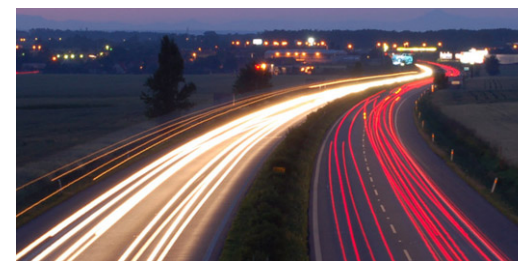# Large Graph Mining - Patterns, Explanations and Cascade Analysis

## *Christos Faloutsos*

## CMU

# Roadmap



➡ • A case for cross-disciplinarity

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

• Part#1: Patterns in graphs

• Part#2: Cascade analysis

• Conclusions
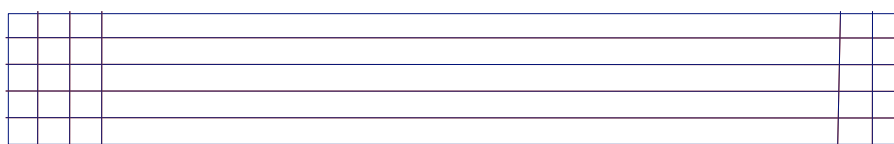
# Problem dfn

Measurement 1 ← 3yrs, every 5' →

~1000

Measurement M

Voltage 1 →

?

Voltage N →

time

Average Hourly Load, PJM Mid-Atlantic Region

# **Problem dfn**

Measurement 1

Measurement M

Voltage 1

Voltage N

?

time

Direct solution:
Slow (Kirchoff's eq.)

# **Problem dfn**

Measurement 1

Measurement M

Voltage 1

Voltage N

?

time

Look for near-neighbors
And use *their* voltages

# Problem dfn

Measurement 1

Measurement M

Voltage 1

Voltage N

T

time

But sequential scan
Is slow, too (MxT)
Can we do better?

Look for near-neighbors
And use *their* voltages

?

# **Problem dfn**

Measurement 1

Measurement M

Voltage 1

Voltage N

T

time

But sequential scan
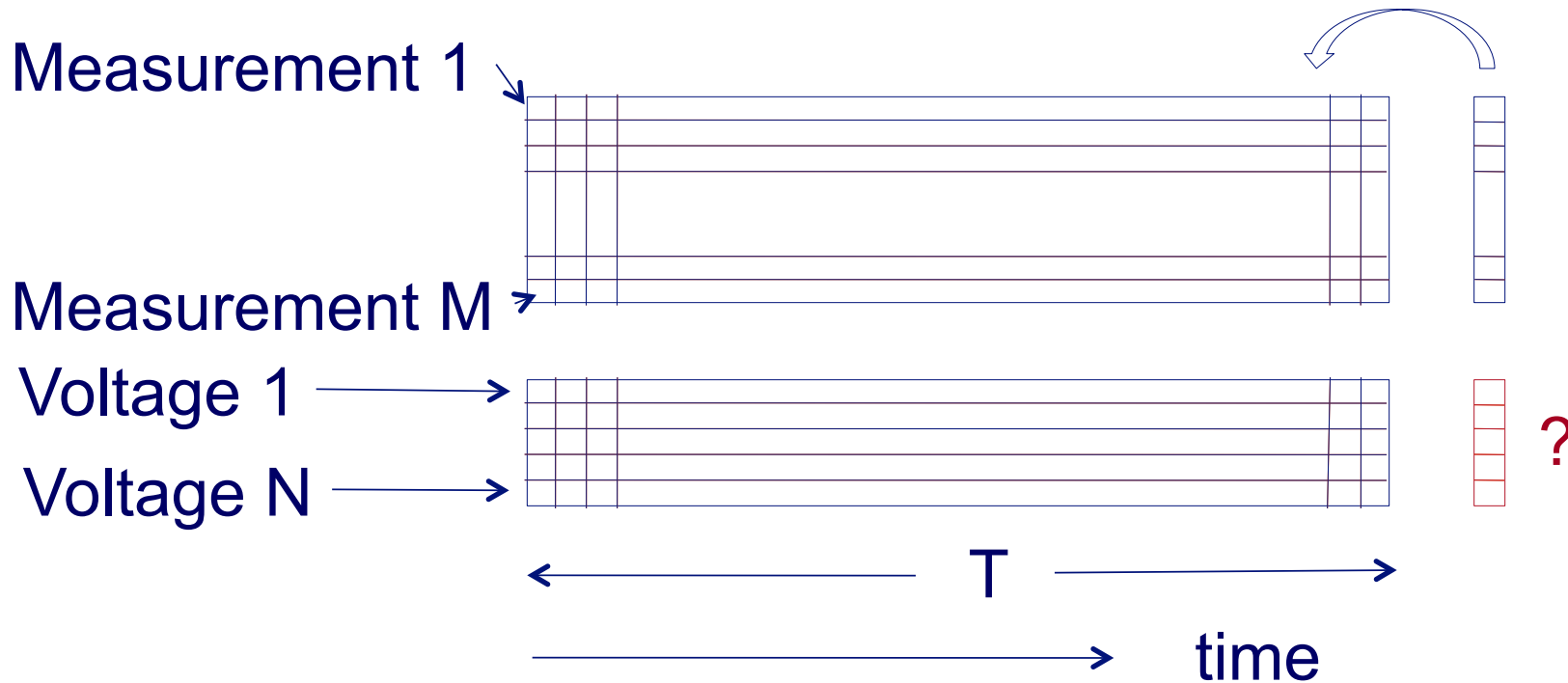Is slow, too (MxT)
Can we do better?

A: yes!
We can reduce both
•T, and
•M

?

