

# Exploring Twitter Networks in Parallel Computing

Bo Xu<sup>1</sup>, Yun Huang<sup>2</sup>, and Noshir Contractor<sup>2</sup>

<sup>1</sup> Northeastern University, China

<sup>2</sup> Science of Networks in Communities (SONIC) research group  
Northwestern University

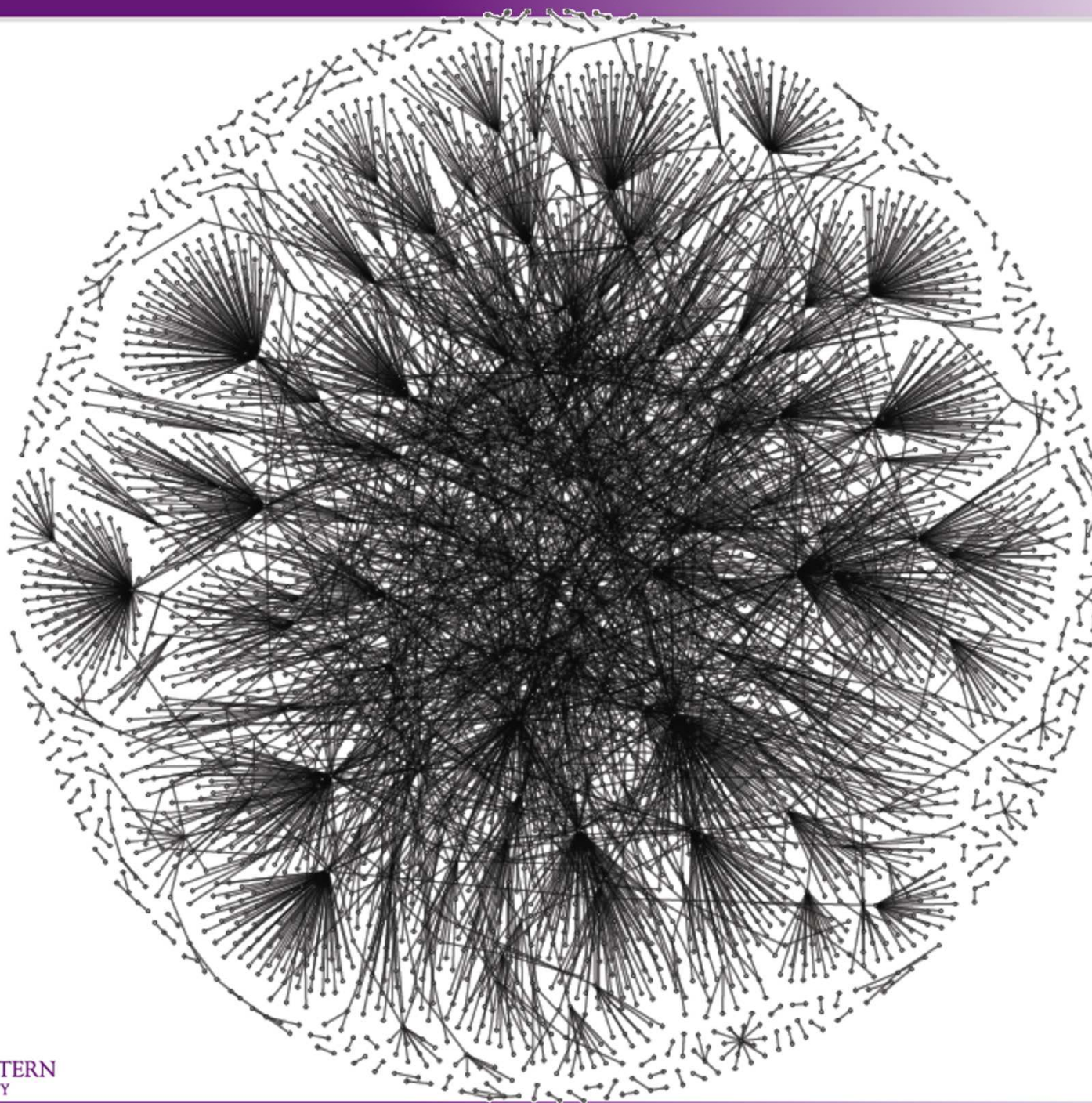
XSEDE

Extreme Science and Engineering  
Discovery Environment



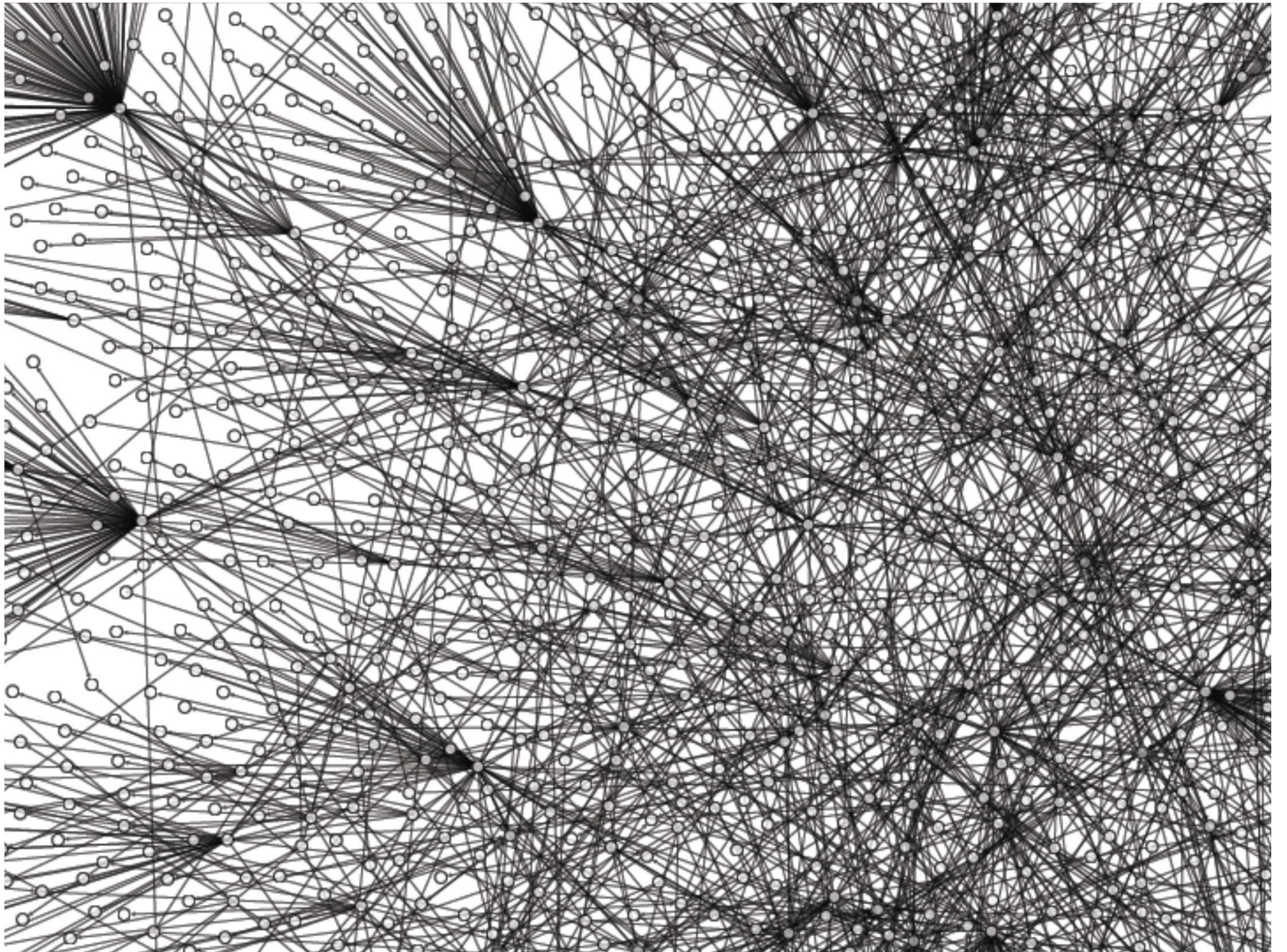
This research was supported by grants from the National Science Foundation (Grant No. CNS-1010904, OCI-0904356, OCI-1053575, & IIS-0841583) and Army Research Lab (W911NF-09-02-0053)





NORTHWESTERN  
UNIVERSITY





# Motivation

- Using parallel computing to explore big networks
  - Visualize network structures of Twitter
  - Quantify the process of link formation
- Media or social networks
  - Heterogeneous nodes and relations of Twitter
- Unfollow behavior
  - Using unfollow to reveal the structure of Twitter networks
  - Identify important factors that influence unfollow behavior using Exponential Random Graph Models



