# Talk 1: Graph Mining – patterns and generators

Christos Faloutsos
CMU



#### Our goal:

Open source system for mining huge graphs:

PEGASUS project (PEta GrAph mining System)

- www.cs.cmu.edu/~pegasus
- code and papers

Project Pegasus

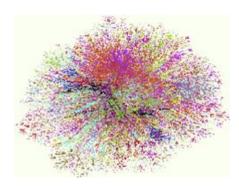
#### **Outline**



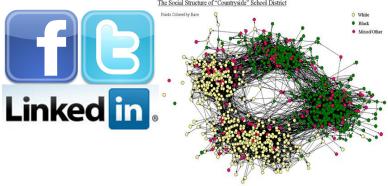
- Introduction Motivation
  - Talk#1: Patterns in graphs; generators
  - Talk#2: Tools (Ranking, proximity)
  - Talk#3: Tools (Tensors, scalability)
  - Conclusions



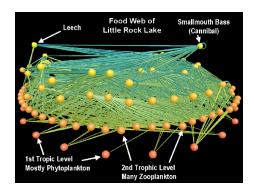
#### Graphs - why should we care?



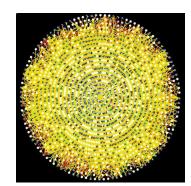
Internet Map [lumeta.com]



Friendship Network [Moody '01]



Food Web [Martinez '91]

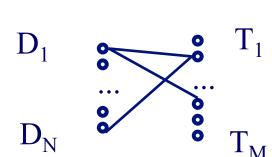


Protein Interactions [genomebiology.com]



#### Graphs - why should we care?

• IR: bi-partite graphs (doc-terms)



web: hyper-text graph

• ... and more:



#### Graphs - why should we care?

- network of companies & board-of-directors members
- 'viral' marketing
- web-log ('blog') news propagation
- computer network security: email/IP traffic and anomaly detection

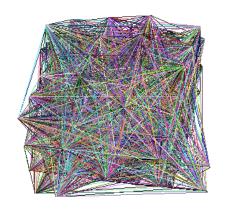
•

#### **Outline**

- Introduction Motivation
- Patterns in graphs
  - Patterns in Static graphs
  - Patterns in Weighted graphs
  - Patterns in Time evolving graphs
  - Generators



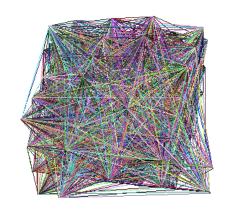
### Network and graph mining



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



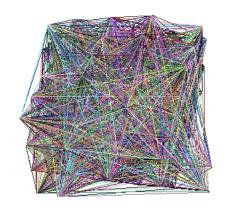
#### Network and graph mining



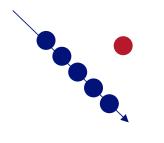
- How does the Internet look like?
- How does FaceBook look like?
- What is 'normal'/'abnormal'?
- which patterns/laws hold?
  - To spot anomalies (rarities), we have to discover patterns



#### Network and graph mining



- How does the Internet look like?
- How 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...

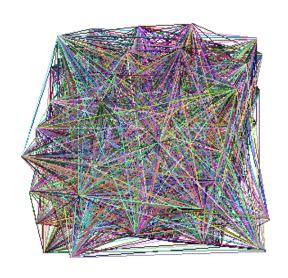


**KAIST-2011** 



## **Topology**

How does the Internet look like? Any rules?



(Looks random – right?)

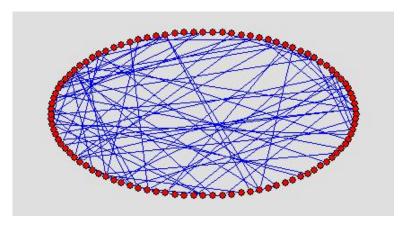


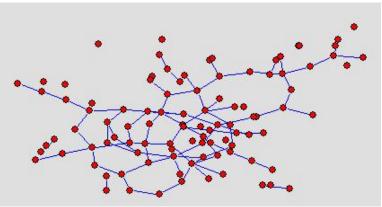
#### Are real graphs random?

- random (Erdos-Renyi)
   graph 100 nodes, avg
   degree = 2
- before layout
- after layout
- No obvious patterns

#### (generated with: pajek

http://vlado.fmf.uni-lj.si/pub/networks/pajek/







### Graph mining

• Are real graphs random?

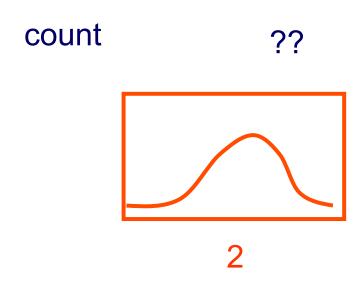
#### Laws and patterns

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



#### Laws – degree distributions

• Q: avg degree is ~2 - what is the most probable degree?

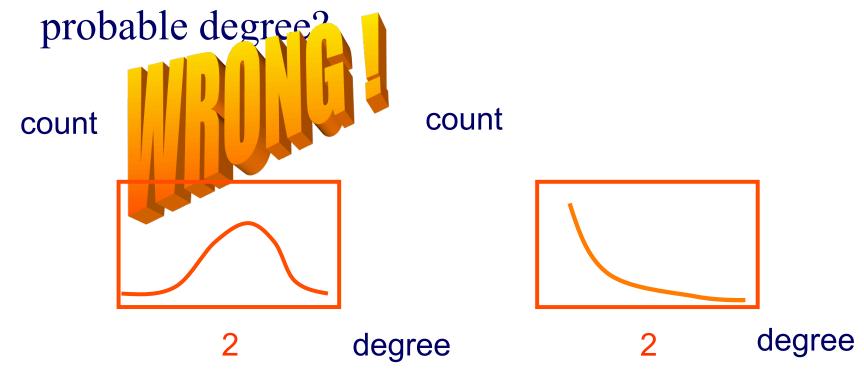


degree



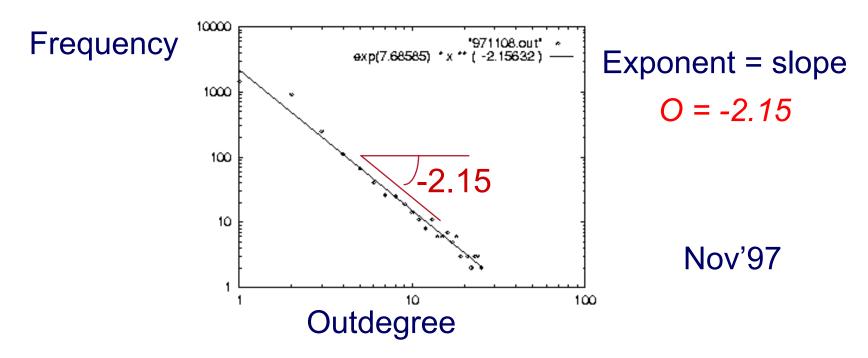
#### Laws – degree distributions

• Q: avg degree is  $\sim 2$  - what is the most





### Solution S1 .Power-law: outdegree O



The plot is linear in log-log scale [FFF'99]

KAIST-2011

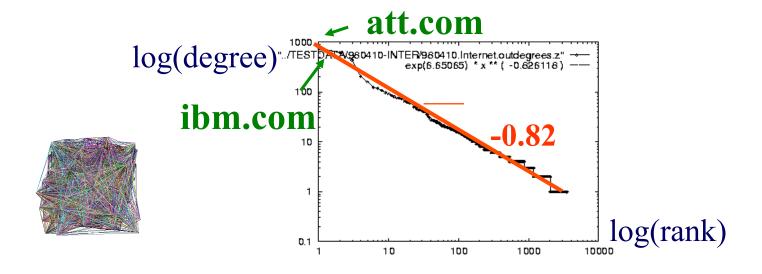
$$freq^{(C)} \stackrel{\text{2011, C}}{=} degree^{(-2.15)}$$



#### Solution# S.1'

• Power law in the degree distribution [SIGCOMM99]

#### internet domains

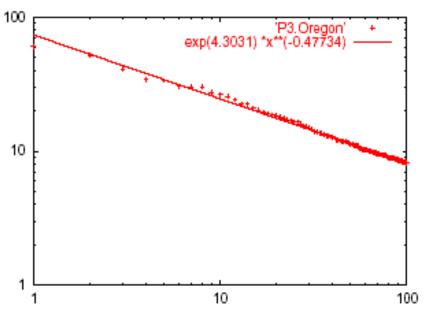


**KAIST-2011** 



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





Exponent = slope

E = -0.48

May 2001

Rank of decreasing eigenvalue

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

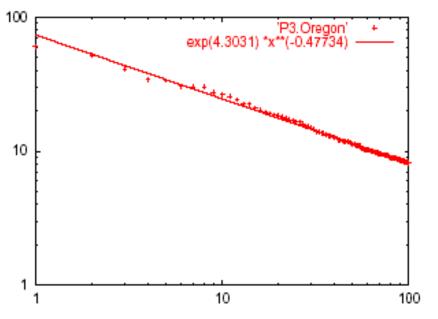
**KAIST-2011** 

(C) 2011, C. Faloutsos



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





Exponent = slope

E = -0.48

May 2001

Rank of decreasing eigenvalue

• [Mihail, Papadimitriou '02]: slope is ½ of rank exponent

KAIST-2011

(C) 2011, C. Faloutsos

20



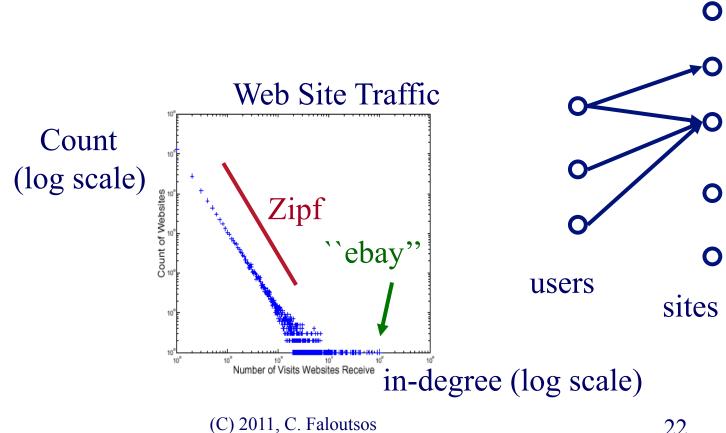
#### **But:**

How about graphs from other domains?



