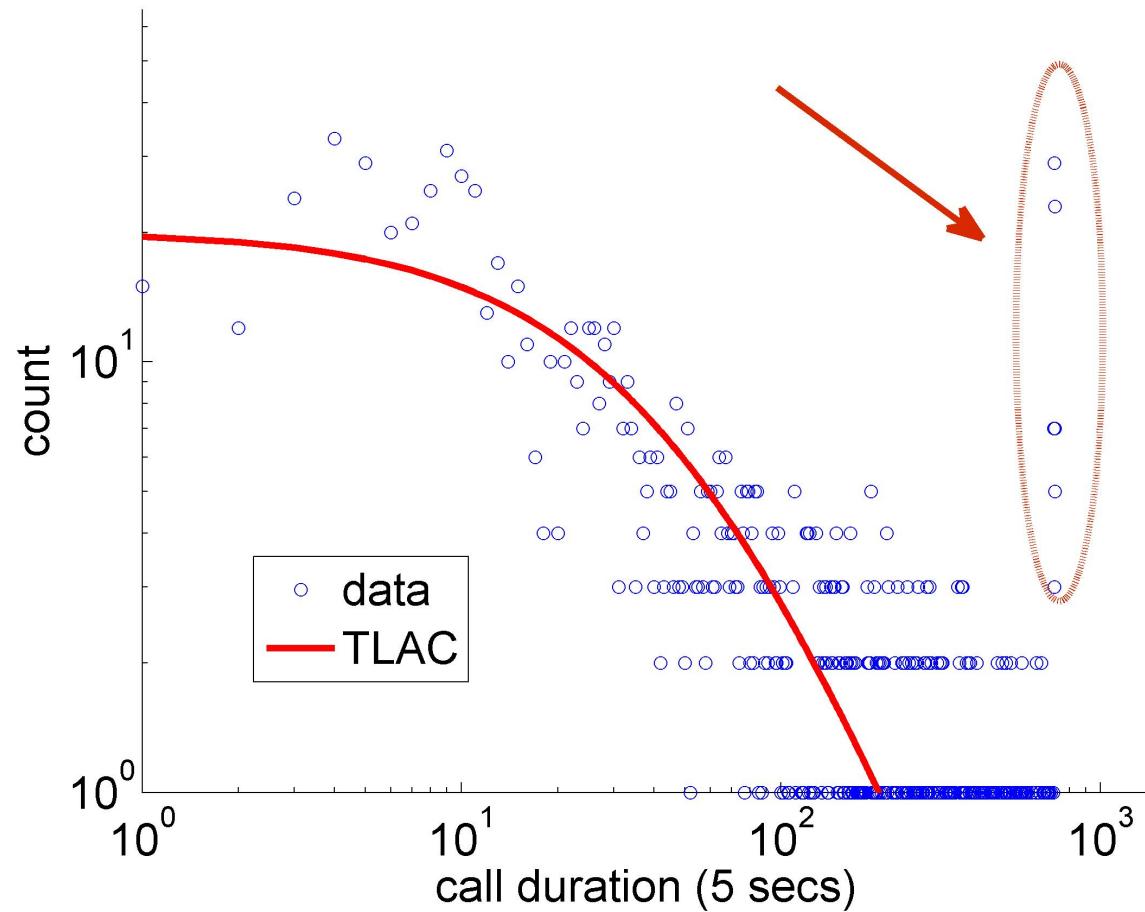
== power law



# Data Description

- Data from a private mobile operator of a large city
  - 4 months of data
  - 3.1 million users
  - more than 1 billion phone records
- Over 96% of ‘talkative’ users obeyed a TLAC distribution (‘talkative’:  $>30$  calls)

# Outliers:



# Real Graph Patterns



	unweighted	weighted
static	<ul style="list-style-type: none"> <li>✓ <b>P01.</b> Power-law degree distribution [Faloutsos et. al. '99, Kleinberg et. al. '99, Chakrabarti et. al. '04, Newman '04]</li> <li>✓ <b>P02.</b> Triangle Power Law [Tsourakakis '08]</li> <li>✓ <b>P03.</b> Eigenvalue Power Law [Siganos et. al. '03]</li> <li>✓ <b>P04.</b> Community structure [Flake et. al. '02, Girvan and Newman '02]</li> <li>✓ <b>P05.</b> Clique Power Laws [Du et. al. '09]</li> </ul>	<ul style="list-style-type: none"> <li><b>P12.</b> Snapshot Power Law [McGlohon et. al. '08]</li> </ul>
dynamic	<ul style="list-style-type: none"> <li>✓ <b>P06.</b> Densification Power Law [Leskovec et. al. '05]</li> <li>✓ <b>P07.</b> Small and shrinking diameter [Albert and Barabási '99, Leskovec et. al. '05, McGlohon et. al. '08]</li> <li>✓ <b>P08.</b> Gelling point [McGlohon et. al. '08]</li> <li>✓ <b>P09.</b> Constant size 2<sup>nd</sup> and 3<sup>rd</sup> connected components [McGlohon et. al. '08]</li> <li><b>P10.</b> Principal Eigenvalue Power Law [Akoglu et. al. '08]</li> <li><b>P11.</b> Bursty/self-similar edge/weight additions [Gomez and Santonja '98, Gribble et. al. '98, Crovella and Bestavros '99, McGlohon et. al. '08]</li> </ul>	<ul style="list-style-type: none"> <li>✓ <b>P13.</b> Weight Power Law [McGlohon et. al. '08]</li> <li>✓ <b>P14.</b> Skewed call duration distributions [Vaz de Melo et. al. '10]</li> </ul>

[RTG: A Recursive Realistic Graph Generator using Random Typing](#)  
Leman Akoglu and Christos Faloutsos. *ECML PKDD'09*.

# Roadmap

- Patterns in graphs
  - Overview
  - Static graphs
  - Weighted graphs
  - Time-evolving graphs
- ➡ • Anomaly Detection
- Application: ebay fraud
- Conclusions