# Simulation Results

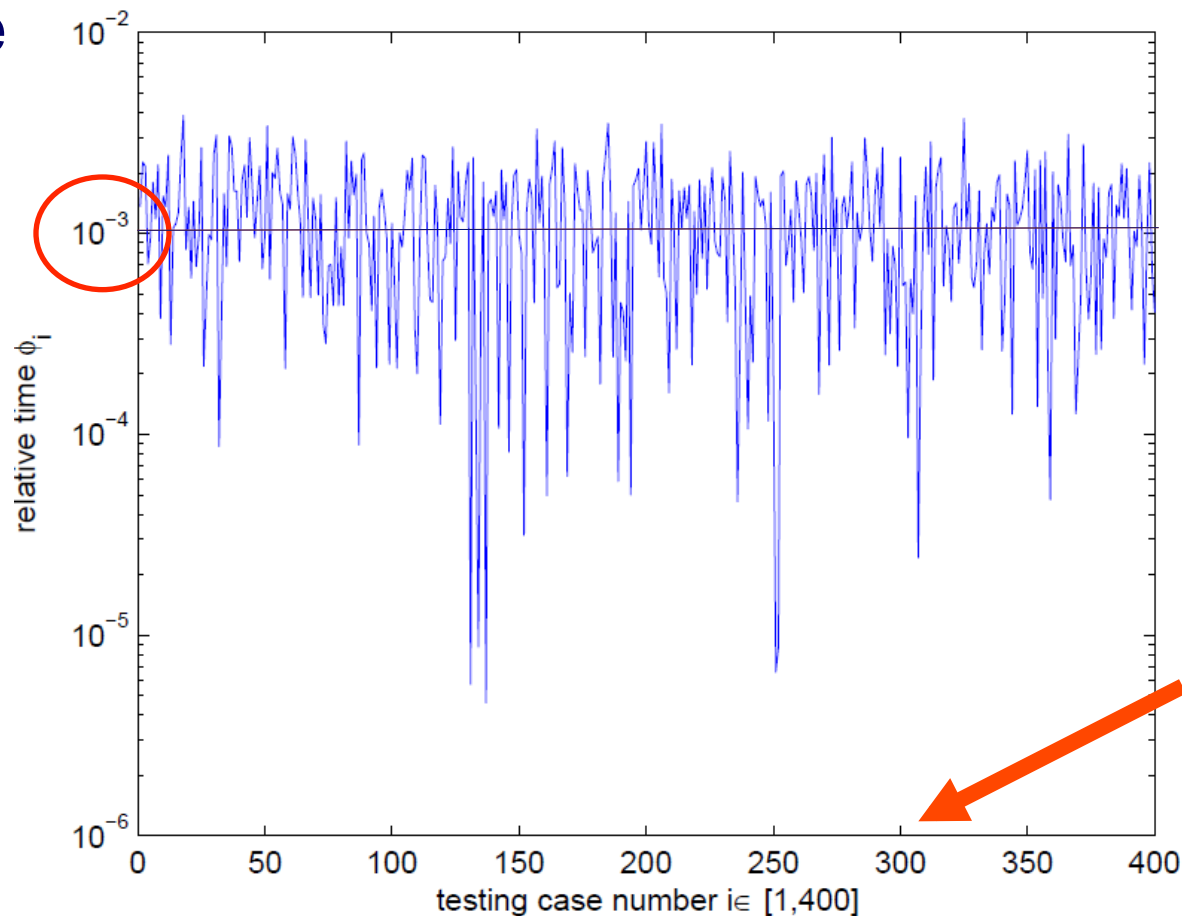- ## Same accuracy, **100x – 100K x faster**

Relative
Search
Time:

## 1000 x

1 sec
vs 15'
vs 1 day



Many
simulations

[1] Yang Weng, Christos Faloutsos, Marija D. Ili´c, and Rohit Negi, Speed up of Data-Driven State Estimation Using Low-Complexity Indexing Method, IEEE PES-General Meeting, (accepted), 2014

# Step1: Reducing dimensionality M

Measurement 1

Measurement M

time

But sequential scan
Is slow, too (MxT)
Can we do better?

A: yes!
We can reduce both
• T, and
• M

SVD

# Step2: Faster than T timeticks
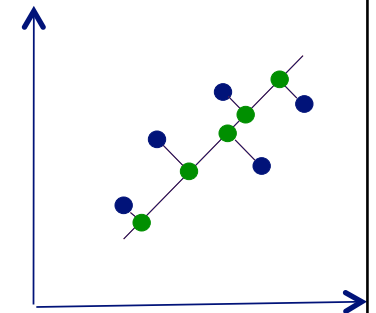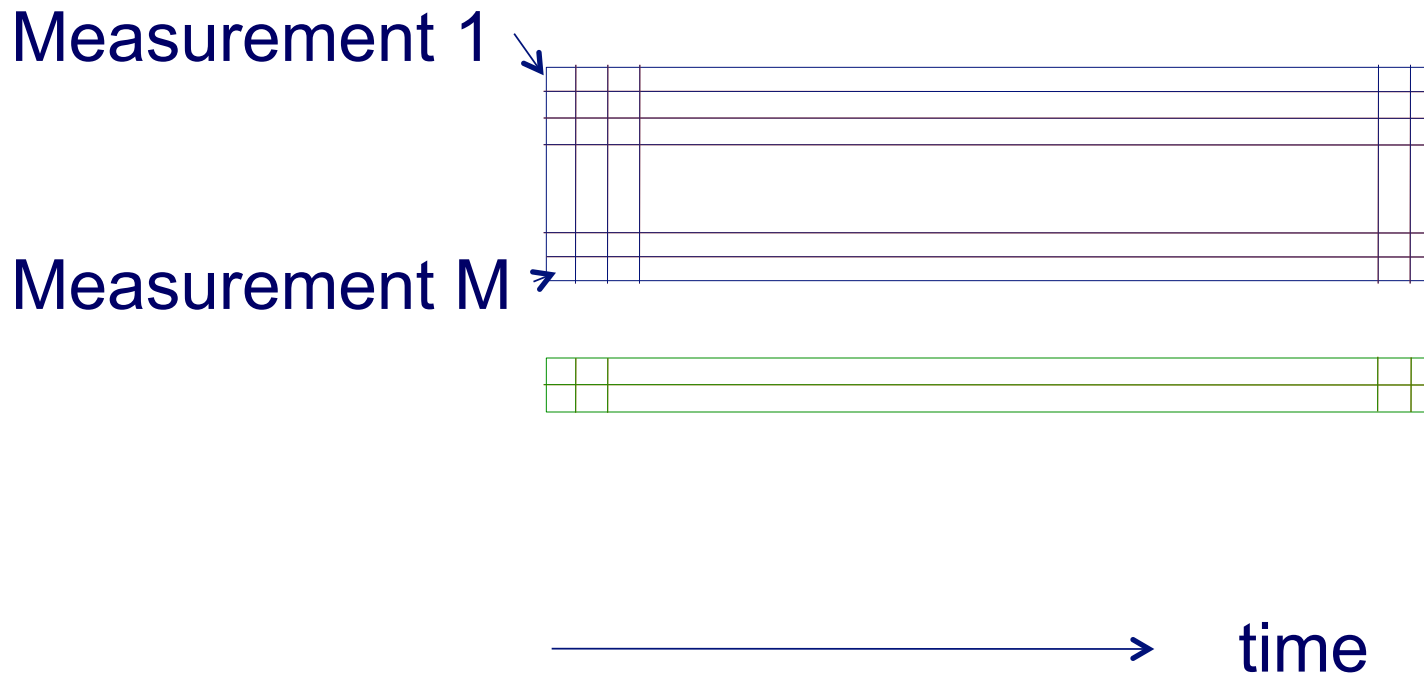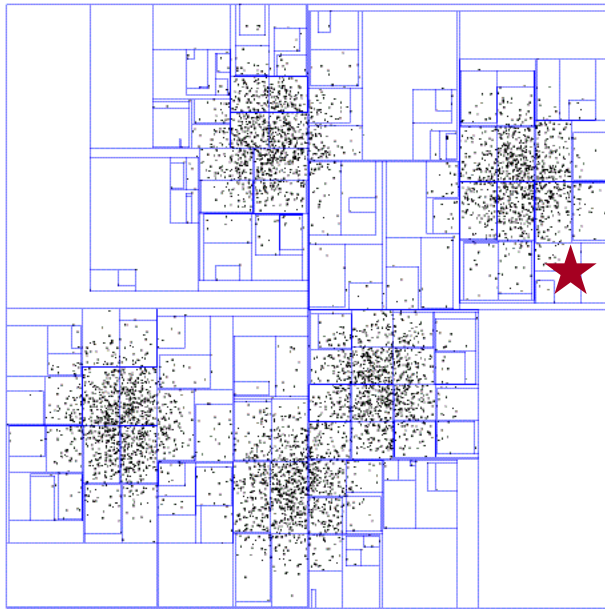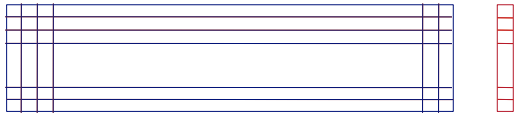
Measurement 1

