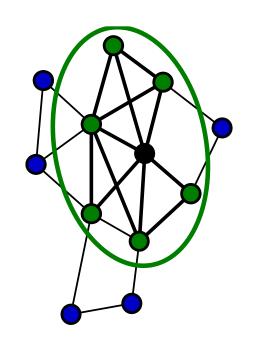
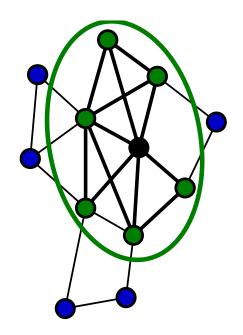
Vertex Neighborhoods, Low Conductance Cuts, and Good Seeds for Local Community Methods (KDD 2012)



DAVID F. GLEICH PURDUE SESH COMANDUR
SANDIA LIVERMORE

Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000."



A vertex neighborhood is a "good" conductance community in a graph with a heavy-tailed degree distribution and large clustering coefficient.

#### Our contributions

1. The previous theorem and its proof. This shows that good communities are *expected* and easy to find in modern networks with heavy-tailed degrees and large clustering.

2. An empirical evaluation of neighborhood communities that shows vertex neighborhoods are the "backbone" of the network community profile.

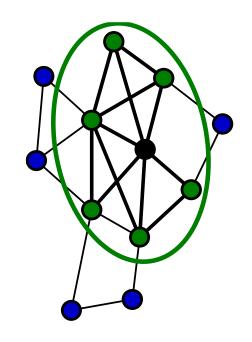
### Formal background for the theorem

- 1. Vertex neighborhoods
- Low conductance cuts
- 3. Clustering coefficients

### Vertex neighborhoods

The set of a vertex and all its neighborhood

Also called an "egonet"



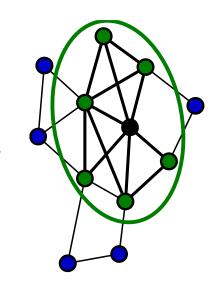
Prior research on egonets of social networks from the "structural holes" perspective [Burt95, Kleinberg08].

Used for anomaly detection [Akoglu10], community seeds [Huang11,Schaeffer11], overlapping communities [Schaeffer07,Rees10].

#### Conductance communities

Conductance is one of the most important community scores [Schaeffer07]

The conductance of a set of vertices is the ratio of edges leaving to total edges:



$$\phi(S) = \frac{\text{cut}(S)}{\min(\text{vol}(S), \text{vol}(\bar{S}))} \frac{\text{(edges leaving the set)}}{\text{cut}(S) = 7} \text{vol}(S) = 33$$

Equivalently, it's the probability that a random edge leaves the set.

$$vol(S) = 33$$

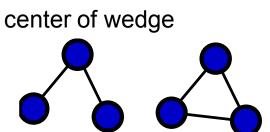
$$vol(\bar{S}) = 11$$

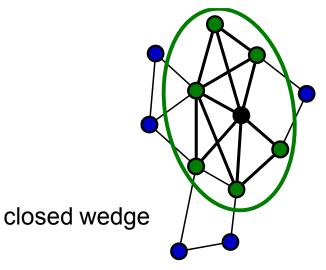
$$\phi(S) = 7/11$$

Small conductance ⇔ Good community

### Clustering coefficients

Wedge





Global clustering coefficient

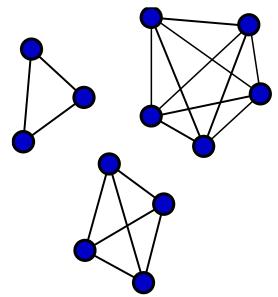
$$\kappa = \frac{\text{number of closed wedges}}{\text{number of wedges}}$$

Probability that a random wedge is closed

#### Simple version of theorem

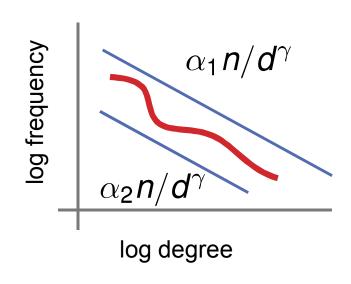
If global clustering coefficient = 1, then the graph is a disjoint union of cliques.