# OddBall: Spotting Anomalies in Weighted Graphs



Leman Akoglu, Mary McGlohon,  
Christos Faloutsos

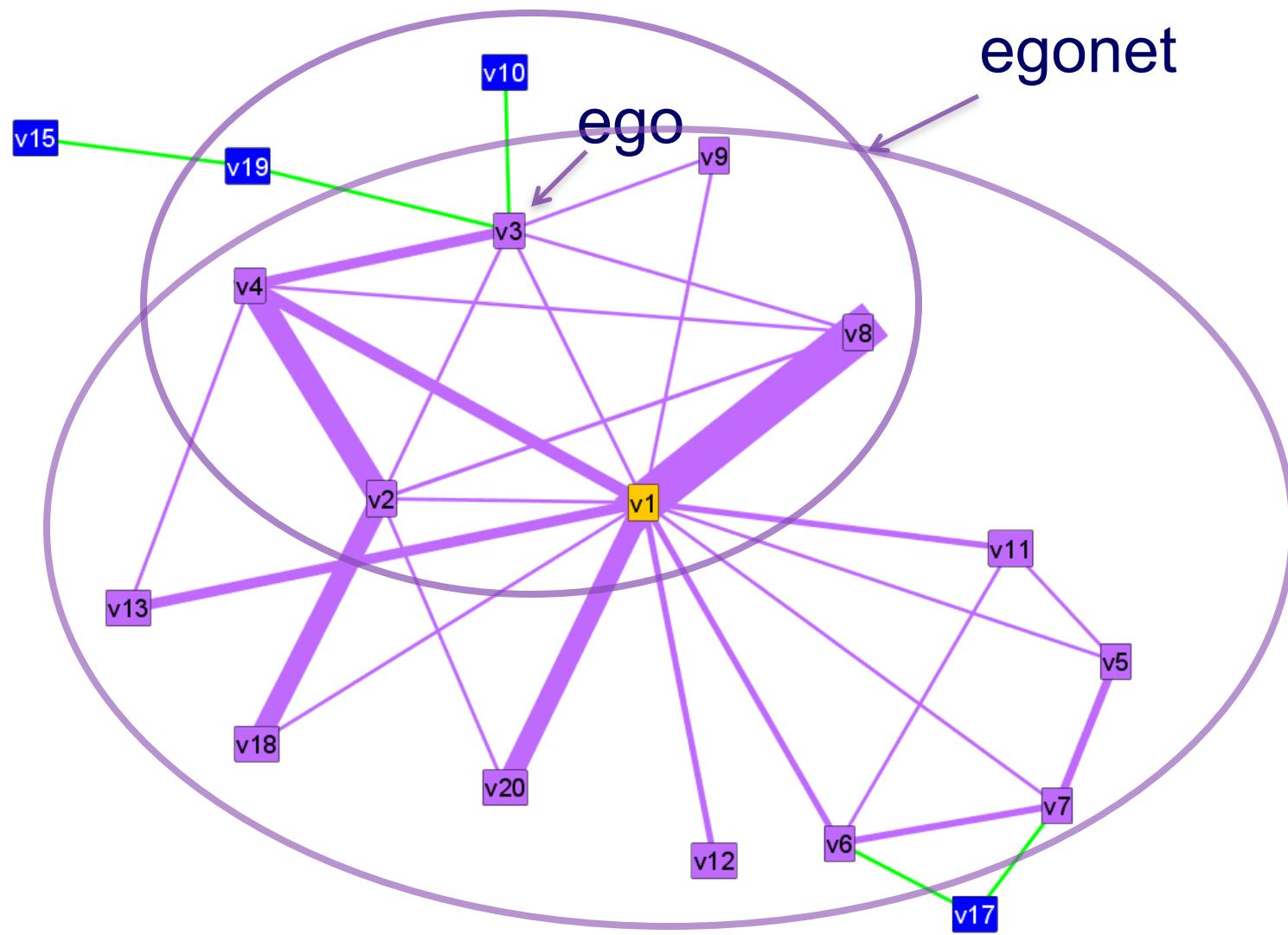
PAKDD 2010, Hyderabad, India

# Main idea

For each node,

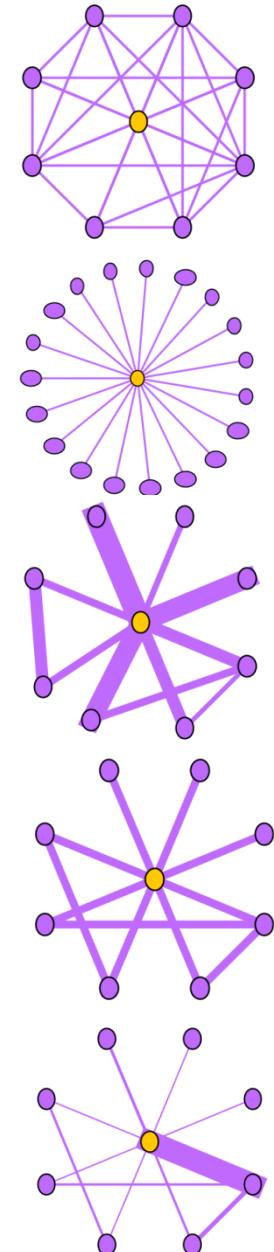
- Extract ‘ego-net’ (=1-step-away neighbors)
- Extract features (#edges, total weight, *etc*)
- Compare with the rest of the population