Measurement M

time

But sequential scan
Is slow, too (MxT)
Can we do better?
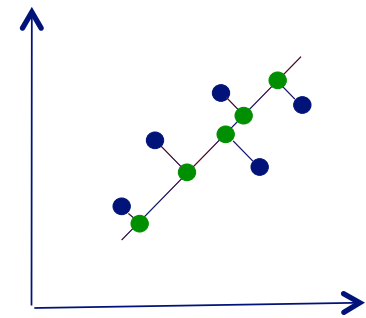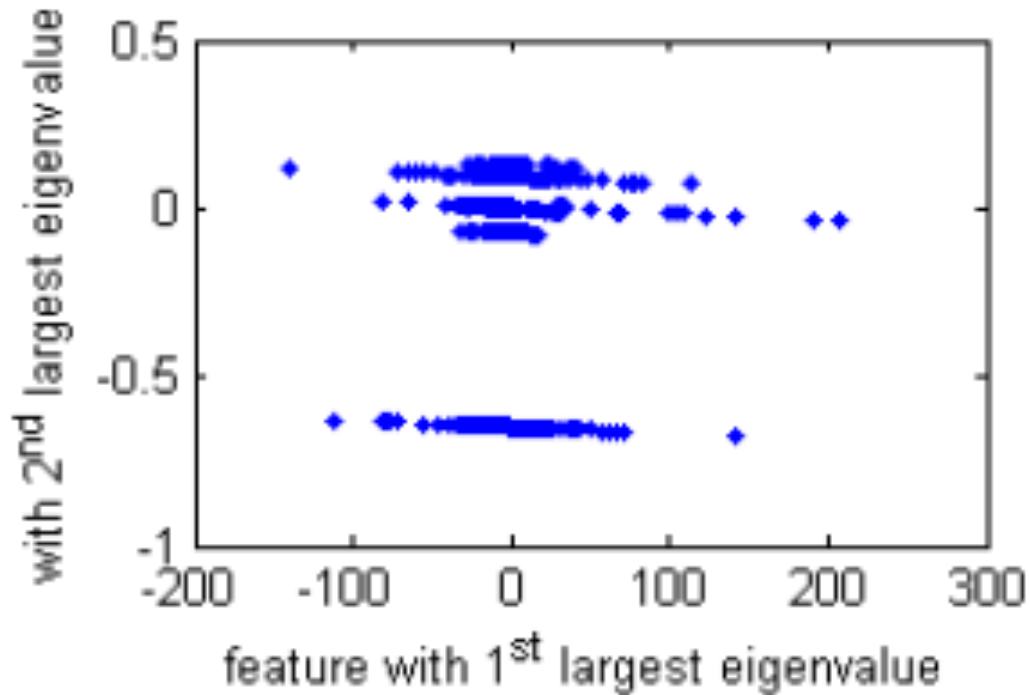
A: yes!
We can reduce both
- T, and
- M

K-d trees
SVD

# Faster than seq. scan: K-d trees

# Thanks to SVD: VISUALIZATION!



- Projection of measurements on to singular vectors of measurement matrix

[1] Yang Weng, Christos Faloutsos, Marija D. Ili´c, and Rohit Negi, Speed up of Data-Driven State Estimation Using Low-Complexity Indexing Method, IEEE PES-General Meeting, (accepted), 2014

# Thanks to SVD: VISUALIZATION!



4 (or 5) groups of behavior!

● Projection of measurements on to singular vectors of measurement matrix

[1] Yang Weng, Christos Faloutsos, Marija D. Ili´c, and Rohit Negi, Speed up of Data-Driven State Estimation Using Low-Complexity Indexing Method, IEEE PES-General Meeting, (accepted), 2014
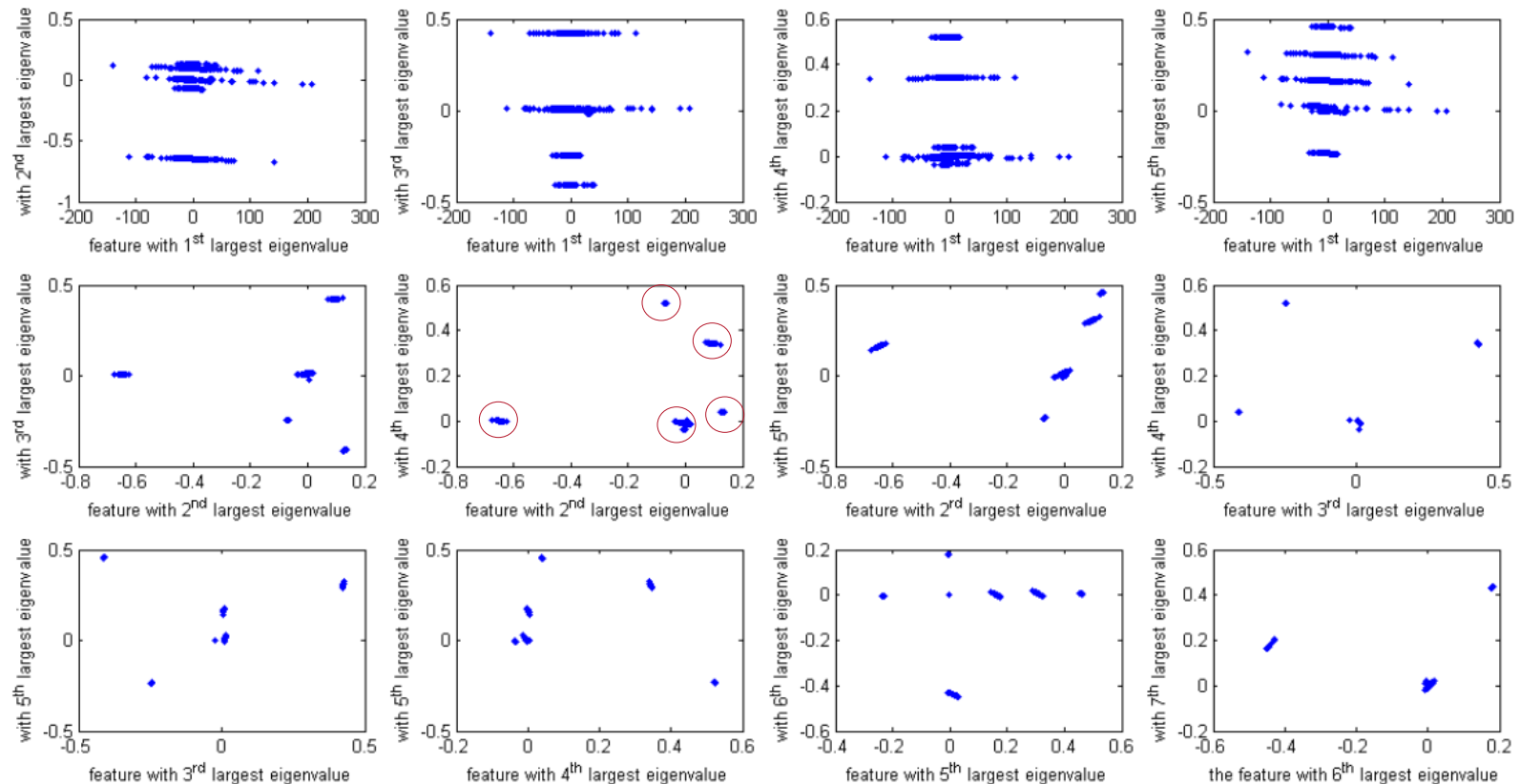
# Thanks to SVD: VISUALIZATION!