Vertex neighborhoods are optimal communities!



#### **Theorem**

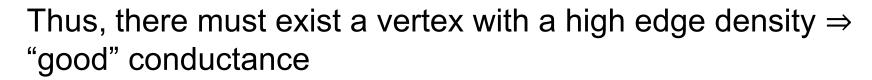
Condition: Let graph G have clustering coefficient  $\kappa$  and have vertex degrees bounded by a power-law function with exponent  $\gamma$  less than 3.



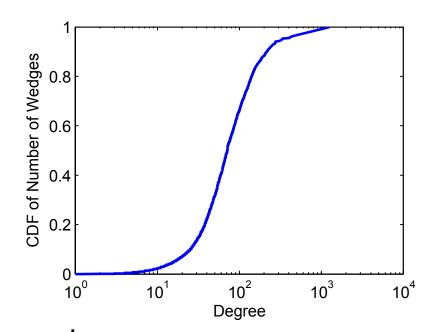
Theorem: Then there exists a vertex neighborhood with conductance  $\leq 4(1 - \kappa)/(3 - 2\kappa)$ 

#### **Proof Sketch**

- 1) Large clustering coefficient
- ⇒ many wedges are closed
- 2) Heavy tailed degree dist
- ⇒ a few vertices have a very large degree
- 3) Large degree  $\Rightarrow$  O( $d^2$ ) wedges  $\Rightarrow$  "most" of wedges



Use the probabilistic method to formalize



# Confession The theory is really weak

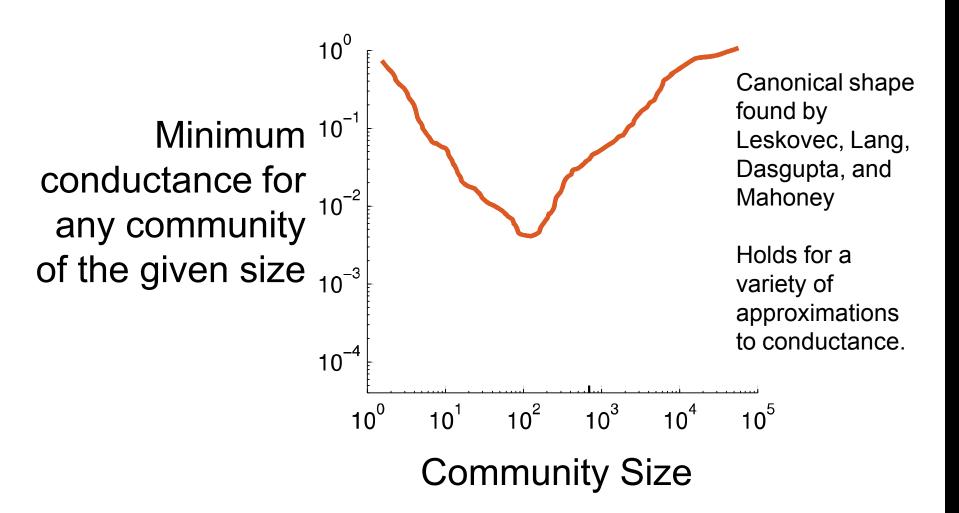
$$\phi(S) \leq 4(1-\kappa)/(3-2\kappa)$$

This bound is useless unless  $\kappa$ 

					> 1 /· )
Graph	Verts	Edges	$\kappa$	Ċ	<del>- 172</del>
ca-AstroPh email-Enron cond-mat-2005 arxiv dblp hollywood-2009	17903 33696 36458 <b>86376</b> 226413 1069126	196972 180811 171735 517563 716460 56306653	0.318 0.085 0.243 0.560 0.383 0.310	0.633 0.509 0.657 0.678 0.635 0.766	Collaboration networks κ ~ [0.1 – 0.5]
fb-Penn94 fb-A-oneyear fb-A soc-LiveJournal1	41536 1138557 3097165 4843953	1362220 4404989 23667394 42845684	0.098 0.038 0.048 0.118	0.212 0.060 0.097 0.274	Social networks $\kappa \sim [0.05 - 0.1]$
oregon2-010526 p2p-Gnutella25 as-22july06 itdk0304	11461 22663 22963 190914	32730 54693 48436 607610	0.037 0.005 0.011 0.061	0.352 0.005 0.230 0.158	Tech. networks $\kappa \sim [0.005 - 0.05]$