# What is an egonet?

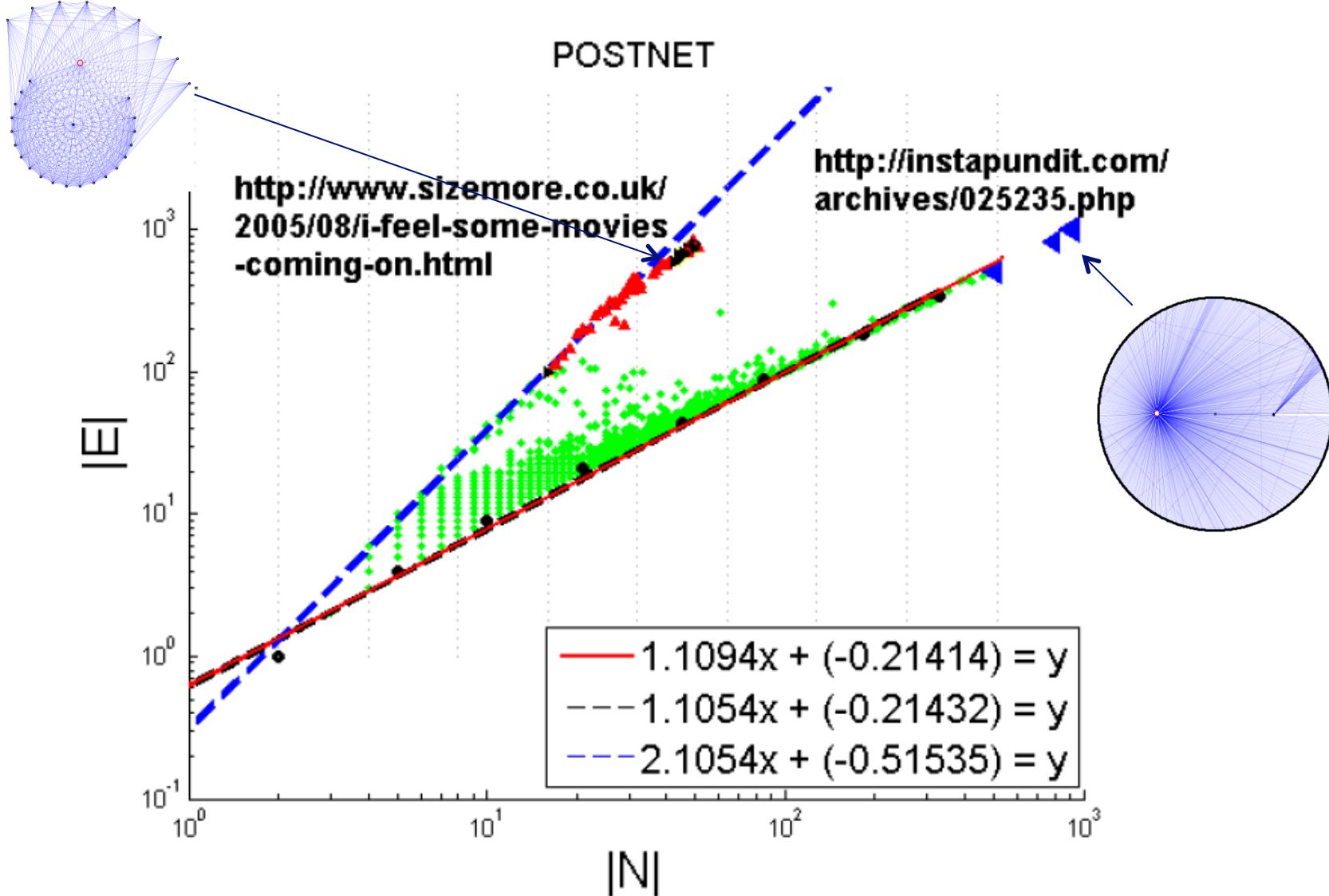


# Selected Features

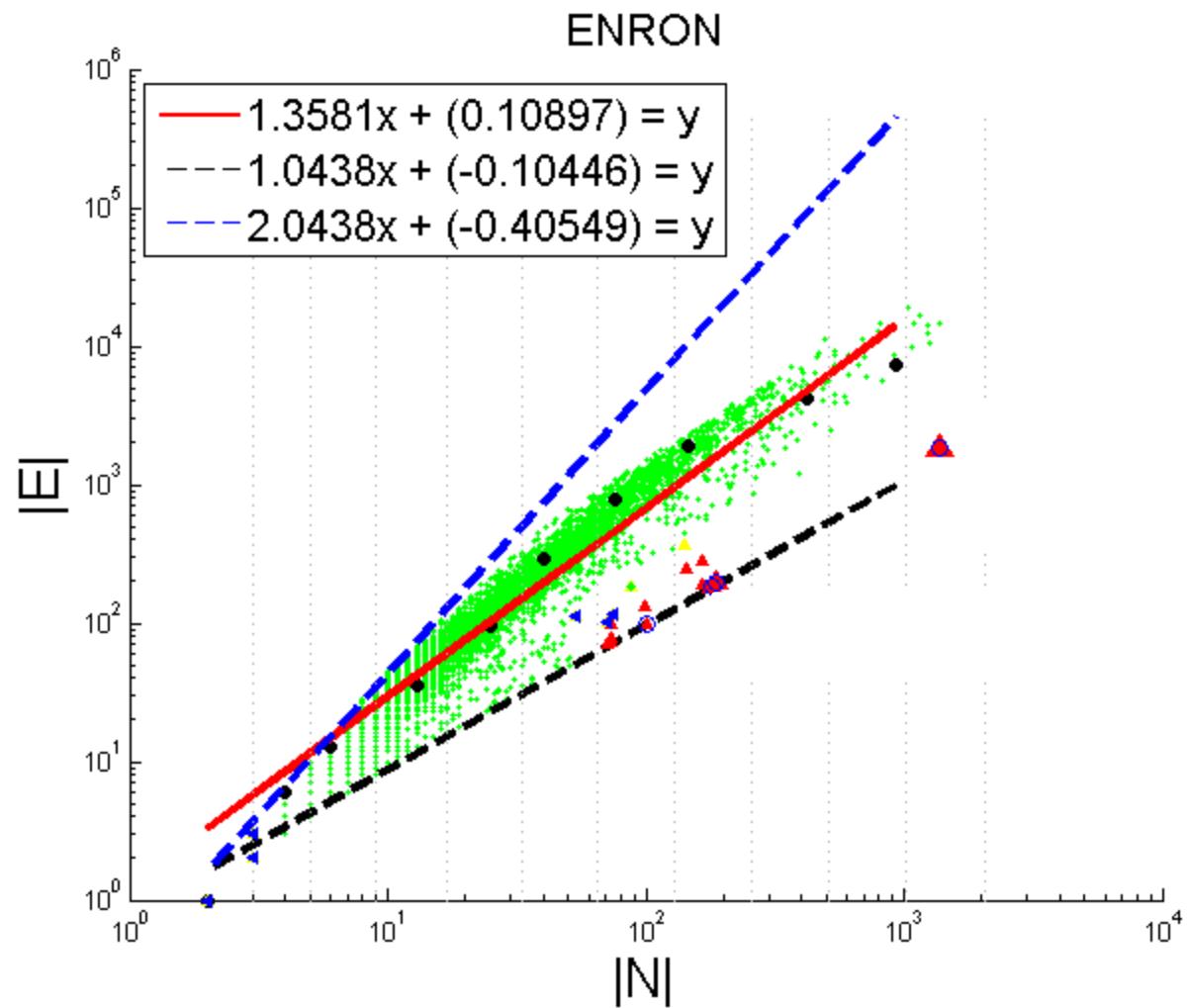
- $N_i$ : number of neighbors (degree) of ego  $i$
- $E_i$ : number of edges in egonet  $i$
- $W_i$ : total weight of egonet  $i$
- $\lambda_{w,i}$ : principal eigenvalue of the **weighted** adjacency matrix of egonet  $I$