# Previous Research on Unfollow

- Autonomous viewpoint
  - Each link in a network is independent of other nodes and links
- Homogenous analysis
  - Performing analysis on entire population
- Findings
  - Both relational and informational factors are important in maintaining relations in Twitter
  - Results do not give much insights on the network



# Hypotheses and Network Structures

Relational factors	Informational factors
<b>Reciprocity:</b>	<b>Informativeness:</b>
H1. mutual following - →	H6. frequent conversations - →
H2. unfollowed by others + →	
<b>Social status:</b>	
H3.1. more followers + → / - ←	
H3.2. more followees - → / + ←	
<b>Social embeddedness:</b>	<b>Topic homophily:</b>
H4. share common followees - →	H5. share similar topics - →



# Data Description

- 697,628 Korean Twitter users w/ two snapshots
  - June 25<sup>th</sup> and September 3<sup>rd</sup> 2010
- Construct unfollow networks
  - Dissolution of a tie between snapshots as an unfollow relation: a directed link from user A to user B if A followed B but stopped following at the next time point.

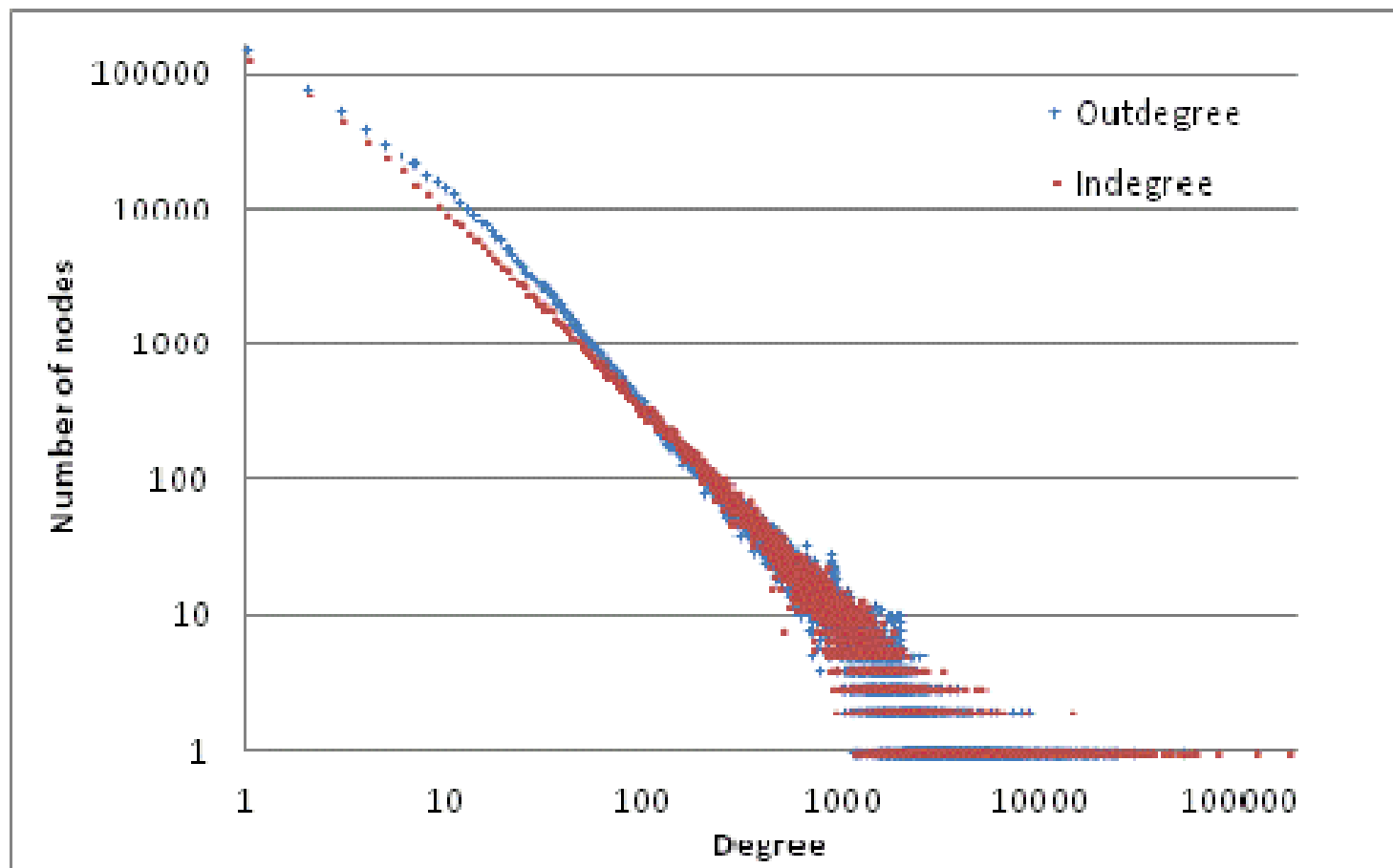


# Data Statistics

	Follow network Time 1	Unfollow network Time 2
<b>Vertices</b>	697,628	211,263 involved
<b>Arcs</b>	34,429,170	858,702 unfollow
<b>Mutual pairs</b>	12,676,988 (73.6% of arcs)	33,015 reciprocal unfollow (7.7% of all unfollow ties)