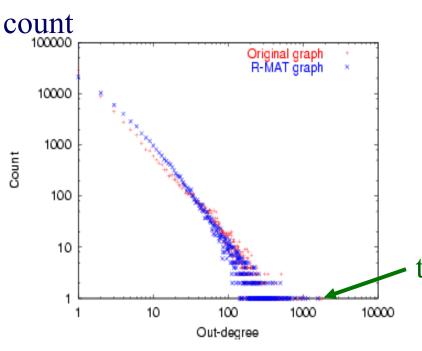
#### More power laws:

• web hit counts [w/ A. Montgomery]



**KAIST-2011** 22

#### epinions.com



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

trusts-2000-people user

(out) degree

#### 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'

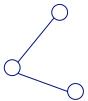
#### **Outline**

- Introduction Motivation
- Patterns in graphs
  - Patterns in Static graphs
    - Degree
    - Triangles
    - •
  - Patterns in Weighted graphs
  - Patterns in Time evolving graphs
- Generators





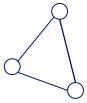
### Solution# S.3: Triangle 'Laws'



Real social networks have a lot of triangles



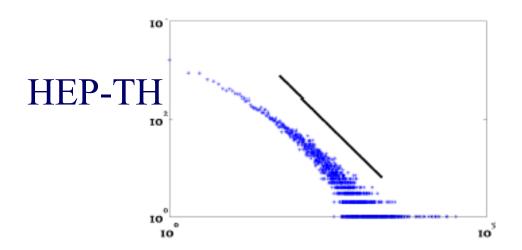
### Solution# S.3: Triangle 'Laws'

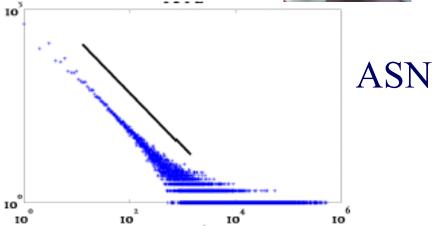


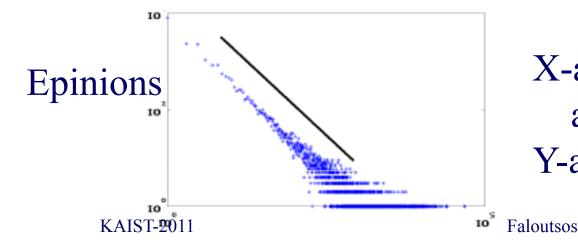
- Real social networks have a lot of triangles
  - Friends of friends are friends
- Any patterns?

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







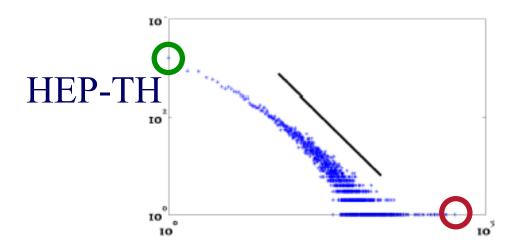


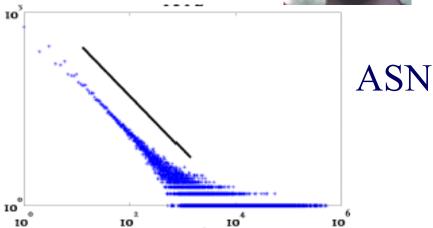
X-axis: # of Triangles a node participates in Y-axis: count of such nodes

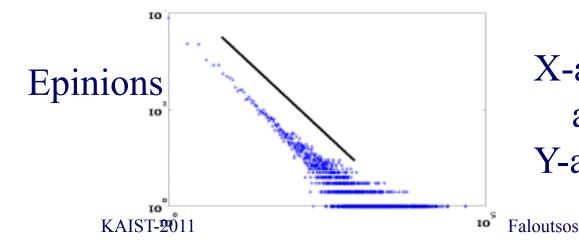
28

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





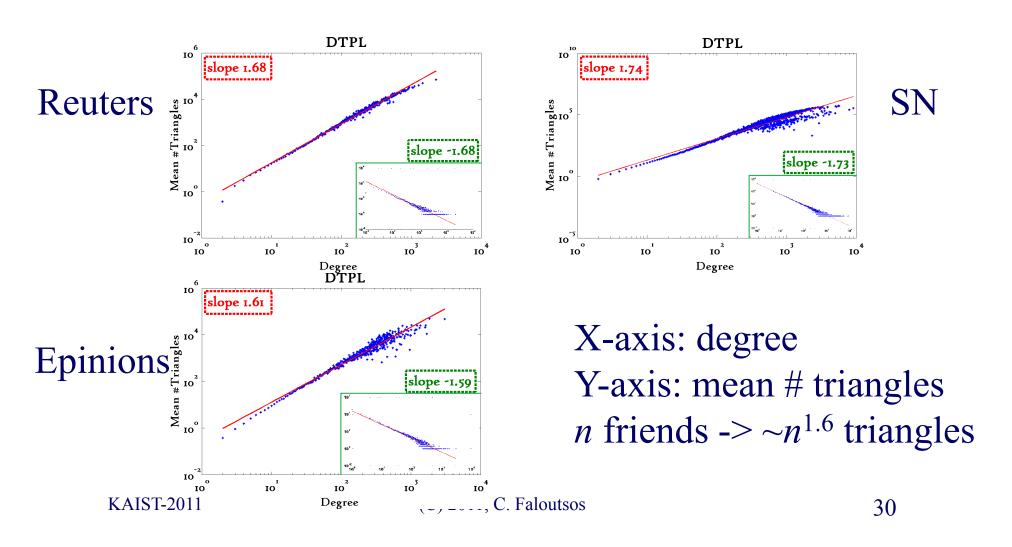




X-axis: # of Triangles a node participates in Y-axis: count of such nodes

29

## Triangle Law: #S.4 [Tsourakakis ICDM 2008]





## Triangle Law: Computations [Tsourakakis ICDM 2008]

details

But: triangles are expensive to compute (3-way join; several approx. algos) Q: Can we do that quickly?



## Triangle Law: Computations [Tsourakakis ICDM 2008]

But: triangles are expensive to compute (3-way join; several approx. algos)

Q: Can we do that quickly?

A: Yes!

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

(and, because of skewness, we only need the top few eigenvalues!

**KAIST-2011** 

(C) 2011, C. Faloutsos

details

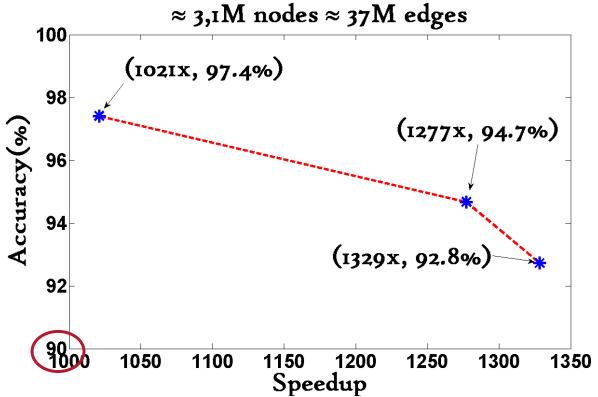




## Triangle Law: Computations

#### [Tsourakakis ICDM 2008]

Wikipedia graph 2006-Nov-04



**KAIST-2011** 

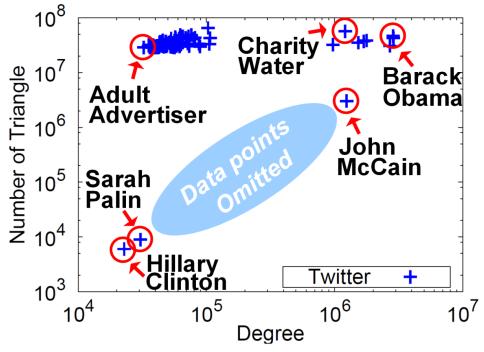


### Triangle counting for large graphs?

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]



## How about cliques?

## Large Human Communication Networks Patterns and a Utility-Driven Generator

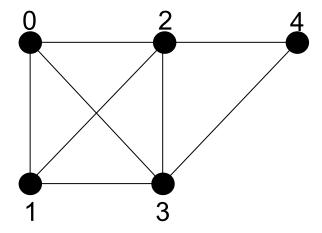
Nan Du, Christos Faloutsos, Bai Wang, Leman Akoglu KDD 2009





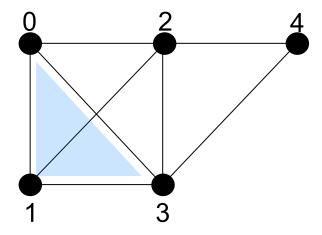
## Cliques

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



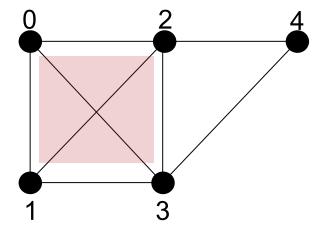
## Clique

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



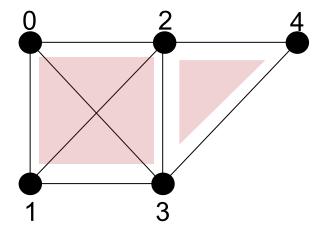
## Clique

- 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;
- {0,1,2,3} and {2,3,4} are the maximal cliques.