- Projection of measurements on to singular vectors of measurement matrix

[1] Yang Weng, Christos Faloutsos, Marija D. Ili´c, and Rohit Negi, Speed up of Data-Driven State Estimation Using Low-Complexity Indexing Method, IEEE PES-General Meeting, (accepted), 2014

# Roadmap



- A case for cross-disciplinarity
➡ - Introduction – Motivation
  - Why study (big) graphs?
- Part#1: Patterns in graphs
- Part#2: Cascade analysis
- Conclusions

# Graphs - why should we care?

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

(c) 2014, C. Faloutsos 18

# Graphs - why should we care?



Food Web
[Martinez '91]

>$10B revenue

>0.5B users



Internet Map
[lumeta.com]

# Graphs - why should we care?

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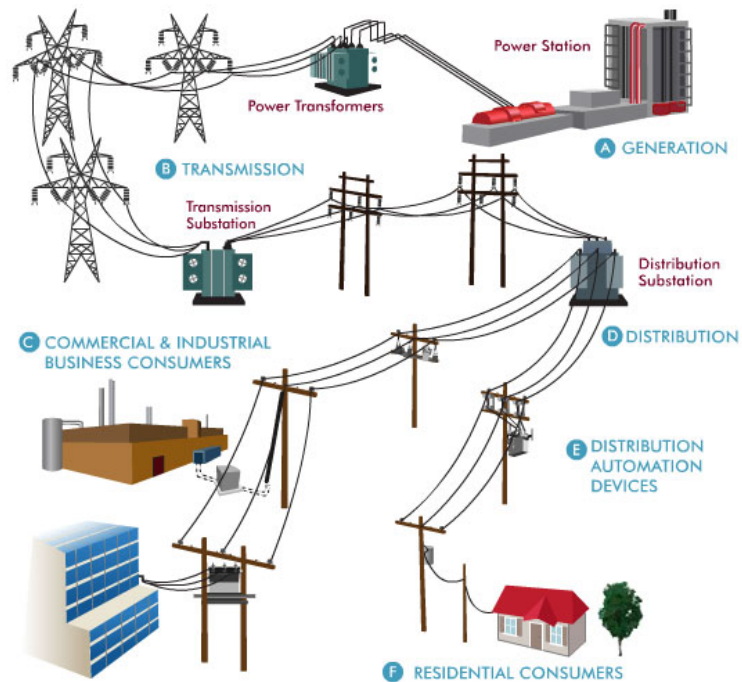- ....

- Many-to-many db relationship -> graph

# Roadmap

- A case for cross-disciplinarity
- Introduction – Motivation
  - Why study (big) graphs?
- ➡ Part#1: Patterns in graphs
- Part#2: Cascade analysis
- Conclusions

# Part 1: Patterns & Laws

# Laws and patterns

- Q1: Are real graphs random?

# Laws and patterns

- Q1: Are real graphs random?
- A1: NO!!
  - Diameter
  - in- and out- degree distributions
  - other (surprising) patterns
- Q2: why so many power laws?
- A2: <self-similarity – stay tuned>

- So, let's look at the data

# Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

**internet domains**

**att.com**

log(degree)

**ibm.com**

log(rank)

# Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

**internet domains**

# Solution# S.1

- Q: So what?

**internet domains**

**att.com**

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

**ibm.com**

**-0.82**

log(rank)

(c) 2014, C. Faloutsos

# Solution# S.1

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

**internet domains**

**att.com**

log(degree)
"../TESTDATA/980410-INTER/980410.Internet.outdegrees.z"
exp(6.65065) * x ** ( -0.626118 )

**ibm.com**

**-0.82**

log(rank)

(c) 2014, C. Faloutsos

# **Gaussian trap**

# **Solution# S.1**

- Q: So what?
- A1: # of two-step-away pairs: $O(d\_max \text{ }^2) \sim 10M^2$

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

**internet domains**

**att.com**

log(degree)

**ibm.com**

-0.82

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

~0.8PB ->
a data center(!)

**DCO @ CMU**

**Gaussian trap**

# Solution# S.1

- Q: So what?

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

$\Downarrow$

~0.8PB ->
a data center(!)

**-0.82**

Such patterns ->
New algorithms

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

Eigenvalue



Exponent = slope

$E = -0.48$

May 2001

$A\ x = \lambda\ x$

Rank of decreasing eigenvalue

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

# Roadmap



- Introduction – Motivation
- Problem#1: Patterns in graphs
  - Static graphs
    - degree, diameter, eigen,
    - Triangles
  - Time evolving graphs
- Problem#2: Tools

# 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

slope 1.68 · slope -1.68

SN

slope 1.74 · slope -1.73

Epinions

slope 1.61 · slope -1.59

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

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

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

(c) 2014, C. Faloutsos

**details**

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

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

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

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

(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]

# **Roadmap**



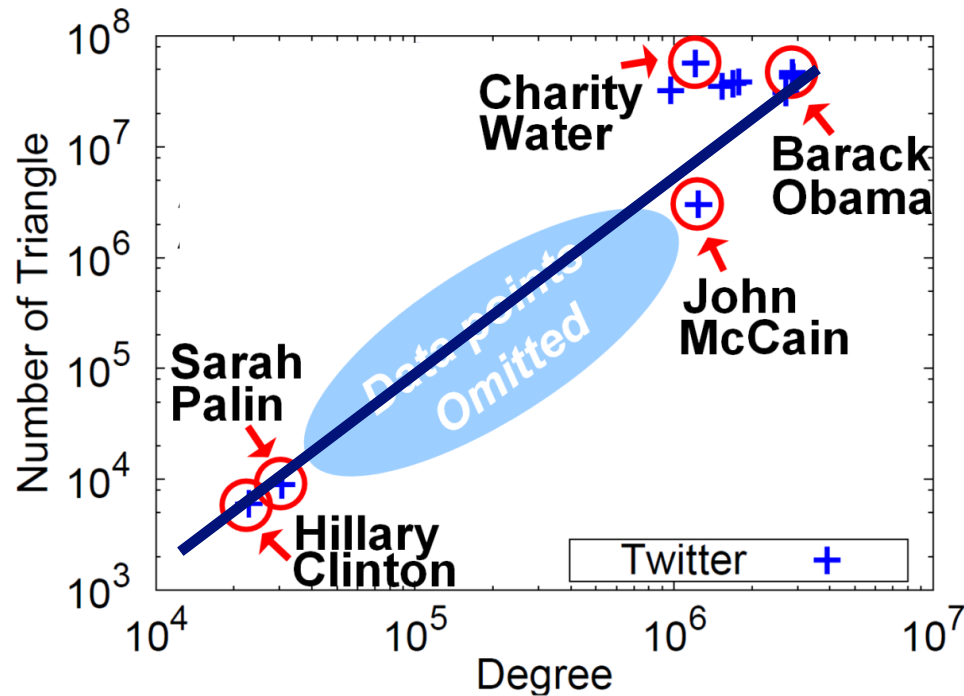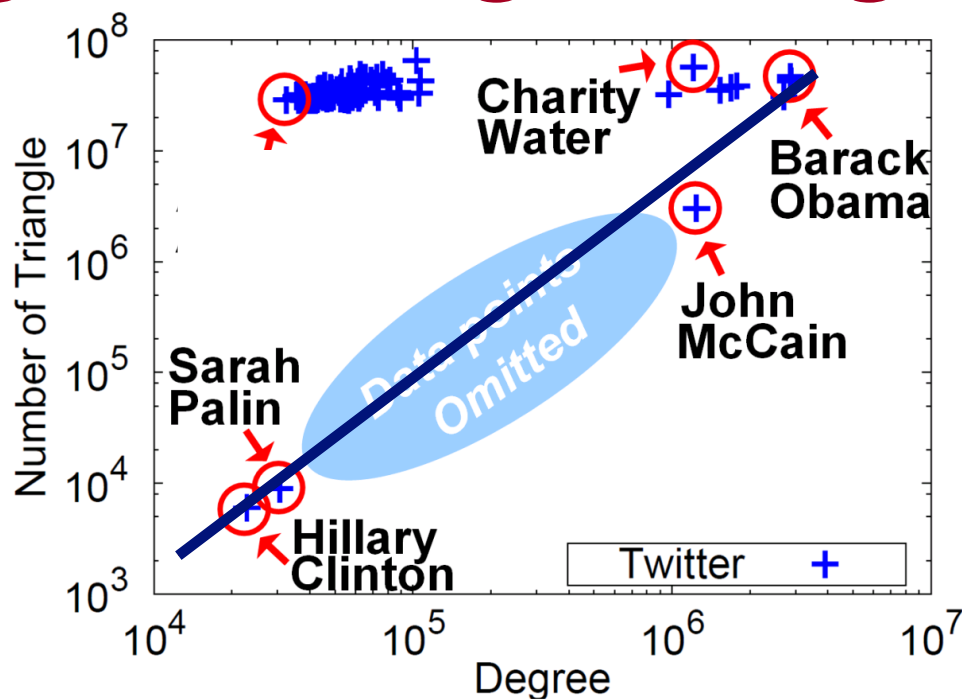- A case for cross-disciplinarity
- Introduction – Motivation
- Part#1: Patterns in graphs
  - Static graphs
  - → Time evolving graphs
- Part#2: Cascade analysis
- Conclusions