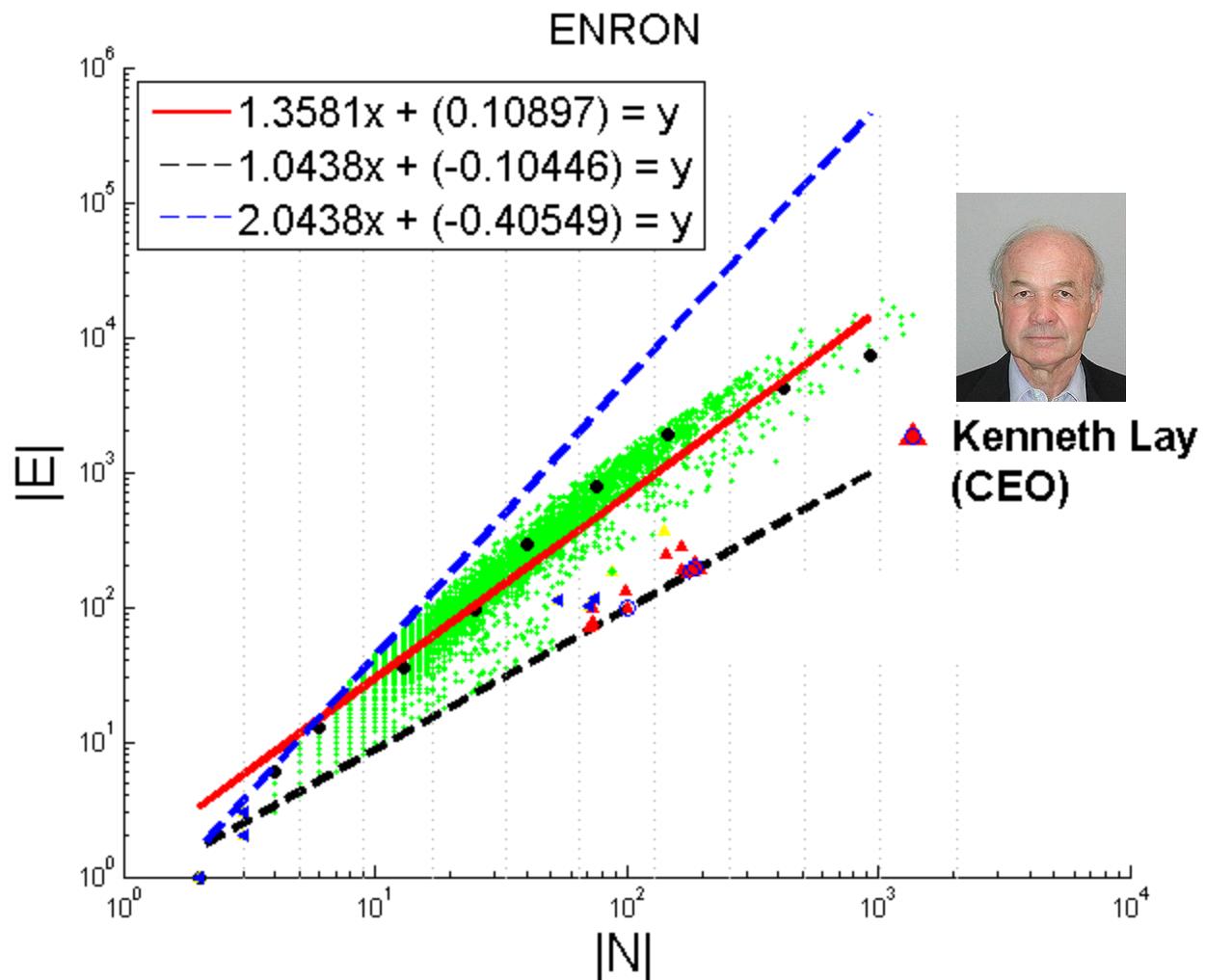
# Near-Clique/Star



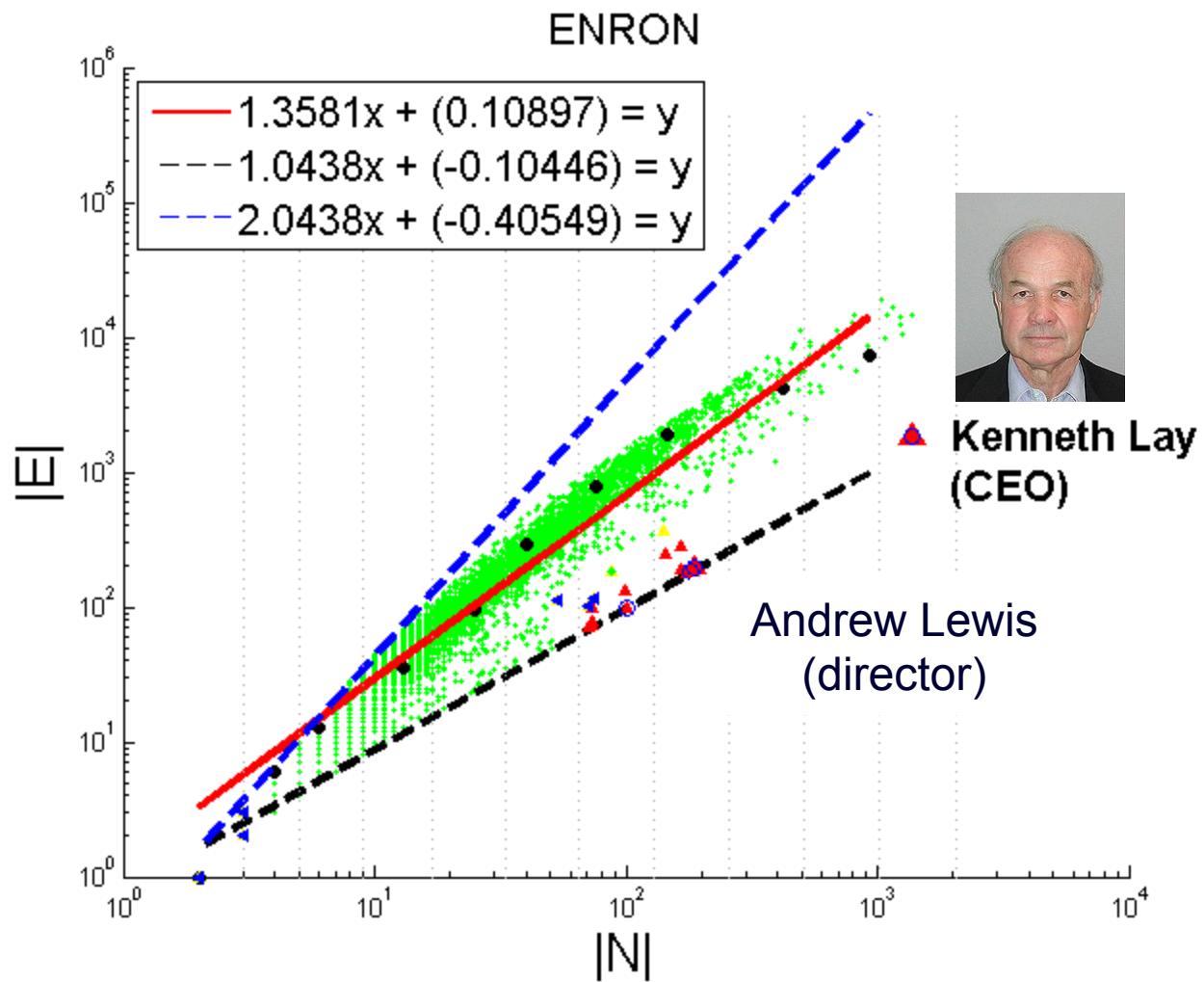
# Near-Clique/Star



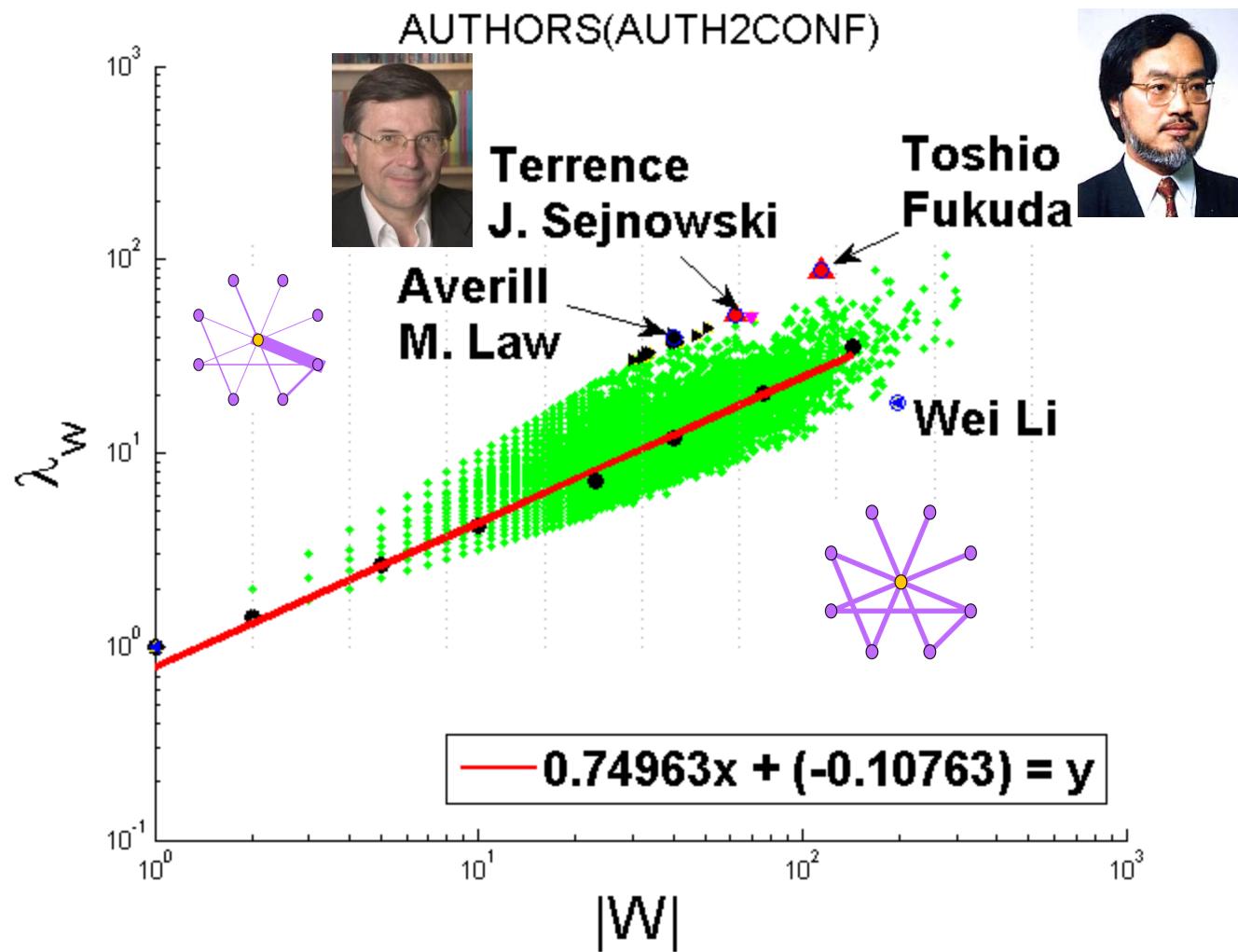
# Near-Clique/Star



# Near-Clique/Star



# Dominant Heavy Link



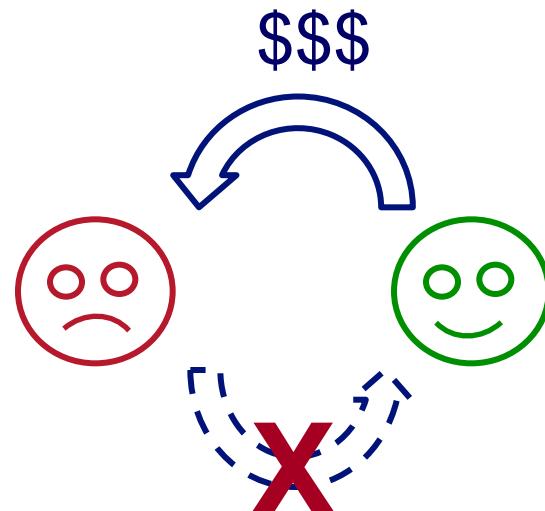
# Roadmap

- Patterns in graphs
  - overview
  - Static graphs
  - Weighted graphs
  - Time-evolving graphs
- Anomaly Detection
- • Application: ebay fraud
- Conclusions