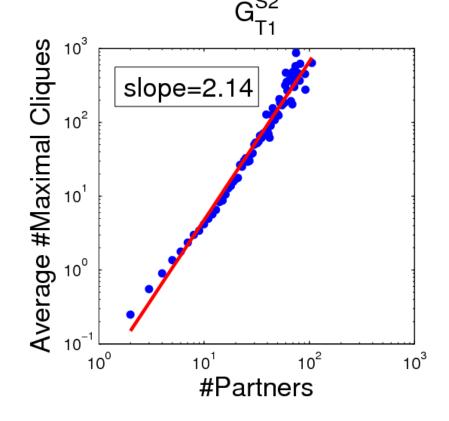
### **S5:** Clique-Degree Power-Law

• Power law:

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

# maximal cliques of node i

degree of node i

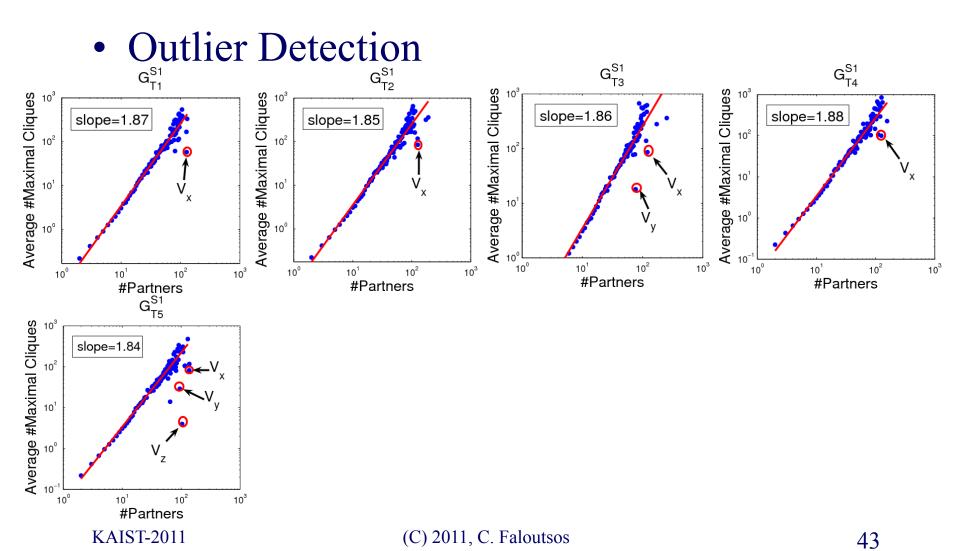


 $\alpha$  is the power law exponent  $\alpha \in [1.8, 2.2]$  for S1~S3

More friends, even more social circles!

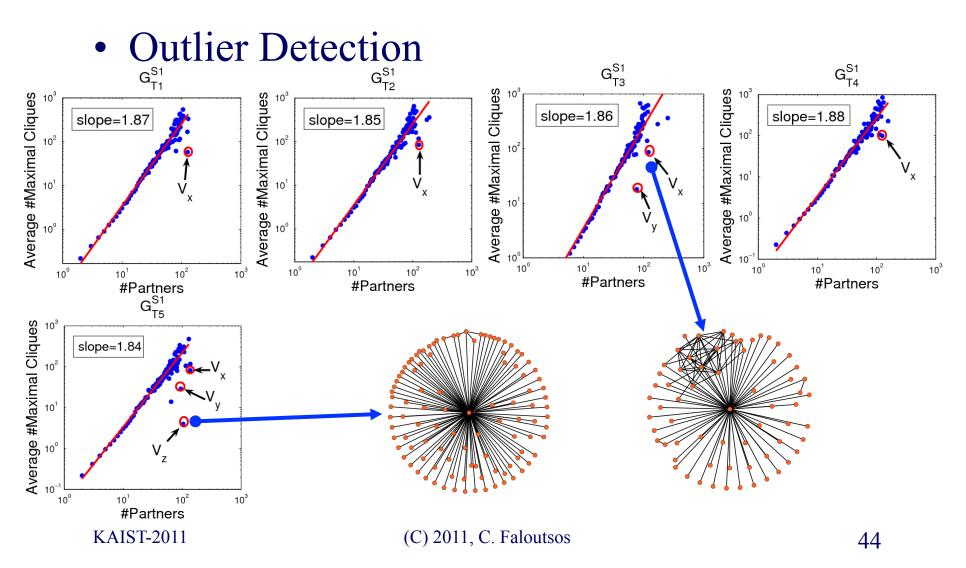


## **S5:** Clique-Degree Power-Law





### **S5:** Clique-Degree Power-Law



#### **Outline**

- Introduction Motivation
- Patterns in graphs
  - Patterns in Static graphs
    - Degree, eigenvalues
    - Triangles, cliques
    - Other observations
  - Patterns in Weighted graphs
  - Patterns in Time evolving graphs
- Generators





Yes!

KAIST-2011 (C) 2011, C. Faloutsos 46

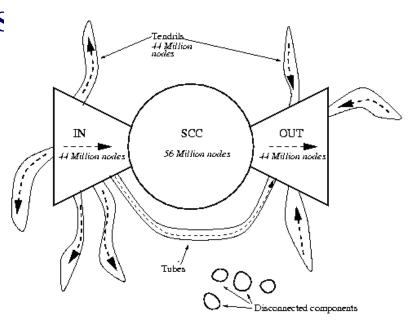
#### Yes!

- Small diameter (~ constant!)
  - six degrees of separation / 'Kevin Bacon'
  - small worlds [Watts and Strogatz]



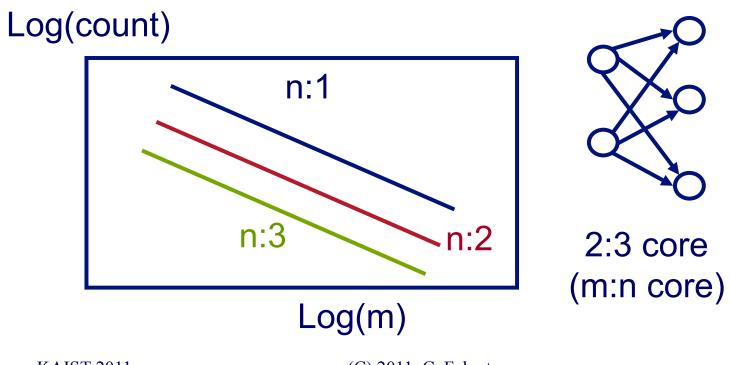
- Bow-tie, for the web [Kumar+ '99]
- IN, SCC, OUT, 'tendrils'

disconnected components





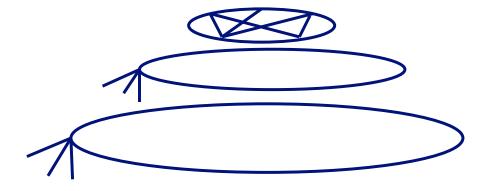
• power-laws in communities (bi-partite cores) [Kumar+, '99]



KAIST-2011 (C) 2011, C. Faloutsos 49



- "Jellyfish" for Internet [Tauro+ '01]
- core: ~clique
- ~5 concentric layers
- many 1-degree nodes



#### **Outline**

- Introduction Motivation
- Patterns in graphs
  - Patterns in Static graphs
    - Degree, eigenvalues
    - Triangles, cliques
    - Other observations
  - Patterns in Weighted graphs
  - Patterns in Time evolving graphs
- Generators



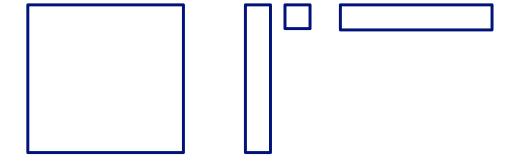




B. Aditya Prakash, Mukund Seshadri, Ashwin Sridharan, Sridhar Machiraju and Christos Faloutsos: *EigenSpokes: Surprising Patterns and Scalable Community Chipping in Large Graphs*, PAKDD 2010, Hyderabad, India, 21-24 June 2010.

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

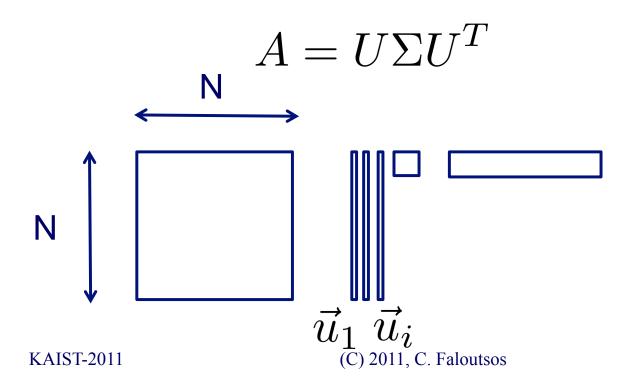
$$A = U\Sigma U^T$$



**KAIST-2011** 

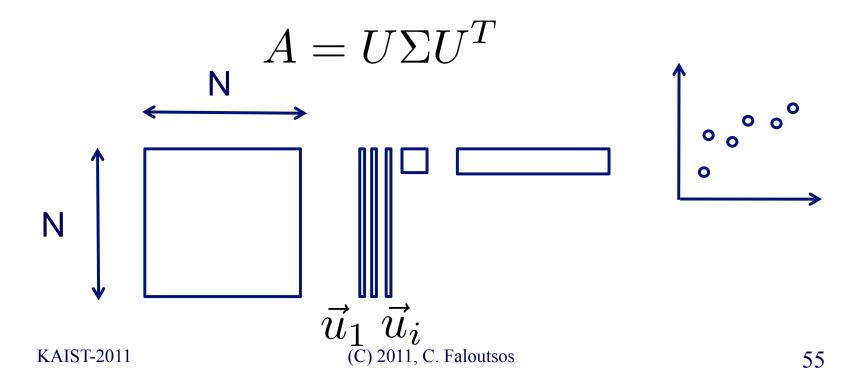


- Eigenvectors of adjacency matrix
  - equivalent to singular vectors (symmetric, undirected graph)