# Problem: Time evolution

- with Jure Leskovec (CMU -> Stanford)



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



Jure Leskovec, Jon Kleinberg and Christos Faloutsos: *Graphs over Time: Densification Laws, Shrinking Diameters and Possible Explanations*, KDD 2005

# T.1 Evolution of the Diameter

- Prior work on Power Law graphs hints at **slowly growing diameter**:
  - [diameter ~ O( $N^{1/3}$)]
  - diameter ~ O(log N)
  - diameter ~ O(log log N)
- What is happening in real data?

diameter

# T.1 Evolution of the Diameter

- Prior work on Power Law graphs hints at **slowly growing diameter**:
    - [diameter ~ O( $N^{1/3}$ )]
    - diameter ~ O (log N)
    - diameter ~ O(log log N)
- What is happening in real data?
- Diameter **shrinks** over time

# T.1 Diameter – "Patents"

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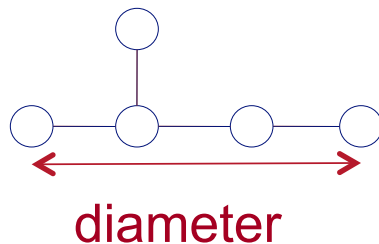- 25 years of data
- @1999
  - 2.9 M nodes
  - 16.5 M edges



diameter

time [years]

# T.2 Temporal Evolution of the Graphs

- N(t) … nodes at time t
- E(t) … edges at time t
- Suppose that

  $$N(t+1) = 2 * N(t)$$

  Say, $k$ friends on average

- Q: what is your guess for

  $$E(t+1) =? 2 * E(t)$$



$k$

# 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

    $$E(t+1) =?\ 2 * E(t)$$

- A: over-doubled! ~ 3x

    – But obeying the ``Densification Power Law''

**Gaussian trap**

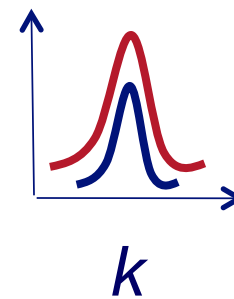Say, *k* friends on average

(c) 2014, C. Faloutsos

# 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
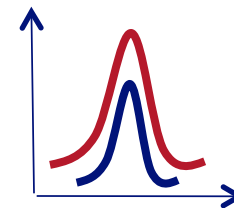
  $$E(t+1) =?~2 * E(t)$$

- A: over-doubled! ~ 3x

  – But obeying the ``Densification Power Law''

**Gaussian trap**

Say, *k* friends on average

(c) 2014, C. Faloutsos

# T.2 Densification – Patent Citations

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

E(t)

1.66

N(t)

Number of edges

* Edges
— = 0.0002 x$^{1.66}$ R$^2$=0.99

Number of nodes

1999

1975

# MORE Graph Patterns

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

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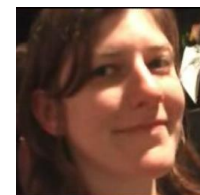
# MORE Graph Patterns

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

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

# MORE Graph Patterns

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

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

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

# Roadmap



- A case for cross-disciplinarity
- Introduction – Motivation
- Part#1: Patterns in graphs
  - …
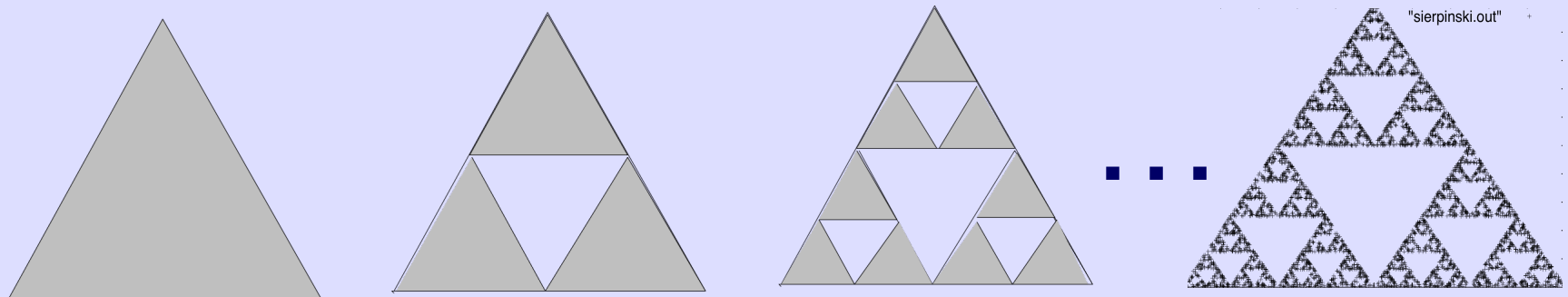  - Why so many power-laws?
- Part#2: Cascade analysis
- Conclusions

# Why so many P.L.?

- Possible answer: self-similarity / fractals

(c) 2014, C. Faloutsos

# 20'' intro to fractals

- Remove the middle triangle; repeat

- -> Sierpinski triangle

- (Bonus question - dimensionality?
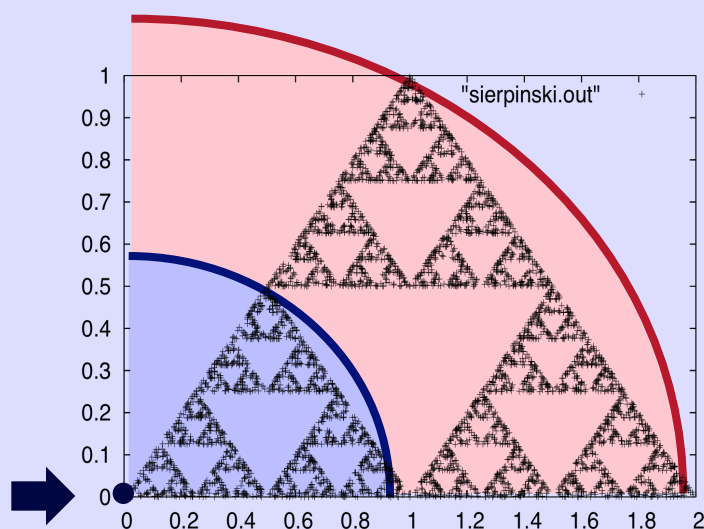  - $>1$ (inf. perimeter – $(4/3)^\infty$ )
  - $<2$ (zero area – $(3/4)^\infty$ )

"sierpinski.out"

. . .