# NetProbe: The Problem

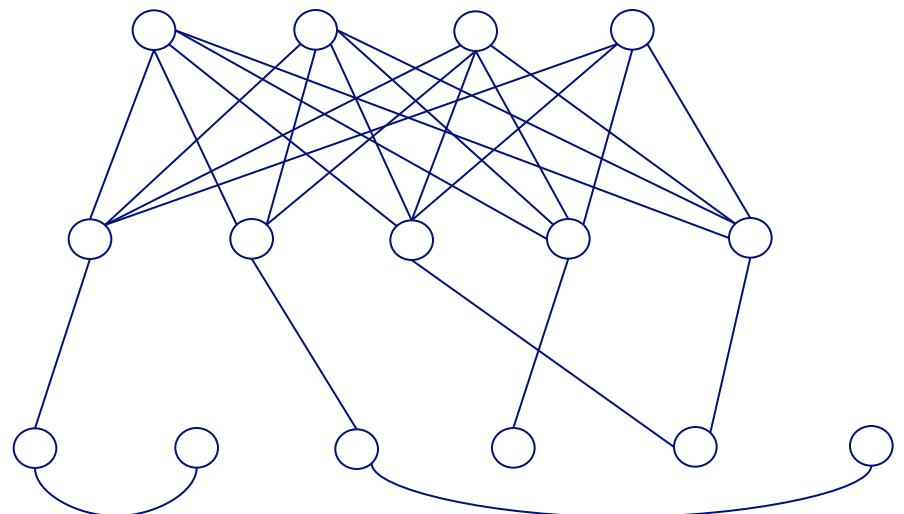
Find **bad sellers (fraudsters)** on eBay  
who don't deliver their (expensive)  
items



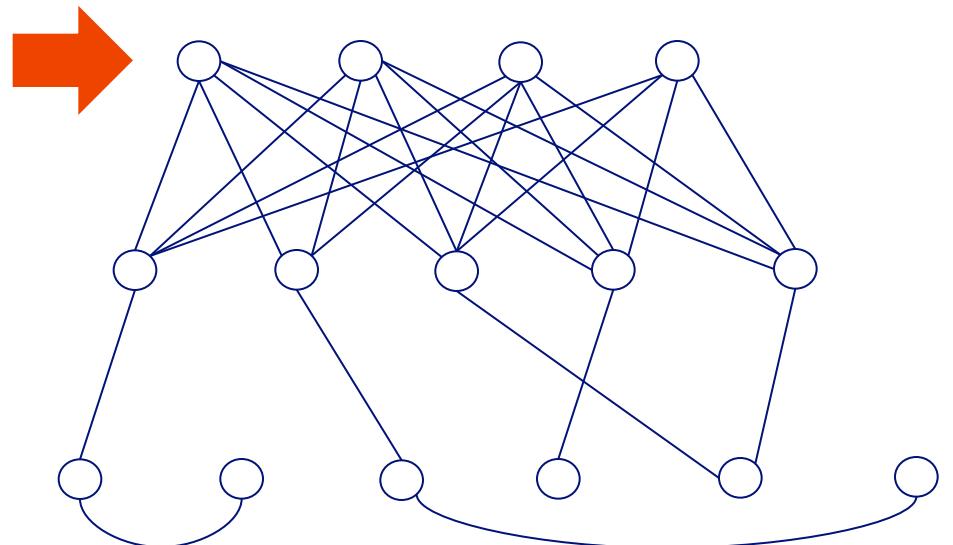
# E-bay Fraud detection



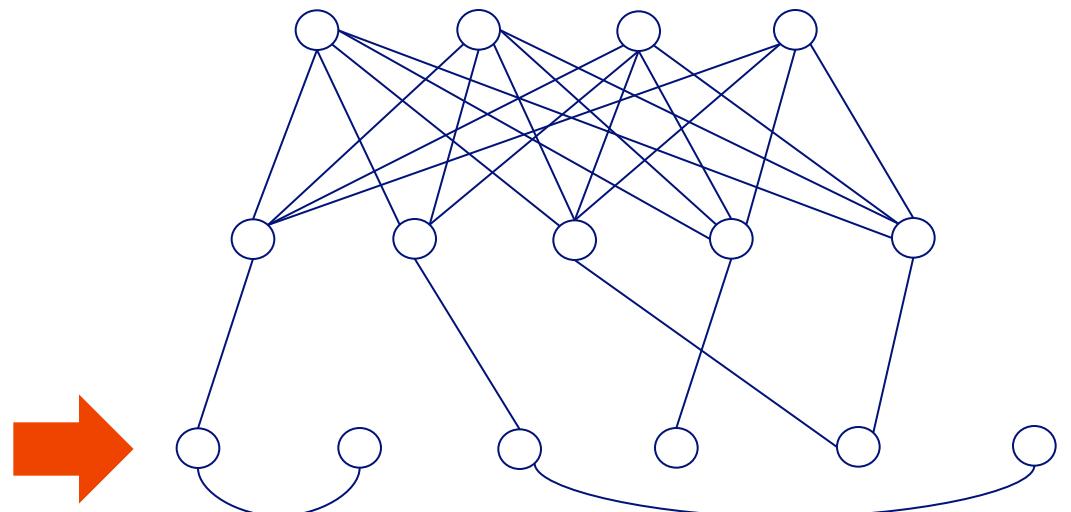
w/ Polo Chau &  
Shashank Pandit, CMU  
[www' 07]