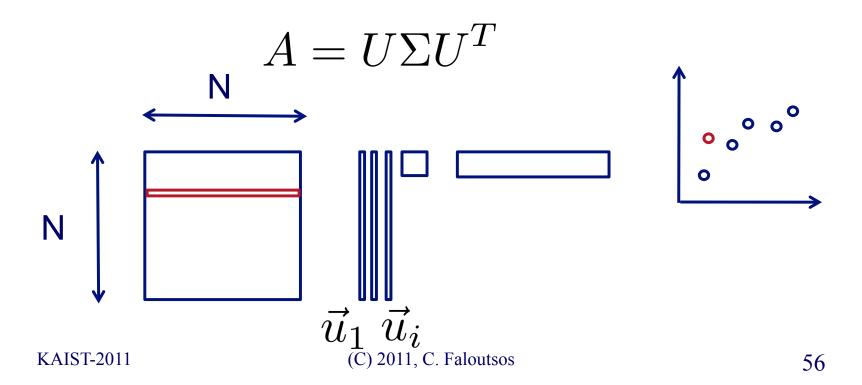


- Eigenvectors of adjacency matrix
  - equivalent to singular vectors (symmetric, undirected graph)



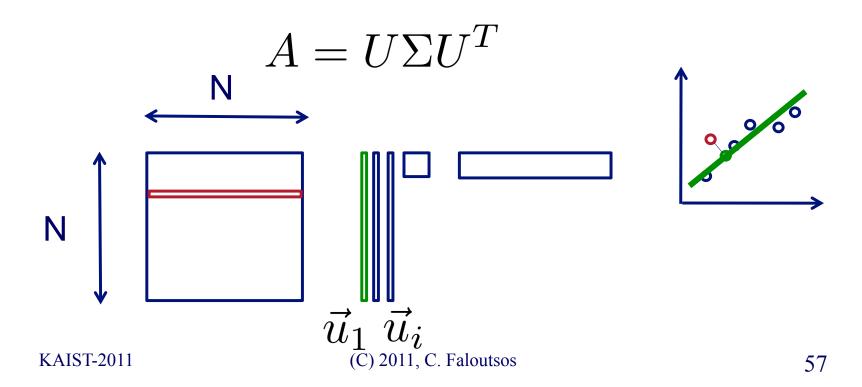


- Eigenvectors of adjacency matrix
  - equivalent to singular vectors (symmetric, undirected graph)





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



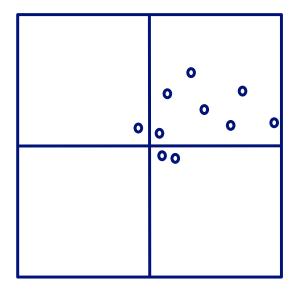
#### Carnegie Mellon

## **EigenSpokes**

• EE plot:

2<sup>nd</sup> Principal component u2

- Scatter plot of scores of u1 vs u2
- One would expect
  - Many points @origin
  - A few scattered~randomly



u1 1<sup>st</sup> Principal

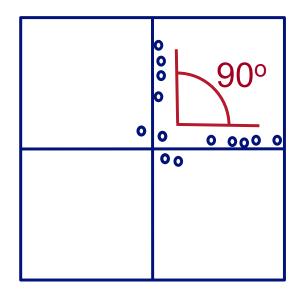
component



**u**2

- EE plot:
- Scatter plot of scores of u1 vs u2
- One would expect
  - Many points @origin

- A few tered



**u**1

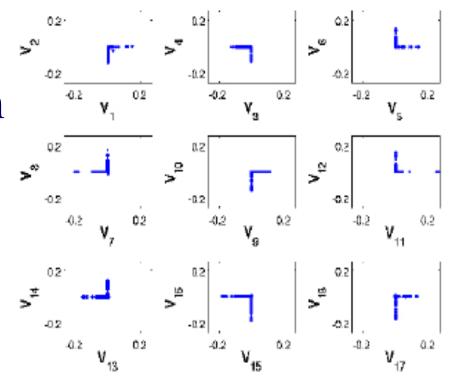
KAIST-2011 (C) 2011, C. Faloutsos 59



## EigenSpokes - pervasiveness

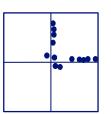
- Present in mobile social graph
  - across time and space

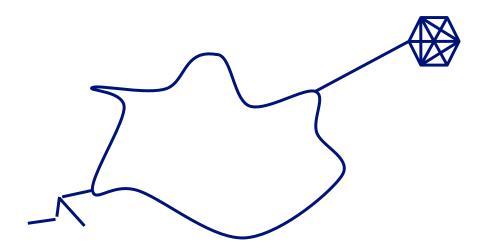
Patent citation graph





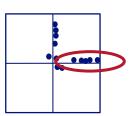
Near-cliques, or nearbipartite-cores, loosely connected

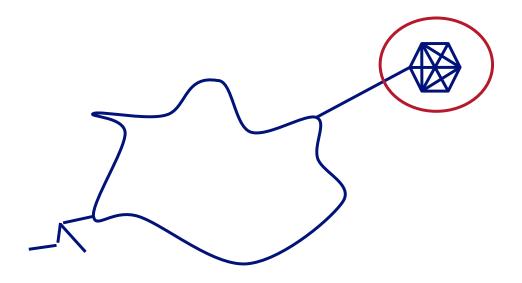






Near-cliques, or nearbipartite-cores, loosely connected



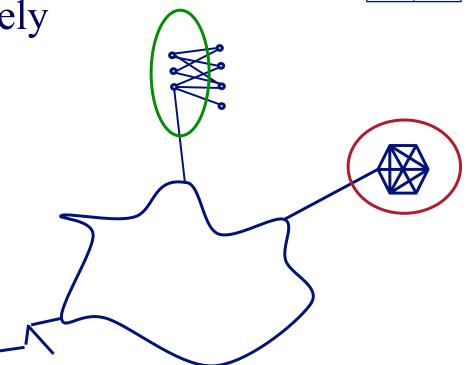




Near-cliques, or near-

bipartite-cores, loosely

connected

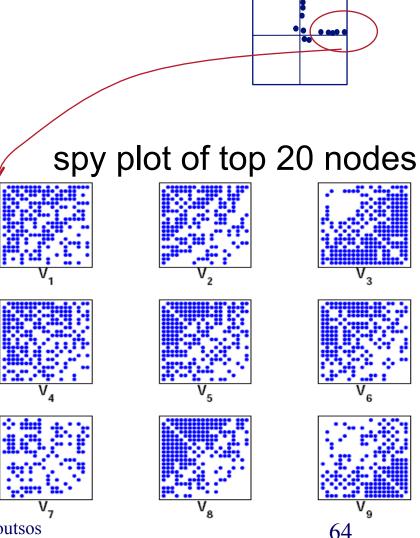




Near-cliques, or nearbipartite-cores, loosely connected

#### So what?

- Extract nodes with high scores
- high connectivity
- Good "communities"



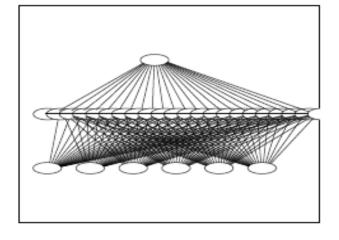


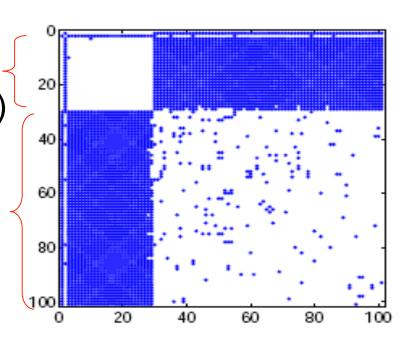
## **Bipartite Communities!**

patents from same inventor(s)

cut-and-paste bibliography!

magnified bipartite community





KAIST-2011 (C) 2011, C. Faloutsos 65

#### **Outline**

- Introduction Motivation
- Patterns in graphs
  - Patterns in Static graphs



- Patterns in Weighted graphs
- Patterns in Time evolving graphs
- Generators



# 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. SIG-KDD 2008

**KAIST-2011** 

#### **Observation W.1: Fortification**

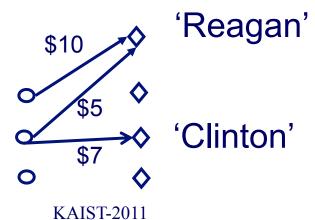
Q: How do the weights of nodes relate to degree?