# 20'' intro to fractals

Self-similarity -> no char. scale

-> power laws, eg:

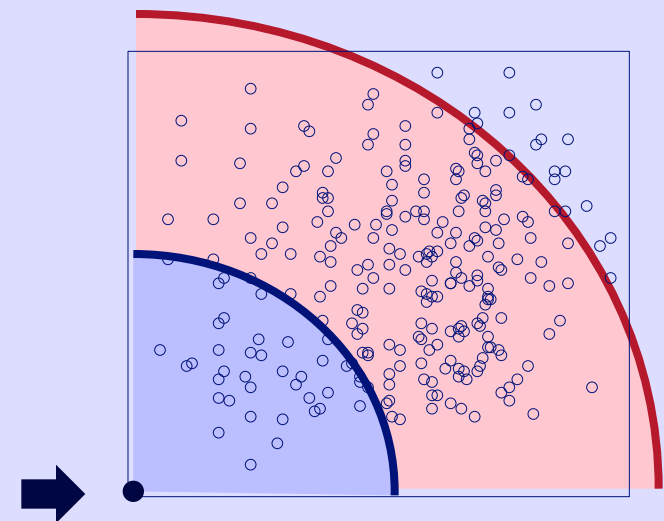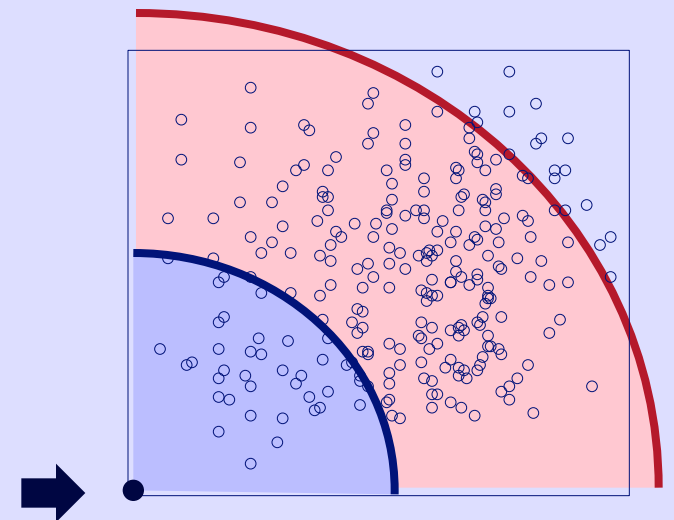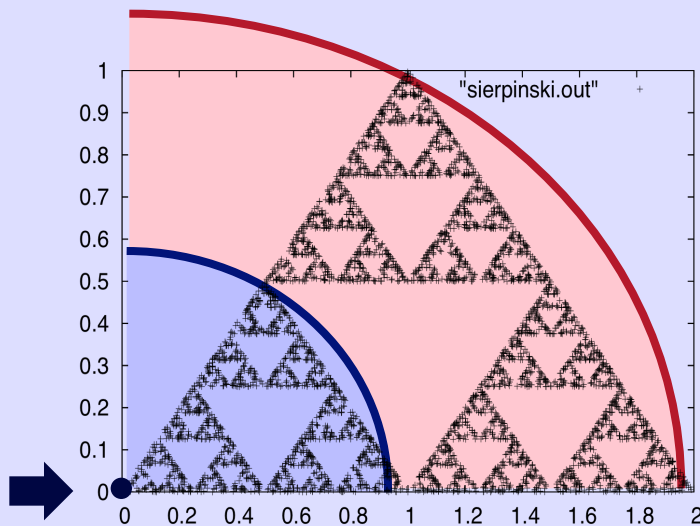2x the radius,

3x the #neighbors nn(r)

$$nn(r) = C\ r^{\,log3/log2}$$



"sierpinski.out"

# 20'' intro to fractals

Self-similarity -> <u>no char. scale</u>
-> power laws, eg:
2x the radius,
3x the #neighbors nn(r)

$$nn(r) = C\ r^{\log 3/\log 2}$$



"sierpinski.out"

# 20'' intro to fractals

Self-similarity -> no char. scale
-> power laws, eg:
2x the radius,
3x the #neighbors
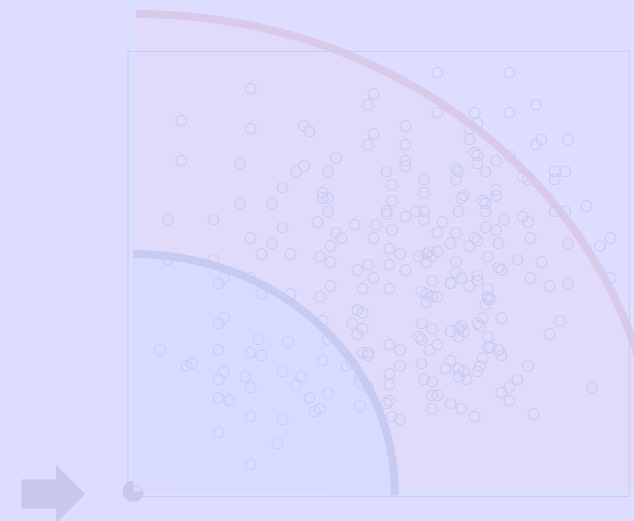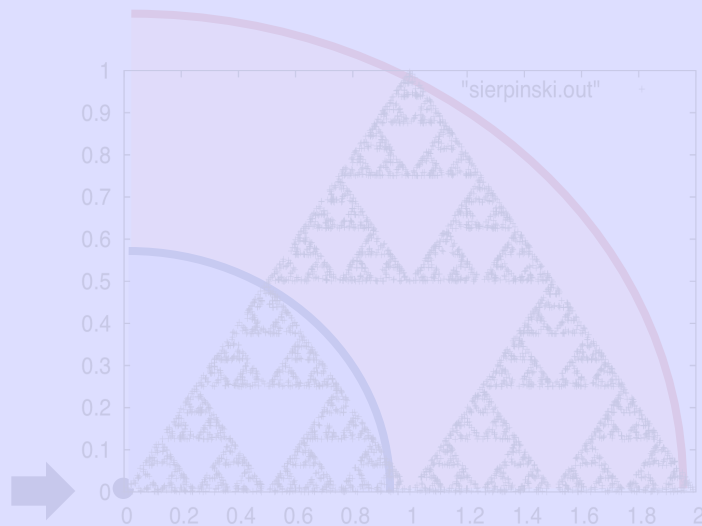$$nn = C \, r^{\,\log 3/\log 2}$$

Reminder:
Densification P.L.
(2x nodes, ~3x edges)

# 20'' intro to fractals

Self-similarity -> no char. scale
-> power laws, eg:
2x the radius,
3x the #neighbors
$nn = C\ r^{\ \log3/\log2}$

2x the radius,
4x neighbors
$nn = C\ r^{\ \log4/\log2} = C\ r^{\ 2}$

# 20'' intro to fractals

Self-similarity -> no char. scale
-> power laws, eg:

2x the radius,

3x the #neighbors

$$nn = C \, r^{\log 3/\log 2} \quad =1.58$$

Fractal dim.