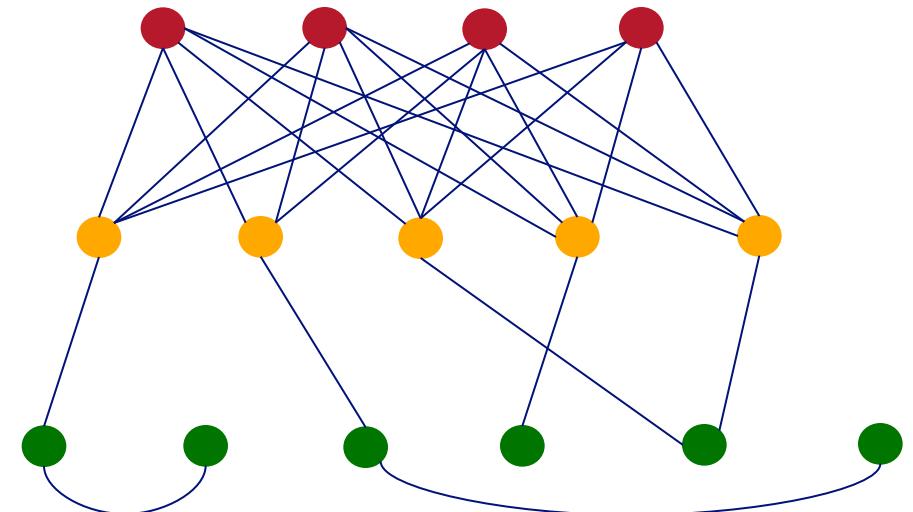
# E-bay Fraud detection



# E-bay Fraud detection

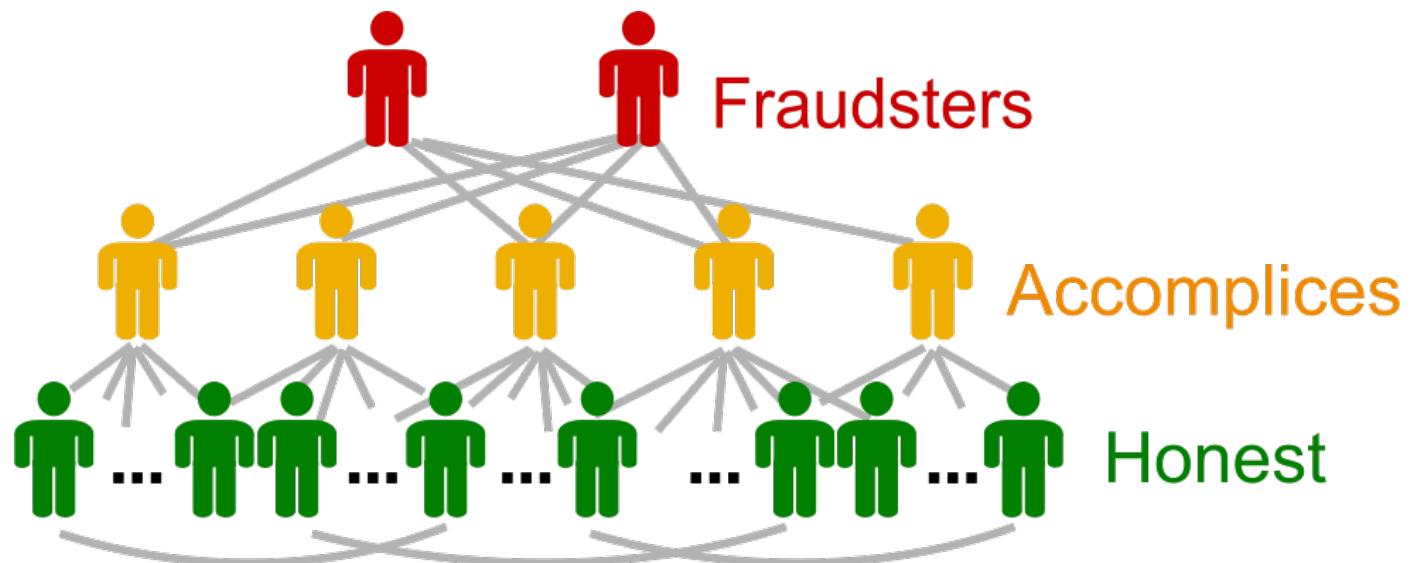


# E-bay Fraud detection - NetProbe



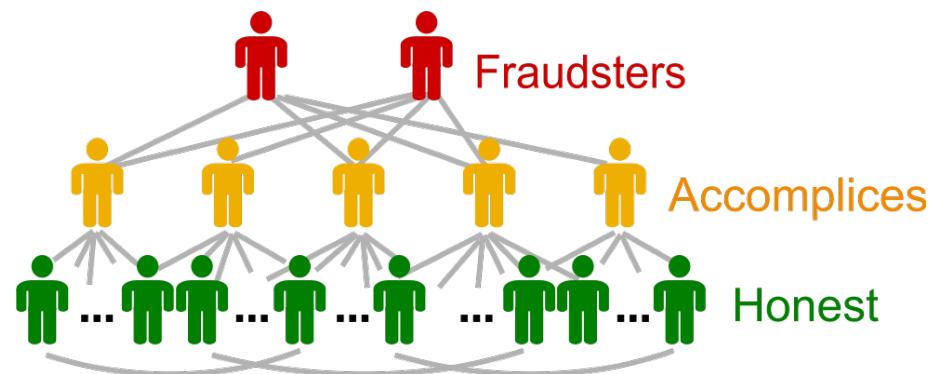
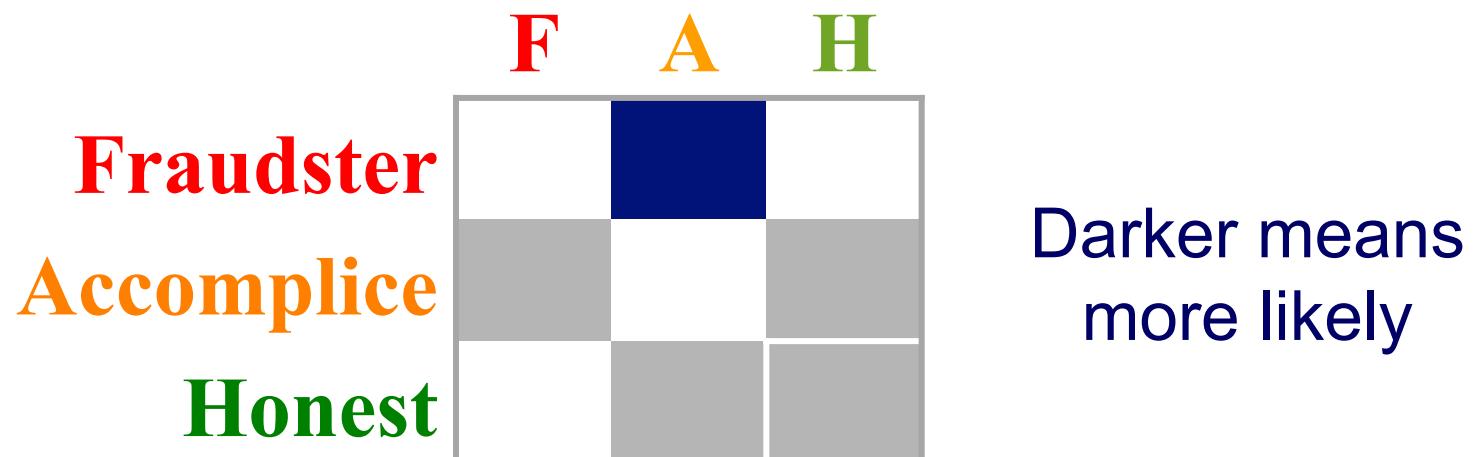
# NetProbe: Key Ideas

- Fraudsters **fabricate their reputation** by “trading” with their accomplices
- Transactions form **near bipartite cores**
- How to detect them?

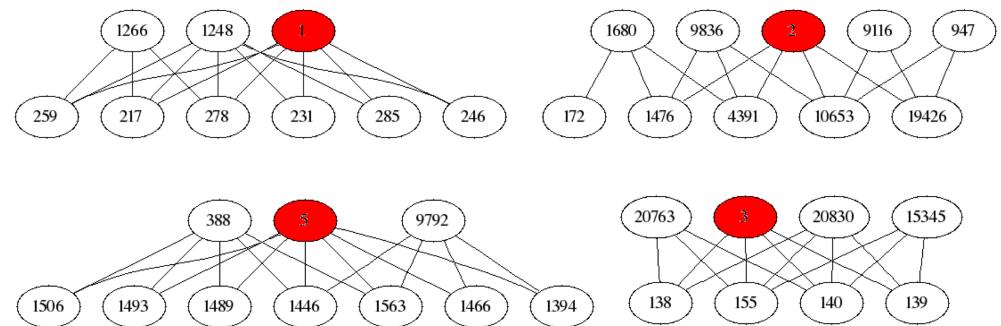
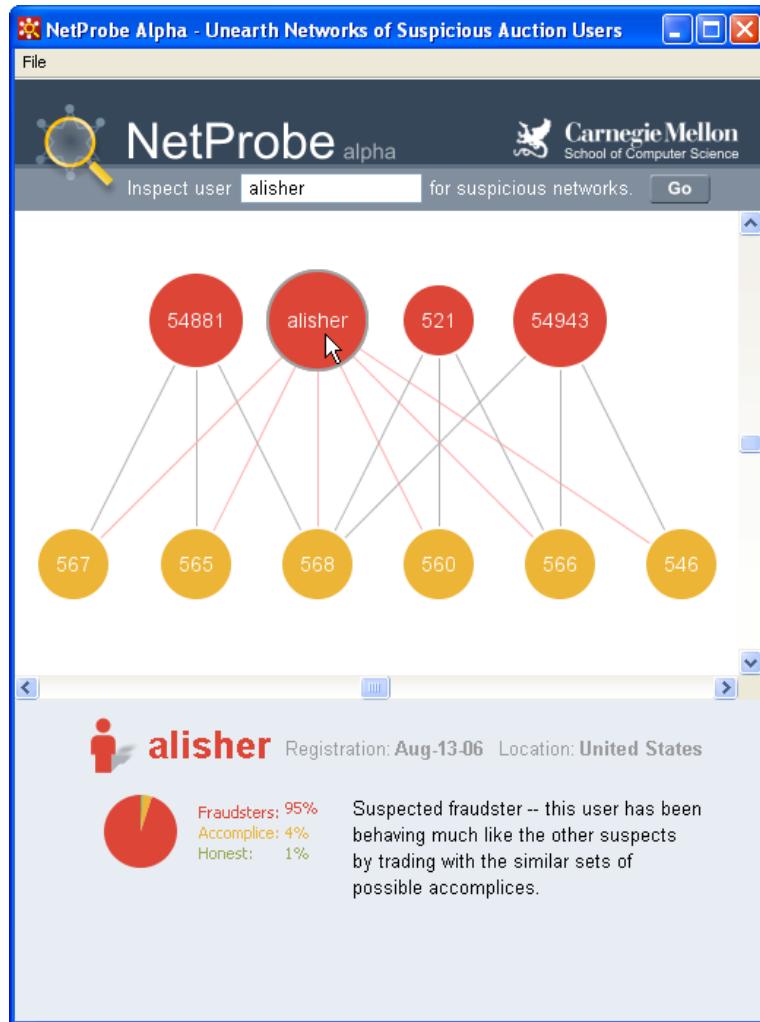