KAIST-2011 (C) 2011, C. Faloutsos 68



#### **Observation W.1: Fortification**

## More donors, more \$ ?

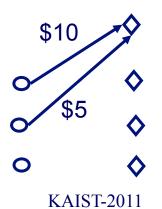




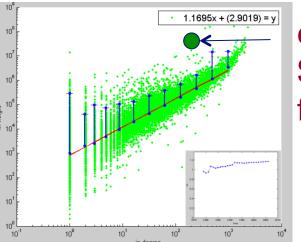
# Observation W.1: fortification: Snapshot Power Law

- 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

(C) 2011, C. Faloutsos

70

**Orgs-Candidates** 

#### **Outline**

- Introduction Motivation
- Patterns in graphs
  - Patterns in Static graphs
  - Patterns in Weighted graphs



- Patterns in Time evolving graphs
- Generators



#### **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 hints 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

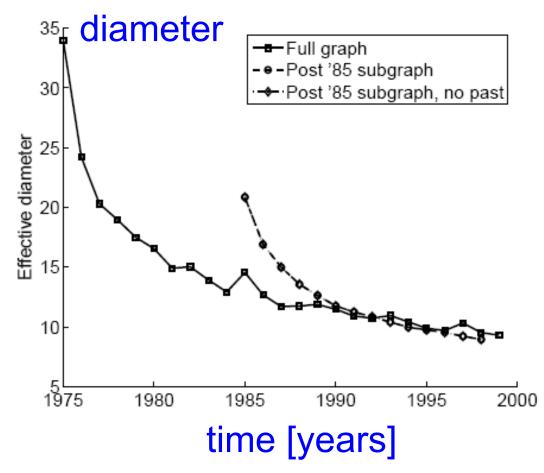
- Prior work on Power Law graphs hints at slowly growing diameter:
  - diameter ~ ((leg ))
  - diameter ~ O(105 log N)



• 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

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

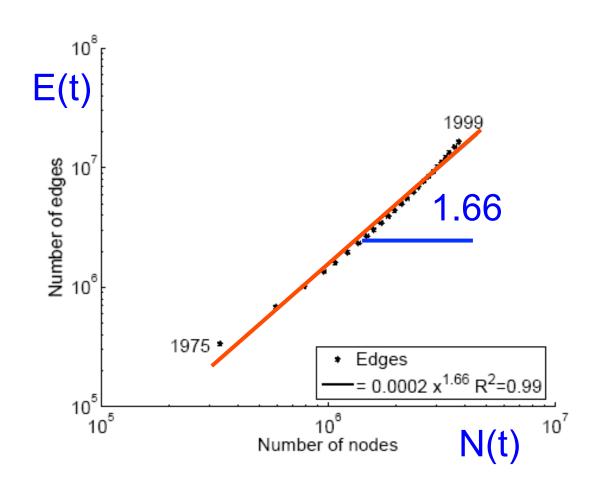
- Q: what is your guess for E(t+1) = (t+1) \* E(t)
- A: over-doubled!
  - But obeying the ``Densification Power Law''

KAIST-2011



# T.2 Densification – Patent Citations

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



KAIST-2011

(C) 2011, C. Faloutsos

#### **Outline**

- Introduction Motivation
- Patterns in graphs
  - Patterns in Static graphs
  - Patterns in Weighted graphs



- Patterns in Time evolving graphs
- Generators



## More on Time-evolving graphs

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

Q: How do NLCC's emerge and join with the GCC?

(``NLCC'' = non-largest conn. components)

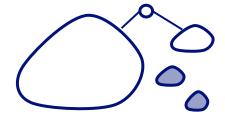
- −Do they continue to grow in size?
- or do they shrink?
- or stabilize?



Q: How do NLCC's emerge and join with the GCC?

(``NLCC'' = non-largest conn. components)

- −Do they continue to grow in size?
- or do they <u>shrink</u>?
- or stabilize?



Q: How do NLCC's emerge and join with the GCC?

(``NLCC'' = non-largest conn. components)

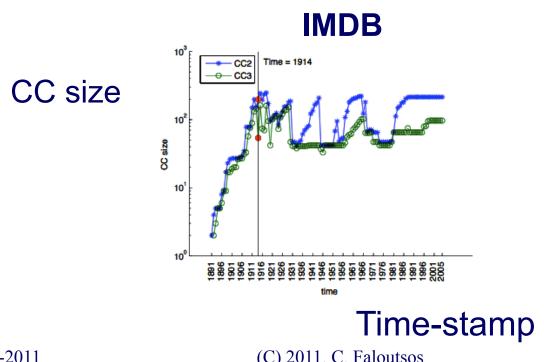
YES – Do they continue to grow in size?

YES – or do they shrink?

YES – or stabilize?



• After the gelling point, the GCC takes off, but NLCC's remain ~constant (actually, oscillate).





## **Timing for Blogs**

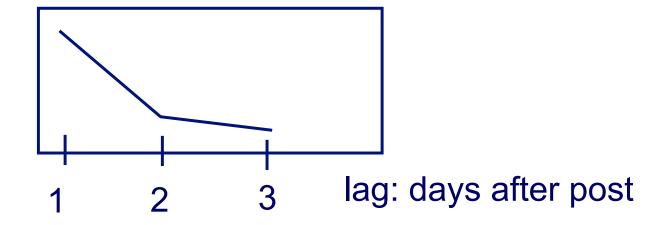
- with Mary McGlohon (CMU)
- Jure Leskovec (CMU->Stanford)
- Natalie Glance (now at Google)
- Mat Hurst (now at MSR)

[SDM'07]

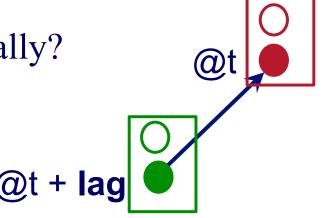


# T.4: popularity over time

# in links



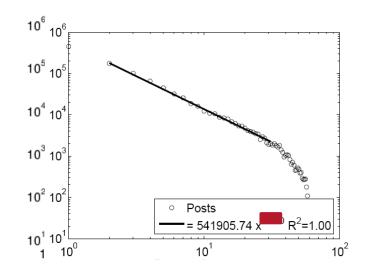
Post popularity drops-off – exponentially?





## T.4: popularity over time

# in links (log)



days after post (log)

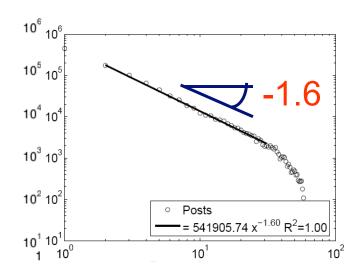
Post popularity drops-off – exporentally? POWER LAW!

Exponent?



# T.4: popularity over time

# in links (log)

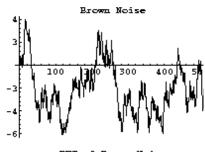


days after post (log)

Post popularity drops-off – exportentially? POWER LAW!

Exponent? -1.6

- close to -1.5: Barabasi's stack model
- and like the zero-crossings of a random walk KAIST-2011 (C) 2011, C. Faloutsos



DFT of Brown Noise

## -1.5 slope

J. G. Oliveira & A.-L. Barabási Human Dynamics: The Correspondence Patterns of Darwin and Einstein.

Nature 437, 1251 (2005). [PDF]

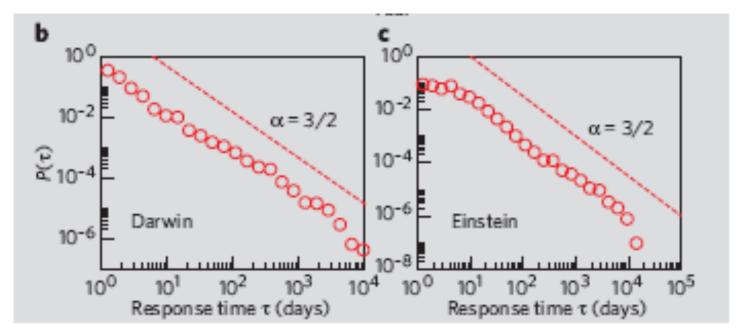


Figure 1 | The correspondence patterns of Darwin and Einstein.

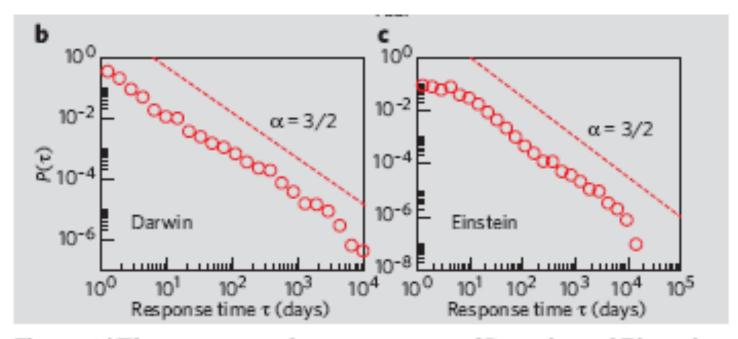


Figure 1 | The correspondence patterns of Darwin and Einstein.



## T.5: duration of phonecalls

Surprising Patterns for the Call Duration Distribution of Mobile Phone Users



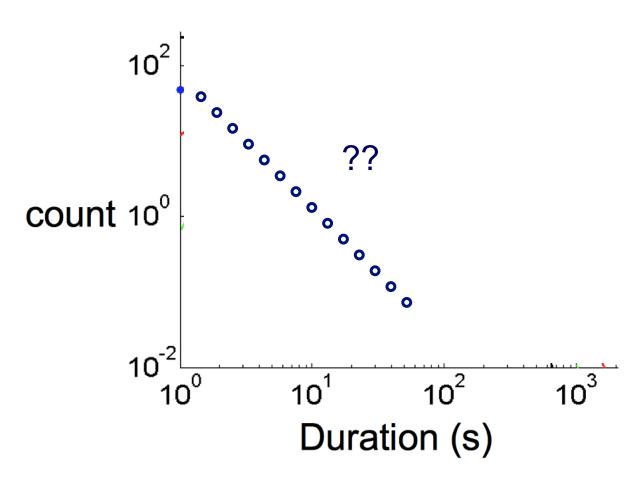
Pedro O. S. Vaz de Melo, Leman

Akoglu, Christos Faloutsos, Antonio

A. F. Loureiro

PKDD 2010

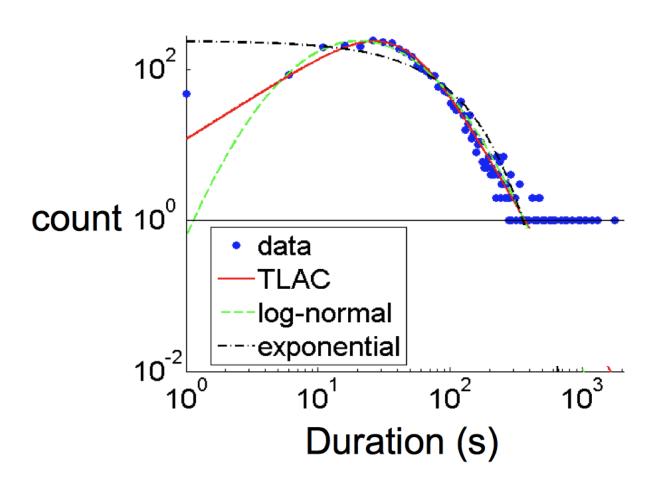
# Probably, power law (?)