2x the radius,

4x neighbors

$$nn = C \, r^{\log 4/\log 2} = C \, r^2$$



"sierpinski.out"

# 20'' intro to fractals

**Self-similarity** -> no char. scale
-> **power laws**, eg:

2x the radius,
3x the #neighbors
   $nn = C \, r^{\log 3/\log 2}$

2x the radius,
4x neighbors
   $nn = C \, r^{\log 4/\log 2} = C \, r^2$

Fractal dim.



"sierpinski.out"

# How does self-similarity help in graphs?

- A: RMAT/Kronecker generators
  - With self-similarity, we get all power-laws, automatically,
  - And small/shrinking diameter
  - And `no good cuts'

*R-MAT: A Recursive Model for Graph Mining*, by D. Chakrabarti, Y. Zhan and C. Faloutsos, SDM 2004, Orlando, Florida, USA

*Realistic, Mathematically Tractable Graph Generation and Evolution, Using Kronecker Multiplication,* by J. Leskovec, D. Chakrabarti, J. Kleinberg, and C. Faloutsos, in PKDD 2005, Porto, Portugal

# Graph gen.: Problem dfn

- Given a growing graph with count of nodes $N_1$, $N_2$, ...
- Generate a realistic sequence of graphs that will obey all the patterns
  - Static Patterns
    - S1 Power Law Degree Distribution
    - S2 Power Law eigenvalue and eigenvector distribution
      - Small Diameter
  - Dynamic Patterns
    - T2 Growth Power Law (2x nodes; 3x edges)
    - T1 Shrinking/Stabilizing Diameters

# Kronecker Graphs



$X_1$

$X_2$

$X_3$

| 1 | 1 | 0 |
|---|---|---|
| 1 | 1 | 1 |
| 0 | 1 | 1 |

$G_1$

Adjacency matrix

# Kronecker Graphs



Intermediate stage

$$\begin{array}{|c|c|c|} \hline 1 & 1 & 0 \\ \hline 1 & 1 & 1 \\ \hline 0 & 1 & 1 \\ \hline \end{array}$$

$G_1$

Adjacency matrix

# Kronecker Graphs



$X_1$  $X_2$  $X_3$

$X_{1,1}$  $X_{1,2}$  $X_{1,3}$
$X_{2,1}$  $X_{2,3}$
$X_{3,1}$  $X_{3,2}$  $X_{3,3}$

Intermediate stage

| 1 | 1 | 0 |
|---|---|---|
| 1 | 1 | 1 |
| 0 | 1 | 1 |

$G_1$

| $G_1$ | $G_1$ | 0 |
|---|---|---|
| $G_1$ | $G_1$ | $G_1$ |
| 0 | $G_1$ | $G_1$ |

$$G_2 = G_1 \otimes G_1$$

Adjacency matrix                     Adjacency matrix

# Kronecker Graphs

- Continuing multiplying with $G_1$ we obtain $G_4$ and so on …



$G_4$ adjacency matrix

# Kronecker Graphs

- Continuing multiplying with $G_1$ we obtain $G_4$ and so on …



$G_4$ adjacency matrix

(c) 2014, C. Faloutsos

# Kronecker Graphs

- Continuing multiplying with $G_1$ we obtain $G_4$ and so on …

$G_4$ adjacency matrix

(c) 2014, C. Faloutsos

# Kronecker Graphs

- Continuing multiplying with $G_1$ we obtain $G_4$ and so on …

Holes within holes;
Communities
within communities

$G_4$ adjacency matrix

(c) 2014, C. Faloutsos

71

# Properties:

- We can PROVE that
    - Degree distribution is multinomial ~ power law

new
    - Diameter: constant

    - Eigenvalue distribution: multinomial

    - First eigenvector: multinomial

# Problem Definition

- Given a growing graph with nodes $N_1, N_2, ...$
- Generate a realistic sequence of graphs that will obey all the patterns
  - Static Patterns
    - ✓ Power Law Degree Distribution
    - ✓ Power Law eigenvalue and eigenvector distribution
    - ✓ Small Diameter
  - Dynamic Patterns
    - ✓ Growth Power Law
    - ✓ Shrinking/Stabilizing Diameters
- First generator for which we can **prove** all these properties

# Impact: Graph500

- Based on RMAT (= 2x2 Kronecker)
- Standard for graph benchmarks
- http://www.graph500.org/
- Competitions 2x year, with all major entities: LLNL, Argonne, ITC-U. Tokyo, Riken, ORNL, Sandia, PSC, …

*To iterate is human, to recurse is devine*

*R-MAT: A Recursive Model for Graph Mining*,
by D. Chakrabarti, Y. Zhan and C. Faloutsos,
SDM 2004, Orlando, Florida, USA

# **Summary of Part#1**

- \*many\* patterns in real graphs
  - Small & shrinking diameters
  - Power-laws everywhere
  - Gaussian trap
- Self-similarity (RMAT/Kronecker): good model

# Roadmap



- A case for cross-disciplinarity
- Introduction – Motivation
- Part#1: Patterns in graphs
- ➡ Part#2: Cascade analysis
- Conclusions

# Comic relief:

- What would a barefooted man get if he steps on an electric wire?



http://energyquest.ca.gov/games/jokes/george.html

# **Comic relief:**



- What would a barefooted man get if he steps on an electric wire?
  (Answer) A pair of *shocks*



http://energyquest.ca.gov/games/jokes/george.html

**Carnegie Mellon**

# Part 2: 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



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

[US-MEDICARE NETWORK 2005]

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

# Fractional Asymmetric Immunization

Drug-resistant Bacteria
(like XDR-TB)

Hospital

Another
Hospital

(c) 2014, C. Faloutsos

# **Fractional Asymmetric Immunization**



Hospital

Another Hospital

# Fractional Asymmetric Immunization

Hospital

Another Hospital

# Fractional Asymmetric Immunization

**Problem**:

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

Hospital

Another Hospital

# Fractional Asymmetric Immunization

**Problem**:

Given k units of disinfectant, distribute them

to maximize hospitals saved @ 365 days

Hospital

Another
Hospital

(c) 2014, C. Faloutsos

# Straightforward solution:

Simulation:

1. Distribute resources

2. 'infect' a few nodes

3. Simulate evolution of spreading
   – (10x, take avg)

4. Tweak, and repeat step 1

# Straightforward solution:

Simulation:

1. Distribute resources

2. 'infect' a few nodes

3. Simulate evolution of spreading
   – (10x, take avg)

4. Tweak, and repeat step 1

# **Straightforward solution:**

Simulation:

1. Distribute resources

2. 'infect' a few nodes

3. Simulate evolution of spreading

   – (10x, take avg)

4. Tweak, and repeat step 1

# Straightforward solution:

Simulation:

1. Distribute resources

2. 'infect' a few nodes

3. Simulate evolution of spreading

   – (10x, take avg)

➡ 4. Tweak, and repeat step 1

**Running Time**

Wall-Clock Time

better

> 1 week

> 30,000x speed-up!