# Degree Distributions



DB: out-in.txt

Scatter

Var: outdegree,

indegree,

n

7.605e+004

5.704e+004

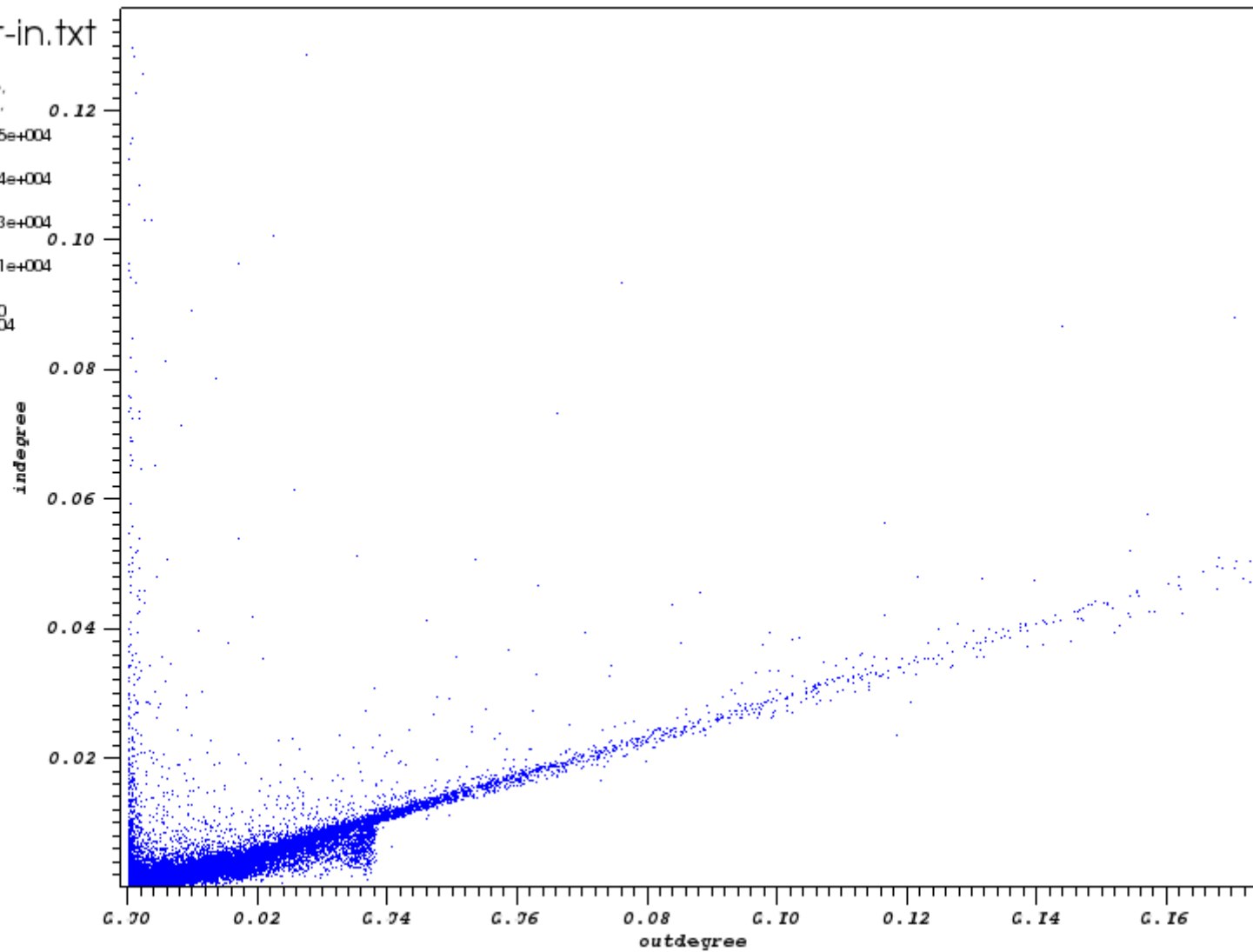
3.803e+004

1.901e+004

1.000