**KAIST-2011** 



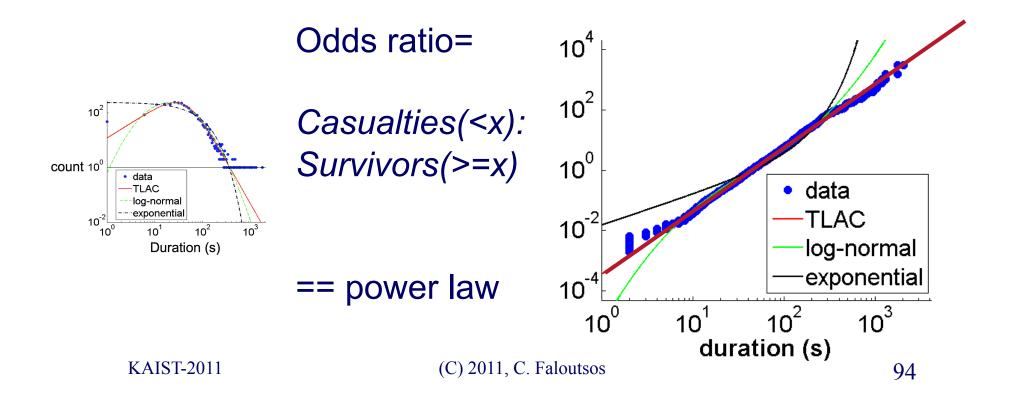
## No Power Law!





## 'TLaC: Lazy Contractor'

- The longer a task (phonecall) has taken,
- The even longer it will take



## **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)

#### **Outline**

- Introduction Motivation
- Patterns in graphs
- Generators



- Erdos-Renyi
- Degree based
- Process based
- Kronecker

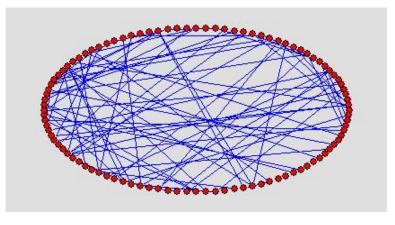
#### Generators

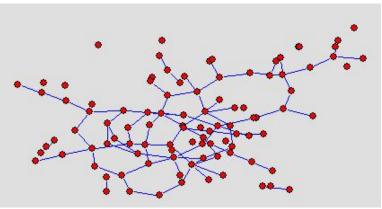
- How to generate random, realistic graphs?
  - Erdos-Renyi model: beautiful, but unrealistic
  - degree-based generators
  - process-based generators
  - recursive/self-similar generators



## **Erdos-Renyi**

- random graph 100
   nodes, avg degree = 2
- Fascinating properties (phase transition)
- But: unrealistic
   (Poisson degree distribution != power law)



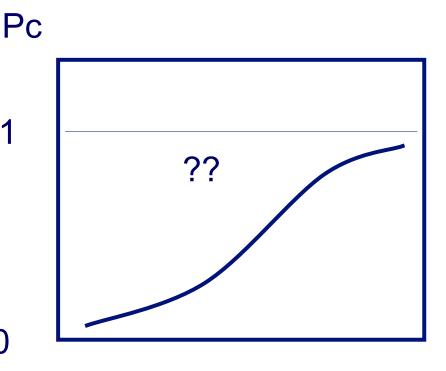






## E-R model & Phase transition

- vary avg degree D
- watch Pc =
   Prob( there is a giant connected component)
- How do you expect it to be?



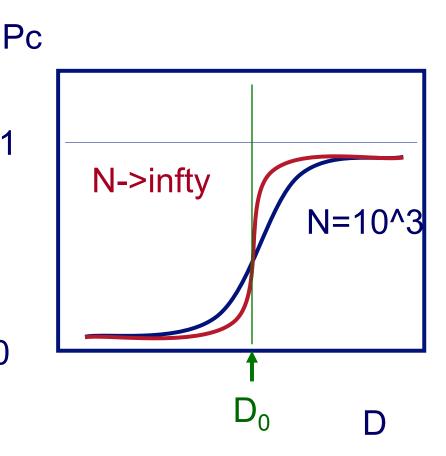
D





## E-R model & Phase transition

- vary avg degree D
- watch Pc =
   Prob( there is a giant connected component)
- How do you expect it to be?



## Degree-based

- Figure out the degree distribution (eg., 'Zipf')
- Assign degrees to nodes
- Put edges, so that they match the original degree distribution



### **Process-based**

- Barabasi; Barabasi-Albert: Preferential attachment -> power-law tails!
  - 'rich get richer'
- [Kumar+]: preferential attachment + mimick
  - Create 'communities'

## Process-based (cont'd)

- [Fabrikant+, '02]: H.O.T.: connect to closest, high connectivity neighbor
- [Pennock+, '02]: Winner does NOT take all

#### **Outline**

- Introduction Motivation
- Patterns in graphs
- Generators
  - Erdos-Renyi
  - Degree based
  - Process based



- Kronecker



## Recursive generators

- (RMAT [Chakrabarti+,'04])
- Kronecker product



## Wish list for a generator:

- Power-law-tail in- and out-degrees
- Power-law-tail scree plots
- shrinking/constant diameter
- Densification Power Law
- communities-within-communities



## Wish list for a generator:

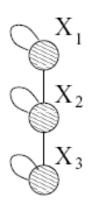
- Power-law-tail in- and out-degrees
- Power-law-tail scree plots
- shrinking/constant diameter
- Densification Power Law
- communities-within-communities

Q: how to achieve all of them?

A: Self-similarity - Kronecker matrix product [Leskovec+05b]
(C) 2011, C. Faloutsos



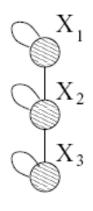
# Kronecker product

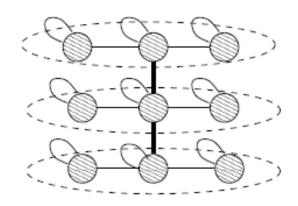


(a) Graph  $G_1$ 



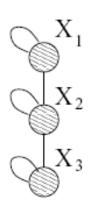
# Kronecker product

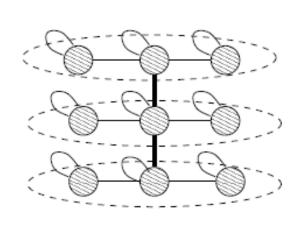


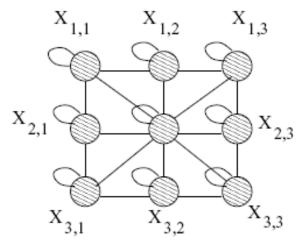


- (a) Graph  $G_1$  (b) Intermediate stage

## Kronecker product



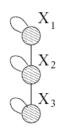


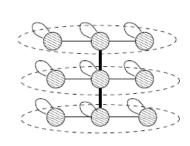


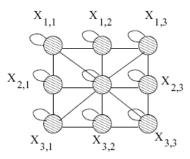
Central node is  $X_{2,2}$ 

- (a) Graph  $G_1$  (b) Intermediate stage (c) Graph  $G_2 = G_1 \otimes G_1$

# Kronecker product







Central node is X 2.2

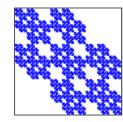
- (a) Graph  $G_1$
- 1
   1
   0

   1
   1
   1

   0
   1
   1

(b) Intermediate stage (c) Graph  $G_2 = G_1 \otimes G_1$ 

$G_1$	$G_1$	0
$G_1$	$G_1$	$G_1$
0	$G_1$	$G_1$



(d) Adjacency matrix of  $G_1$ 



(e) Adjacency matrix (f) Plot of  $G_4$  of  $G_2 = G_1 \otimes G_1$ 



N\*\*4

N\*N

**KAIST-2011** 

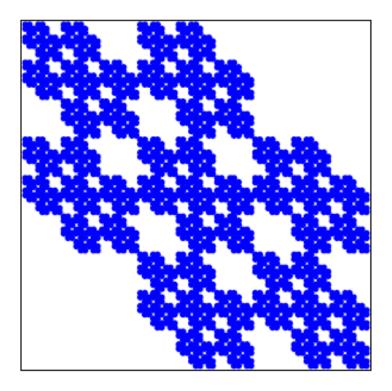
(C) 2011, C. Faloutsos



# Kronecker Product – a Graph

• Continuing multiplying with  $G_1$  we obtain  $G_4$  and

so on ...



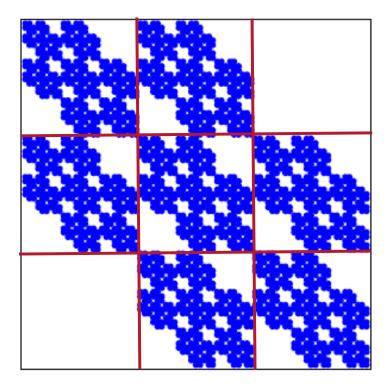
G<sub>4</sub> adjacency matrix (C) 2011, C. Faloutsos



# Kronecker Product – a Graph

• Continuing multiplying with  $G_1$  we obtain  $G_4$  and

so on ...