# NetProbe: Key Ideas

Use ‘Belief Propagation’ and ~heterophily



# NetProbe: Main Results

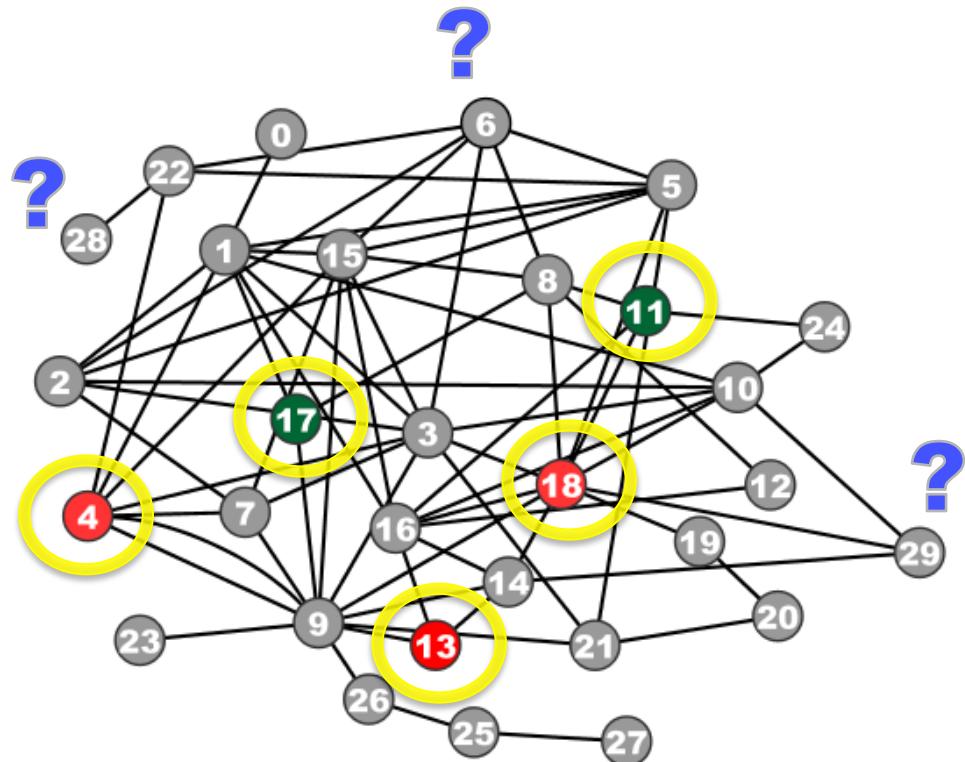


# Roadmap

- Patterns in graphs
- Anomaly Detection
- Application: ebay fraud
- Conclusions



# Guilt-by-Association Techniques



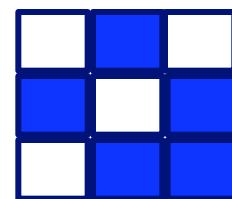
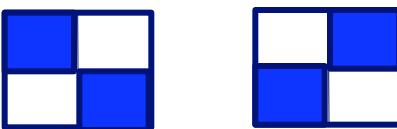
**Given:**

- graph and
- few labeled nodes

**Find:** class (red/green)  
for rest nodes

**Assuming:** network  
effects (homophily/  
heterophily, etc)

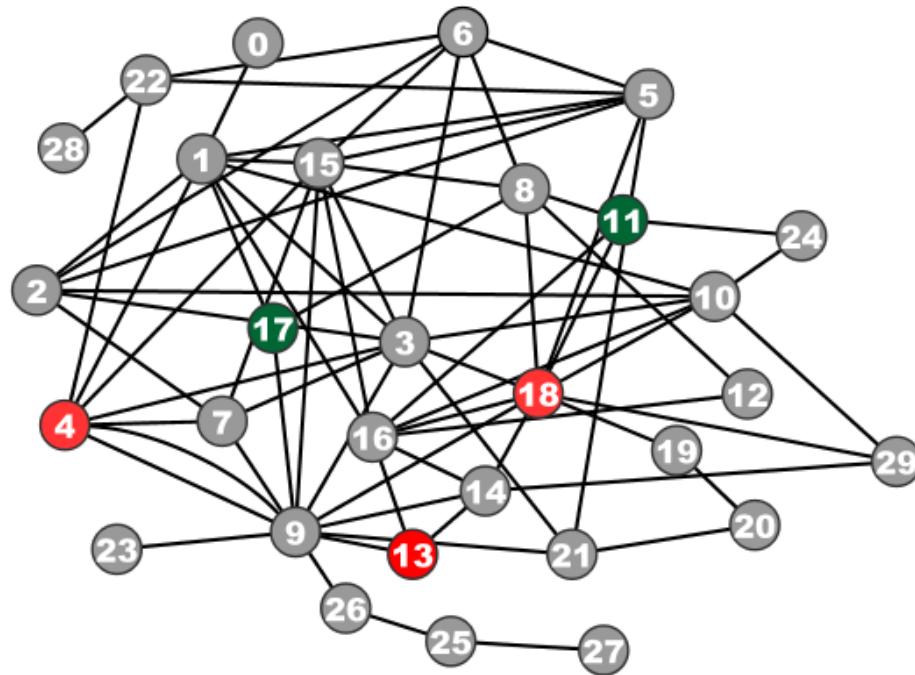
red  
green



F  
A  
H

# Correspondence of Methods

Random Walk with Restarts (RWR)   Google  
Semi-supervised Learning (SSL)  
Belief Propagation (BP)   Bayesian

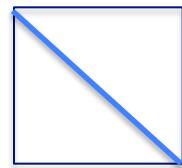
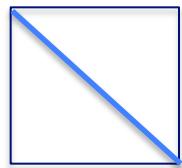


# Correspondence of Methods



Random Walk with Restarts (RWR)  $\approx$   
 Semi-supervised Learning (SSL)  $\approx$   
 Belief Propagation (BP)

Method	Matrix	unknown	known
RWR	$[I - c A D^{-1}]$	$x$	$= (1-c)y$
SSL	$[I + \alpha(D - A)]$	$x$	$= y$
FABP	$[I + \alpha D - c' A]$	$b_h$	$= \phi_h$



$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

$$\boxed{?}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

Unifying Guilt-by-Association Approaches: Theorems and Fast Algorithms. Danai Koutra, et al PKDD'11