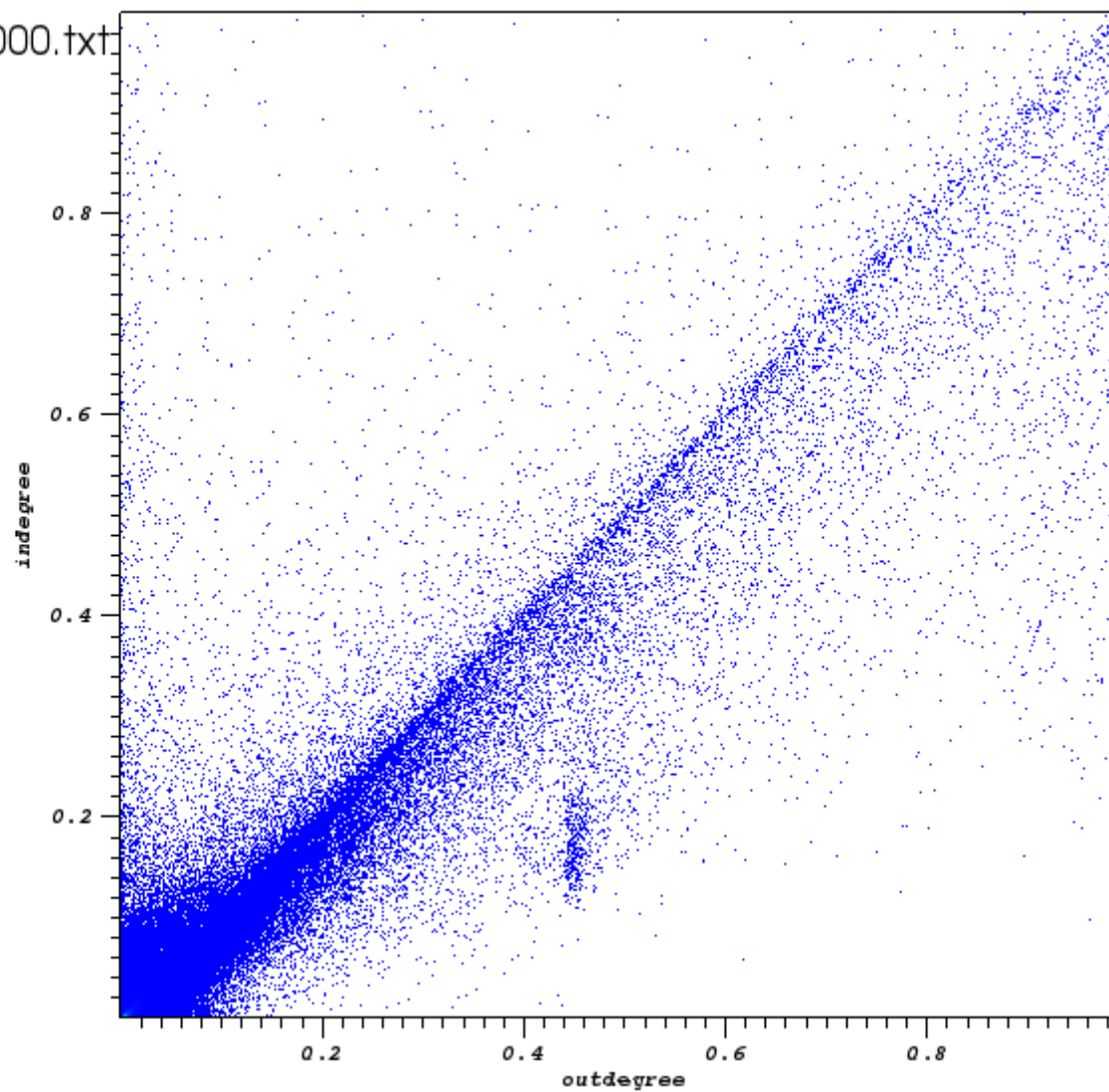
Max: 7.605e+004

Min: 1.000



DB: out-in-2000.txt

Scatter  
Var: outdegree,  
indegree,  
n  
4909.  
3682.  
2455.  
1228.  
1.000  
Max: 4909.  
Min: 1.000