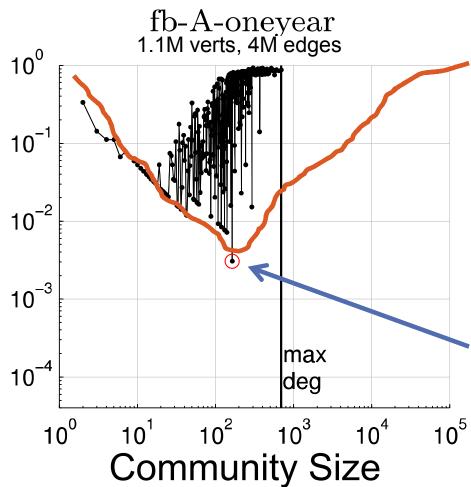
# We view this theory as "intuition for the truth"

## Network Community Profiles



## Network Community Profiles

Minimum 10<sup>-1</sup> conductance for any community neighborhood of 10<sup>-2</sup> the given size



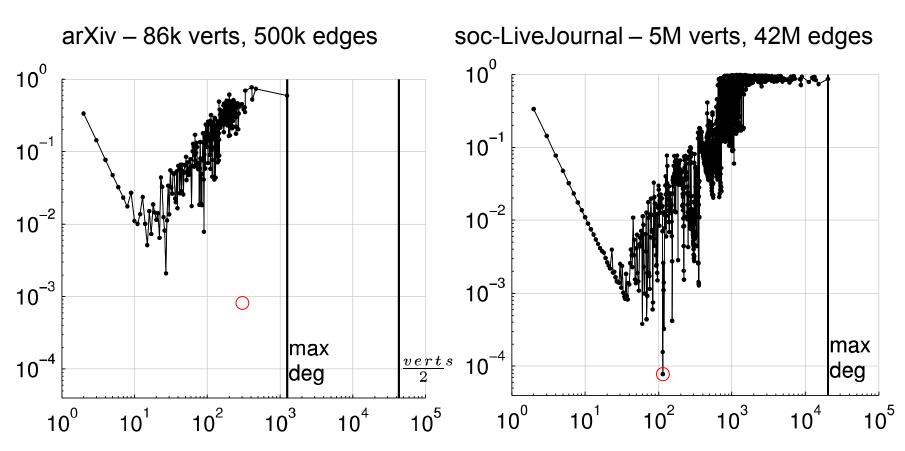
(Degree + 1)

Facebook data from Wilson et al. 2009

"Egonet community profile" shows the same shape, 3 secs to compute.

The Fiedler community computed from the normalized Laplacian is a neighborhood!

### Not just one graph



15 more graphs available www.cs.purdue.edu/~dgleich/codes/neighborhoods

# Communities [Andersen06]

To find the canonical NCP structure, Leskovec et al. used a personalized PageRank based community finder.

These start with a single vertex seed, and then expand the community based on the solution of a personalized PageRank problem.

The resulting community satisfies a local Cheeger inequality.

This needs to run thousands of times for an NCP

# Network Community Profile

Minimum conductance for any community of the given size

This region fills when using the PPR method (like now!)