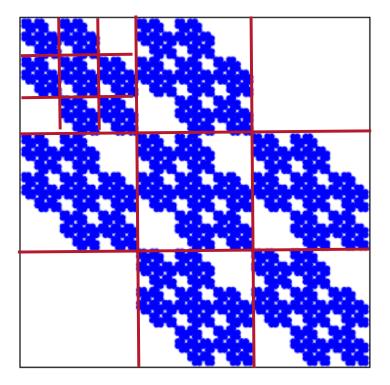
G<sub>4</sub> adjacency matrix (C) 2011, C. Faloutsos



# Kronecker Product – a Graph

• Continuing multiplying with  $G_1$  we obtain  $G_4$  and

so on ...



G<sub>4</sub> adjacency matrix (C) 2011, C. Faloutsos



# Properties of Kronecker graphs:

- Y Power-law-tail in- and out-degrees
- ✓ Power-law-tail scree plots
- ✓ constant diameter
- perfect Densification Power Law
- communities-within-communities



# Properties of Kronecker graphs:

- Power-law-tail in- and out-degrees
- Power-law-tail scree plots
- **constant** diameter
- perfect Densification Power Law
- communities-within-communities

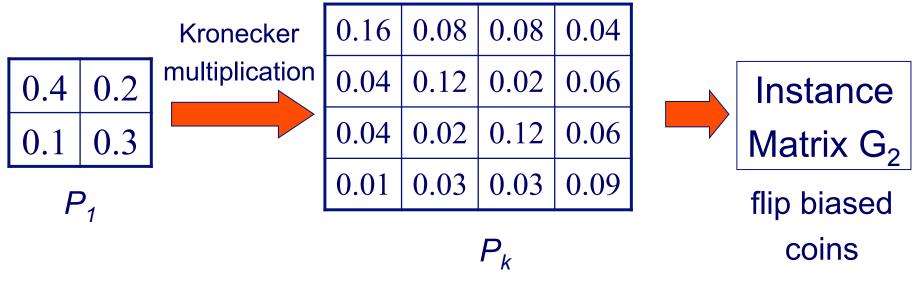
and we can prove all of the above

(first generator that does that)



# Stochastic Kronecker Graphs

- Create  $N_1 \times N_1$  probability matrix  $P_1$
- Compute the  $k^{th}$  Kronecker power  $P_k$
- For each entry  $p_{uv}$  of  $P_k$  include an edge (u,v) with probability  $p_{uv}$



# **Experiments**

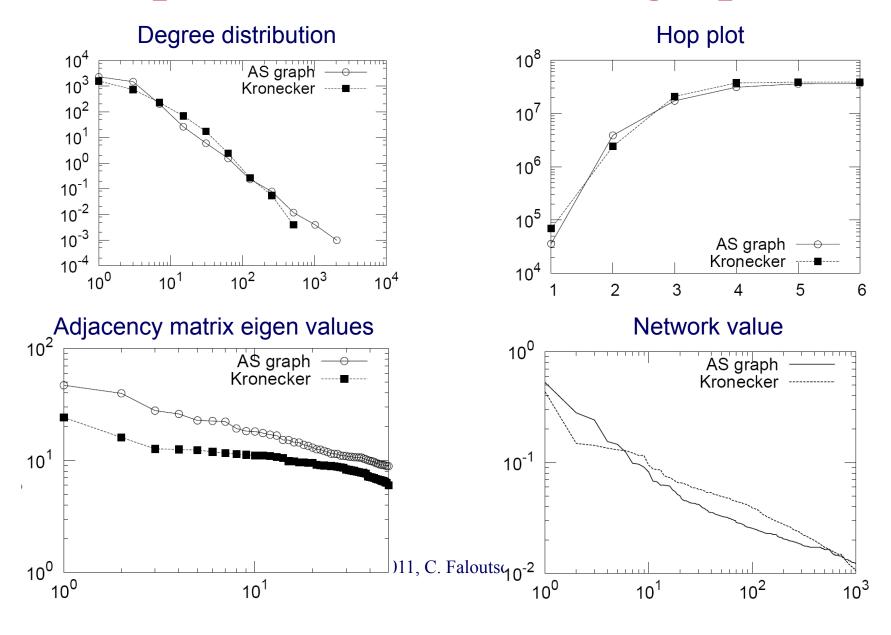
- How well can we match real graphs?
  - Arxiv: physics citations:
    - 30,000 papers, 350,000 citations
    - 10 years of data
  - U.S. Patent citation network
    - 4 million patents, 16 million citations
    - 37 years of data
  - Autonomous systems graph of internet
    - Single snapshot from January 2002
    - 6,400 nodes, 26,000 edges
- We show both static and temporal patterns

# (Q: how to fit the parm's?)

#### A:

- Stochastic version of Kronecker graphs +
- Max likelihood +
- Metropolis sampling
- [Leskovec+, ICML'07]

# Experiments on real AS graph



#### **Conclusions**

- Kronecker graphs have:
  - All the static properties
    - ✓ Heavy tailed degree distributions
    - ✓ Small diameter
    - ✓ Multinomial eigenvalues and eigenvectors
  - All the temporal properties
    - ✓ Densification Power Law
    - ✓ Shrinking/Stabilizing Diameters
  - We can formally prove these results



#### **OVERALL CONCLUSIONS**

- Several new **patterns** (fortification, triangle-laws, conn. components, etc)
- Recursive generators (Kronecker), with provable properties

• 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)

- 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



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



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

- 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

- 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



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

- 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



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



# **Project info**

www.cs.cmu.edu/~pegasus



Chau, Polo



McGlohon, Mary



Tsourakakis, **Babis** 









Akoglu, Leman

Kang, U

Prakash, **Aditya** 

Tong, Hanghang

Thanks to: NSF IIS-0705359, IIS-0534205,

CTA-INARC; Yahoo (M45), LLNL, IBM, SPRINT,

INTEL, HP

# 



# Extra material – why so many power laws?

### At least 6-7 mechanisms (!)

Power laws, Pareto distributions and Zipf's law Contemporary Physics 46, 323-351 (2005)

#### **Outline**

• Generative mechanisms

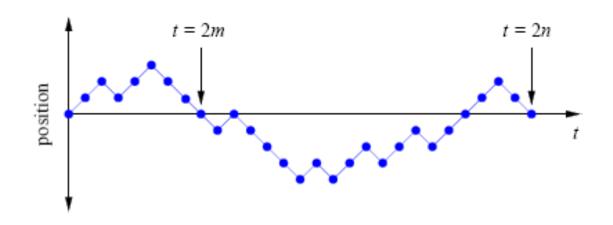


- Random walk
  - Yule distribution = CRP
  - Percolation
  - Self-organized criticality
  - Other

136



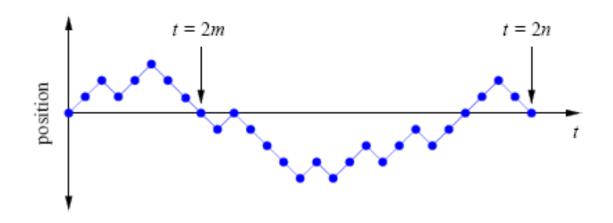
#### Random walks



Inter-arrival times PDF:  $p(t) \sim ??^{2}$ 



#### Random walks



Inter-arrival times PDF:  $p(t) \sim t^{-3/2}$ 



#### Random walks

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

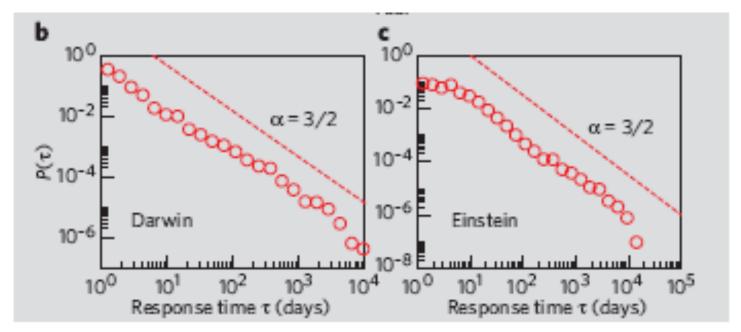


Figure 1 | The correspondence patterns of Darwin and Einstein.

#### **Outline**

- Generative mechanisms
  - Random walk
- Yule distribution = CRP
  - Percolation
  - Self-organized criticality
  - Other



#### Yule distribution and CRP

Chinese Restaurant Process (CRP):

Newcomer to a restaurant

- Joins an existing table (preferring large groups
- Or starts a new table/group of its own, with prob 1/m

a.k.a.: rich get richer; Yule process

#### Yule distribution and CRP

#### Then:

```
Prob( k people in a group) = p_k
= (1 + 1/m) B(k, 2+1/m)
\sim k^{-(2+1/m)}
(since B(a,b) \sim a ** (-b) : power law tail)
```



#### Yule distribution and CRP

- Yule process
- Gibrat principle
- Matthew effect
- Cumulative advantage
- Preferential attachement
- 'rich get richer'

#### **Outline**

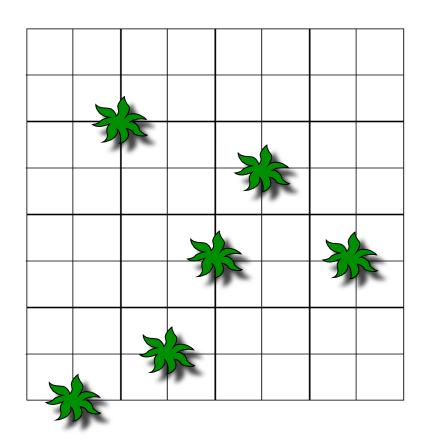
- Generative mechanisms
  - Random walk
  - Yule distribution = CRP



- Percolation
  - Self-organized criticality
  - Other



#### Percolation and forest fires

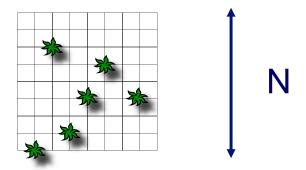


A burning tree will cause its neighbors to burn next.