14 secs

Simulations    SMART-ALLOC

# Experiments

# infected

uniform

$\downarrow$ *better*

~6x

SMART-ALLOC

0    100    200    300    400
Time ticks (days)

# epochs

$K = 120$

# 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)?

> 30,000x
speed-
up!

14
secs

– 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

Avg degree
Max degree
Diameter
Modularity
'Conductance'

Q1: why??

(Q2: dfn & intuition of eigenvalue ? )

# 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 = effective strength*
- *s < 1 : below threshold*



| Models | Effective Strength (s) | Threshold (tipping point) |
|---|---|---|
| SIS, SIR, SIRS, SEIR | $s = \lambda \left( \dfrac{\beta}{\delta} \right)$ | $s = 1$ |
| SIV, SEIV | $s = \lambda \cdot \left( \dfrac{\beta\gamma}{\delta(\gamma + \theta)} \right)$ | |
| $SI_1I_2V_1V_2$ (**H.I.V.**) | $s = \lambda \cdot \left( \dfrac{\beta_1 v_2 + \beta_2 \varepsilon}{v_2(\varepsilon + v_1)} \right)$ | |

# Our thresholds for some models

- *s = effective strength*
- *s < 1 : below threshold*

| Models | Effective Strength | Threshold (tipping point) |
|---|---|---|
| SIS, SIR, SIRS, SEIR | $s = \lambda \cdot \left(\dfrac{\beta}{\delta}\right)$ | |
| SIV, SEIV | $s = \lambda \cdot \left(\dfrac{\beta\gamma}{\delta(\gamma + \theta)}\right)$ | $s = 1$ |
| $SI_1I_2V_1V_2$ (**H.I.V.**) | $s = \lambda \cdot \left(\dfrac{\beta_1 v_2 + \beta_2 \varepsilon}{v_2(\varepsilon + v_1)}\right)$ | |

No immunity

Temp. immunity

w/ incubation

# Roadmap



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

# Intuition for λ

## "Official" definitions:

- *Let **A** be the adjacency matrix. Then λ is the root with the largest magnitude of the characteristic polynomial of A [det(A − λI)].*
- Also: $\mathbf{A}\,\mathbf{x} = \lambda\,\mathbf{x}$

Neither gives much intuition!

## "Un-official" Intuition

- For 'homogeneous' graphs, $\lambda == degree$

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

# Largest Eigenvalue (λ)

## better connectivity ⟶ higher λ



$\lambda \approx 2$                    $\lambda = \sqrt{N}$                    $\lambda = N-1$

(a)Chain                    (b)Star                    (c)Clique

$\lambda \approx 2$                    $\lambda = 31.67$                    $\lambda = 999$

*N* = 1000 nodes

# Largest Eigenvalue (λ)

better connectivity ⟶ higher λ

$\lambda \approx 2$

(a)Chain

$\lambda = \sqrt{N}$

(b)Star

$\lambda = N-1$

(c)Clique

$\lambda \approx 2$

$\lambda = 31.67$

$\lambda = 999$

N = 1000 nodes

# 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*

# 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
- Part#1: Patterns in graphs
- Part#2: Cascade analysis
  - (Fractional) Immunization
  - Epidemic thresholds
- ➡ Acks & Conclusions

# **Thanks**

# Project info: PEGASUS

**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**

**Beutel, Alex**

**Chau, Polo**

**Kang, U**

**Koutra, Danai**

**McGlohon, Mary**

**Prakash, Aditya**

**Papalexakis, Vagelis**

**Tong, Hanghang**

# CONCLUSION#1 – Big data

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

# CONCLUSION#2 – self-similarity

- powerful tool / viewpoint
  - Power laws; shrinking diameters
  - **Gaussian trap** (eg., F.O.F.)
  - RMAT – `graph500` generator

# CONCLUSION#3 – eigen-drop

- Cascades & immunization: G2 theorem & **eigenvalue**

~6x fewer!

[US-MEDICARE NETWORK 2005]

CURRENT PRACTICE

OUR METHOD

> 30,000x speed-up!

14 secs

# References

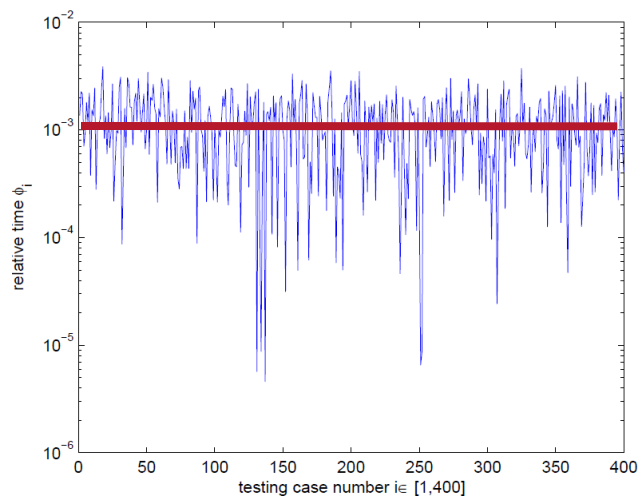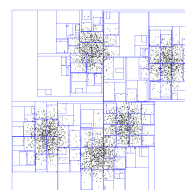- D. Chakrabarti, C. Faloutsos: *Graph Mining – Laws, Tools and Case Studies*, Morgan Claypool 2012
- http://www.morganclaypool.com/doi/abs/10.2200/ S00449ED1V01Y201209DMK006

# TAKE HOME MESSAGE:

# Cross-disciplinarity

**Already started paying off for power grids**
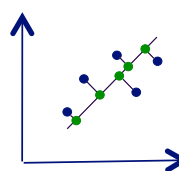
- Same accuracy, **100x – 100K x faster**
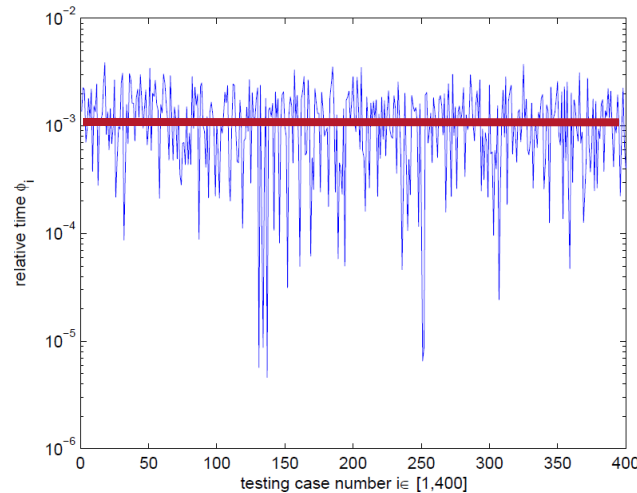
**1000 x**

Kd-tree        SVD

[1] Yang Weng, <u>Christos Faloutsos</u>, Marija D. Ili´c, and Rohit Negi, Speed up of Data-Driven State Estimation Using Low-Complexity Indexing Metho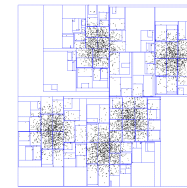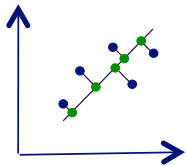d, IEEE PES-General Meeting, (accepted), 2014