7807 seconds

5

Community Size

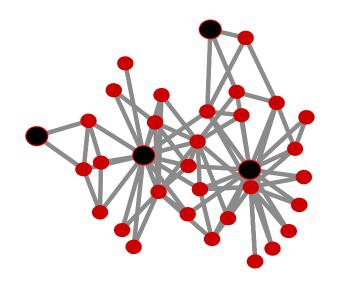
Vertex Neighborhoods, Low Conductance Cuts, and Good Seeds for Local Community Methods

# Locally Minimal Communities

"My conductance is the best locally."

$$\phi(N(v)) \leq \phi(N(w))$$

 $\dot{}$  w adjacent to v

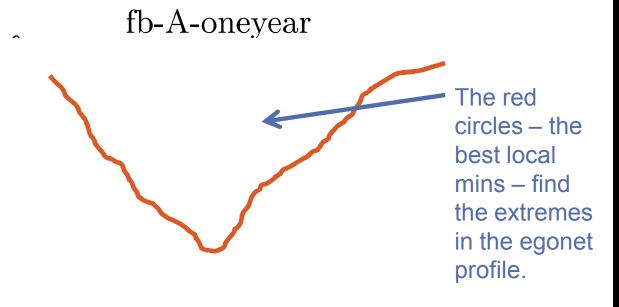


In Zachary's Karate Club network, there are four locally minimal communities

# capture extremal neighborhoods

Red dots are conductance and size of a locally minimal community

Usually about 1% of # of vertices.

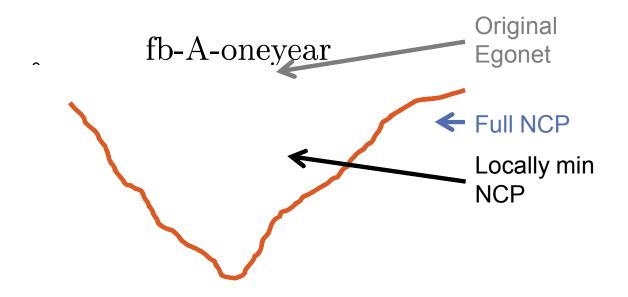


**Community Size** 

5

# Growing locally minimal comm.

Growing only locally minimal communities



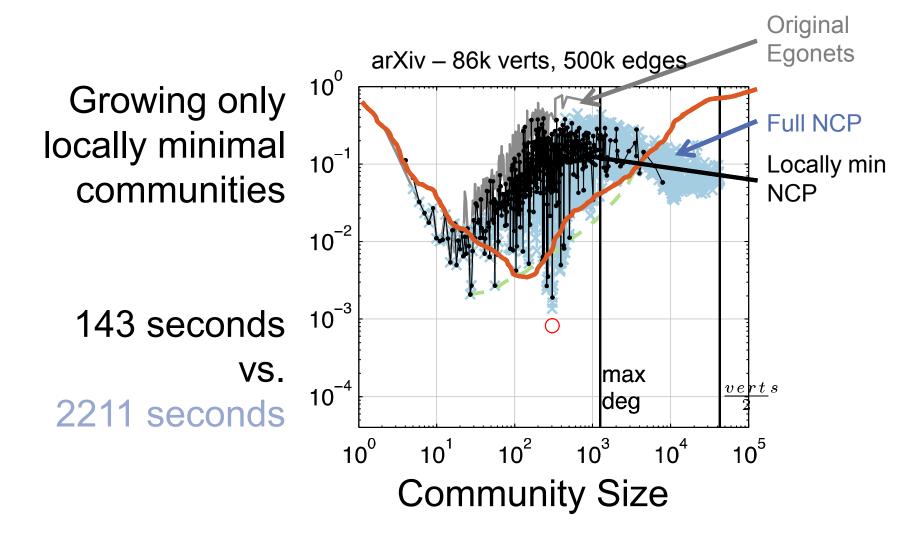
283 seconds vs.

7807 seconds

5

Community Size

## Growing locally minimal comm.



#### Recap

A theorem relating clustering, heavy-tailed degrees, and low-conductance cuts of vertex neighborhoods.

Empirical evaluation of vertex neighborhoods.

More on k-cores in the paper.

- ⇒ Many communities are easy to find!
- ⇒ Explains success of community detection?

Acknowledgements

David supported by NSF CAREER award 1149756-CCF.

Sesh supported by the Sandia LDRD program (project 158477) and the applied mathematics program at the Dept. of Energy.