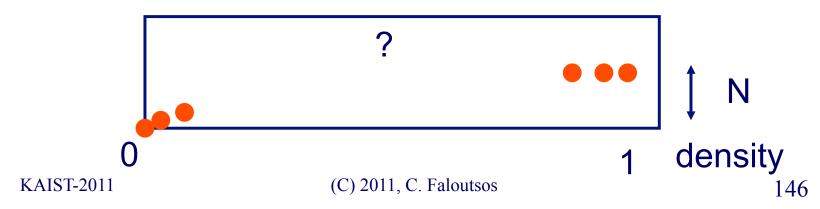
Which tree density *p* will cause the fire to last longest?



### Percolation and forest fires



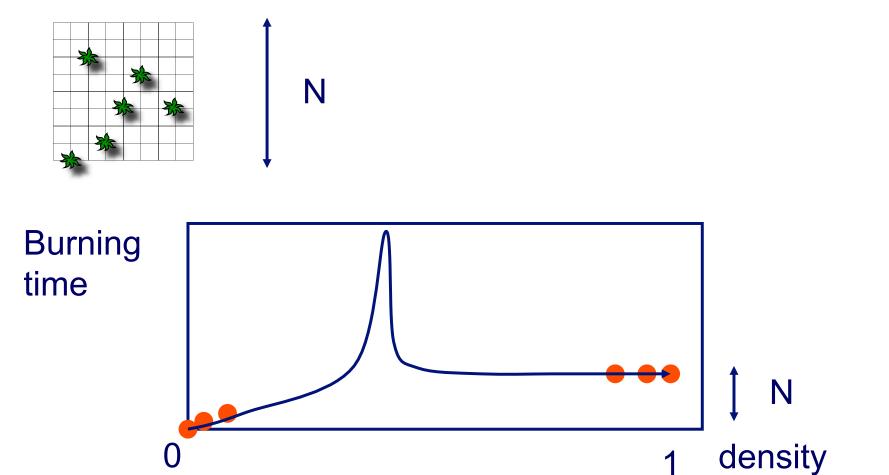
# Burning time





**KAIST-2011** 

### Percolation and forest fires

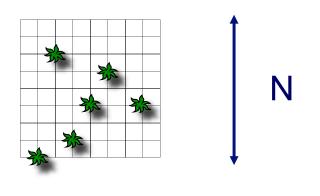


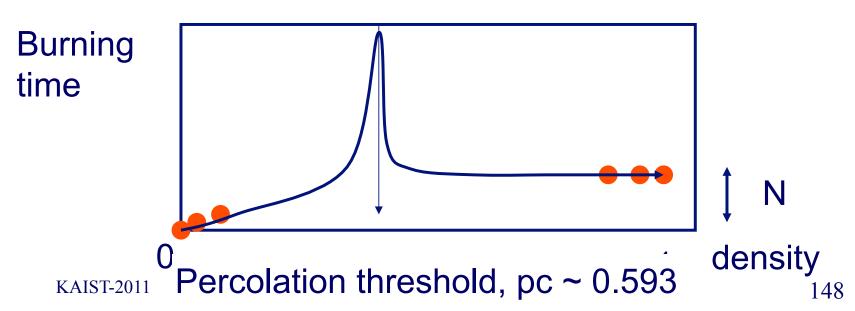
(C) 2011, C. Faloutsos

147



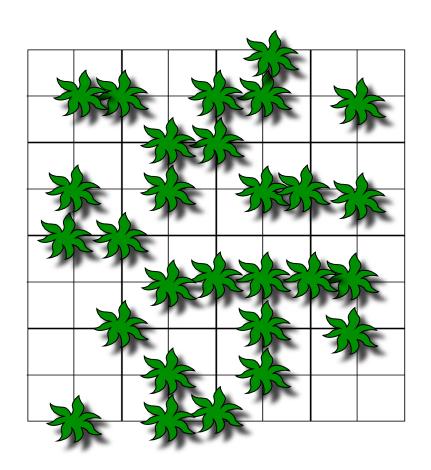
#### Percolation and forest fires







### Percolation and forest fires



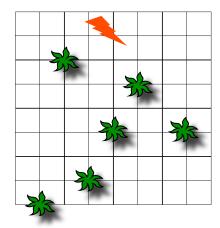
At pc ~ 0.593: No characteristic scale; 'patches' of all sizes; Korcak-like 'law'.

### **Outline**

- Generative mechanisms
  - Random walk
  - Yule distribution = CRP
  - Percolation
- Self-organized criticality
  - Other

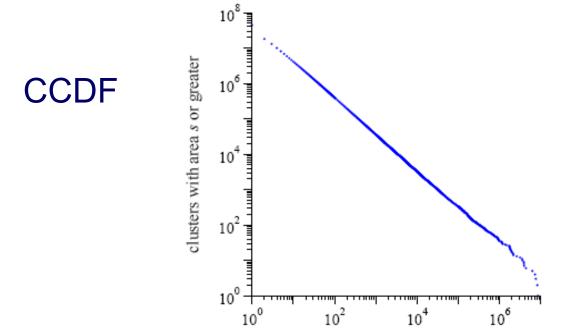


- Trees appear at random (eg., seeds, by the wind)
- Fires start at random (eg., lightning)
- Q1: What is the distribution of size of forest fires?





• A1: Power law-like

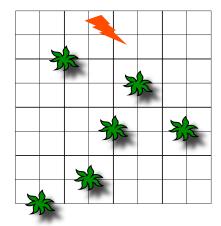


(C) 2011, C. Falout Area of cluster s

area of cluster s



- Trees appear at random (eg., seeds, by the wind)
- Fires start at random (eg., lightning)
- Q2: what is the average density?





• A2: the critical density  $pc \sim 0.593$ 

- [Bak]: size of avalanches ~ power law:
- Drop a grain randomly on a grid
- It causes an avalanche if height(x,y) is >1 higher than its four neighbors

[Per Bak: How Nature works, 1996]

### **Outline**

- Generative mechanisms
  - Random walk
  - Yule distribution = CRP
  - Percolation
  - Self-organized criticality



Other

- Random multiplication
- Fragmentation
- -> lead to lognormals (~ look like power laws)

### Random multiplication:

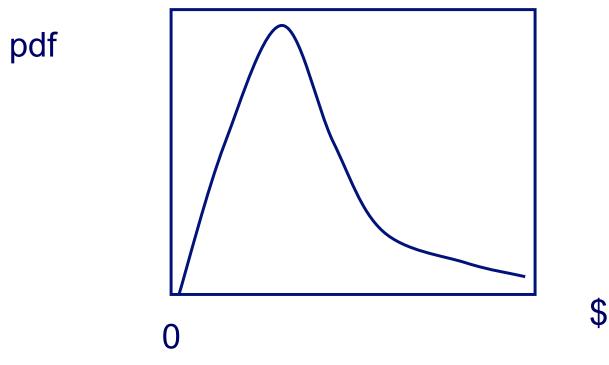
- Start with C dollars; put in bank
- Random interest rate s(t) each year t
- Each year t: C(t) = C(t-1) \* (1+s(t))
- Log(C(t)) = log( C ) + log(..) + log(..) ... -> Gaussian

#### Random multiplication:

• Log(C(t)) = log( C ) + log(..) + log(..) ... -> Gaussian

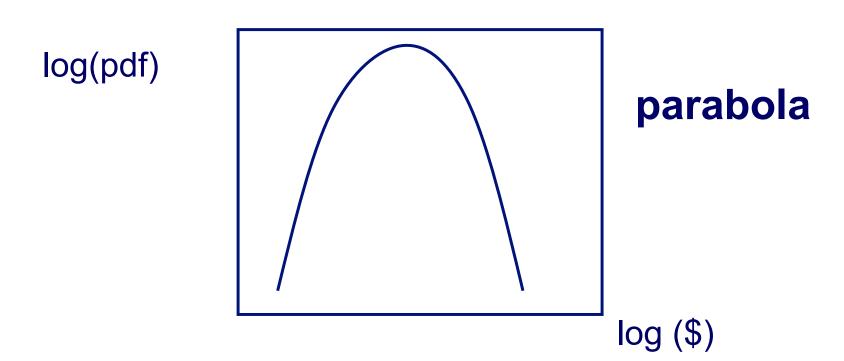
- Thus  $C(t) = \exp(Gaussian)$
- By definition, this is Lognormal

# Lognormal:



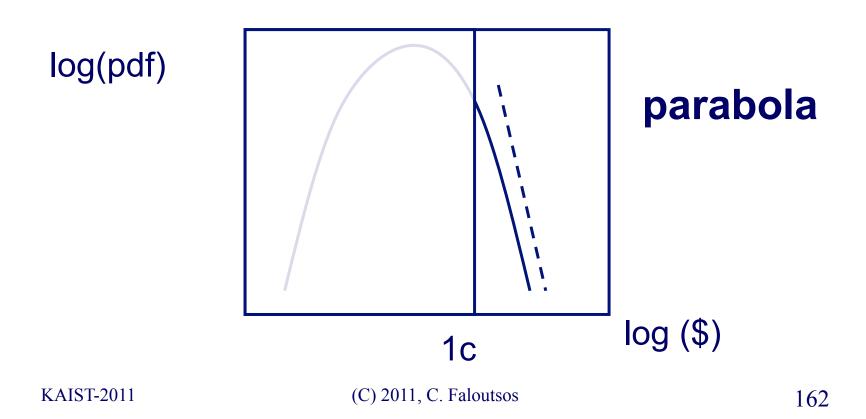


## Lognormal:





## Lognormal:



- Random multiplication
- Fragmentation
- -> lead to lognormals (~ look like power laws)

- Stick of length 1
- Break it at a random point  $x (0 \le x \le 1)$
- Break each of the pieces at random

• Resulting distribution: lognormal (why?)



#### **Conclusions**

• Many, natural mechanisms, may yield power-laws (or log-normals etc)



# **Questions?**

• www.cs.cmu.edu/~christos