# Roadmap

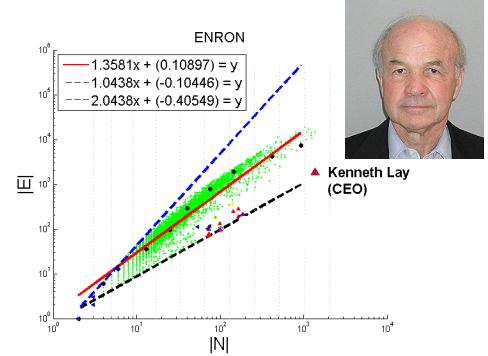
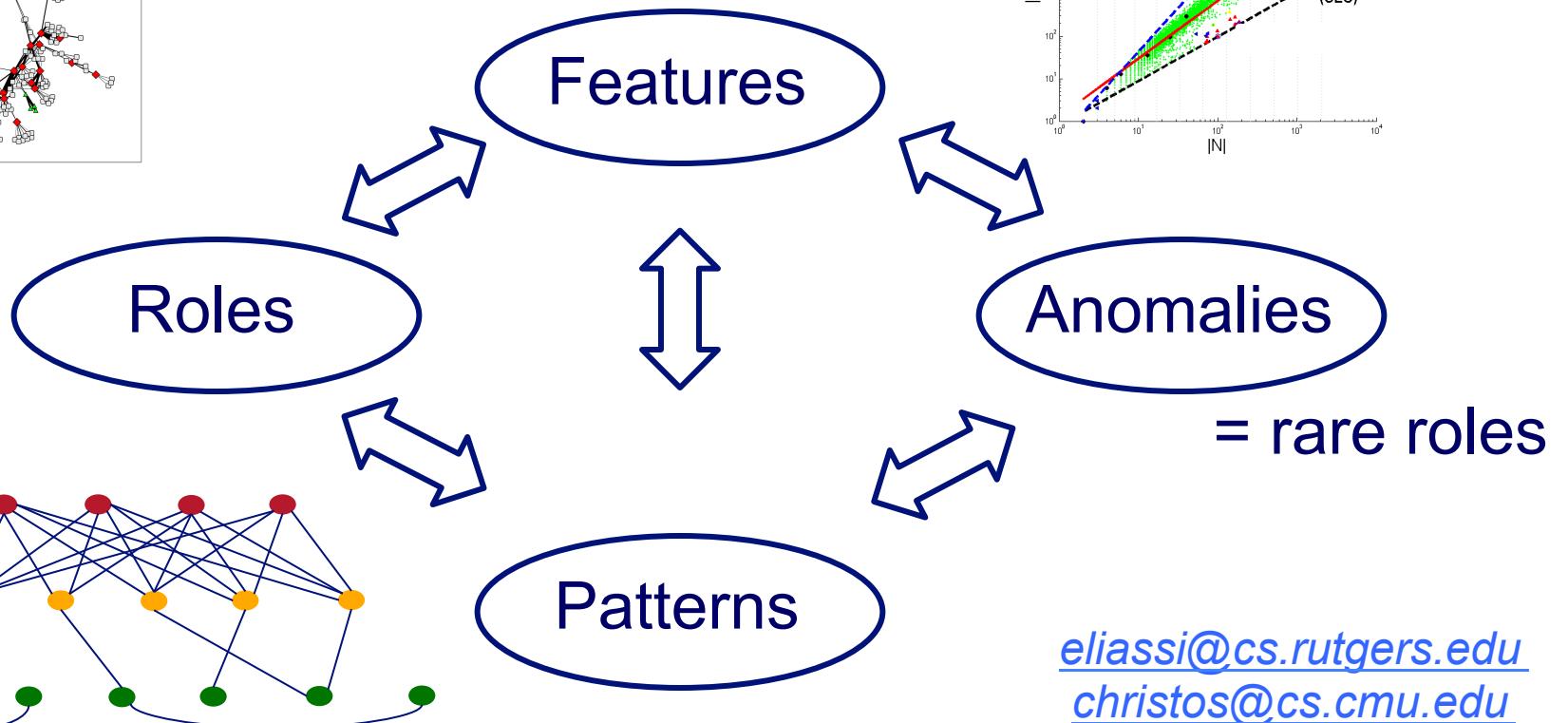
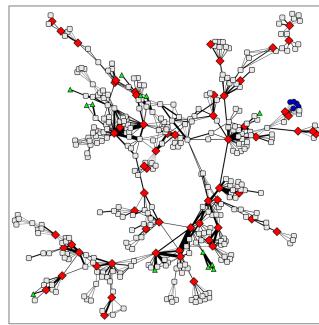
- Patterns in graphs
- Anomaly Detection
- Application: ebay fraud
- • Conclusions



# Overall Conclusions

- Roles
  - Past work in social networks ('regular', 'structural', *etc*)
  - Scalable algorithms' to find such roles
- Anomalies & patterns
  - Static (power-laws, 'six degrees')
  - Weighted (super-linearity)
  - Time-evolving (densification, -1.5 exponent)

# Thank You!



[eliassi@cs.rutgers.edu](mailto:eliassi@cs.rutgers.edu)  
[christos@cs.cmu.edu](mailto:christos@cs.cmu.edu)

# Project info

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



Chau,  
Polo



Koutra,  
Danai



Prakash,  
Aditya



Akoglu,  
Leman

Kang, U

McGlohon,  
Mary

Tong,  
Hanghang

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

# References

- Leman Akoglu, Christos Faloutsos: *RTG: A Recursive Realistic Graph Generator Using Random Typing*. ECML/PKDD (1) 2009: 13-28
- Deepayan Chakrabarti, Christos Faloutsos: *Graph mining: Laws, generators, and algorithms*. ACM Comput. Surv. 38(1): (2006)

# References

- Deepayan Chakrabarti, Yang Wang, Chenxi Wang, Jure Leskovec, Christos Faloutsos: *Epidemic thresholds in real networks*. ACM Trans. Inf. Syst. Secur. 10(4): (2008)
- Deepayan Chakrabarti, Jure Leskovec, Christos Faloutsos, Samuel Madden, Carlos Guestrin, Michalis Faloutsos: *Information Survival Threshold in Sensor and P2P Networks*. INFOCOM 2007: 1316-1324

# References

- Christos Faloutsos, Tamara G. Kolda, Jimeng Sun:  
*Mining large graphs and streams using matrix and tensor tools*. Tutorial, SIGMOD Conference 2007: 1174

# References

- T. G. Kolda and J. Sun. *Scalable Tensor Decompositions for Multi-aspect Data Mining*. In: ICDM 2008, pp. 363-372, December 2008.

# References

- Jure Leskovec, Jon Kleinberg and Christos Faloutsos *Graphs over Time: Densification Laws, Shrinking Diameters and Possible Explanations*, KDD 2005 (Best Research paper award).
- Jure Leskovec, Deepayan Chakrabarti, Jon M. Kleinberg, Christos Faloutsos: *Realistic, Mathematically Tractable Graph Generation and Evolution, Using Kronecker Multiplication*. PKDD 2005: 133-145

# References

- Jimeng Sun, Yinglian Xie, Hui Zhang, Christos Faloutsos. *Less is More: Compact Matrix Decomposition for Large Sparse Graphs*, SDM, Minneapolis, Minnesota, Apr 2007.
- Jimeng Sun, Spiros Papadimitriou, Philip S. Yu, and Christos Faloutsos, *GraphScope: Parameter-free Mining of Large Time-evolving Graphs* ACM SIGKDD Conference, San Jose, CA, August 2007

# References

- Jimeng Sun, Dacheng Tao, Christos Faloutsos: *Beyond streams and graphs: dynamic tensor analysis*. KDD 2006: 374-383

# References

- Hanghang Tong, Christos Faloutsos, and Jia-Yu Pan, *Fast Random Walk with Restart and Its Applications*, ICDM 2006, Hong Kong.
- Hanghang Tong, Christos Faloutsos, *Center-Piece Subgraphs: Problem Definition and Fast Solutions*, KDD 2006, Philadelphia, PA

# References

- Hanghang Tong, Christos Faloutsos, Brian Gallagher, Tina Eliassi-Rad: Fast best-effort pattern matching in large attributed graphs.  
KDD 2007: 737-746