NORTHWESTERN  
UNIVERSITY

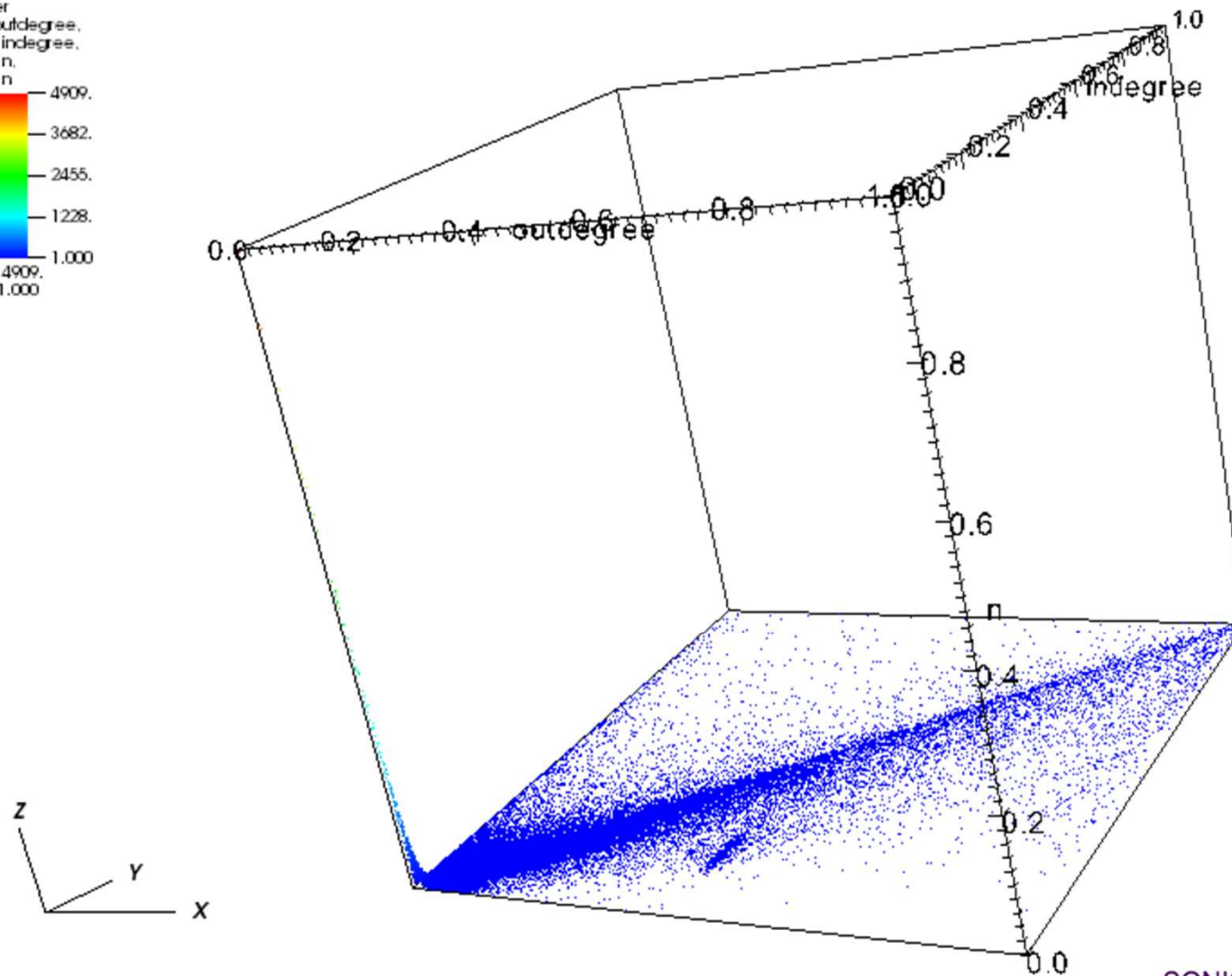
SONIC



advancing the  
science of networks in communities

DB: out-in-2000.txt

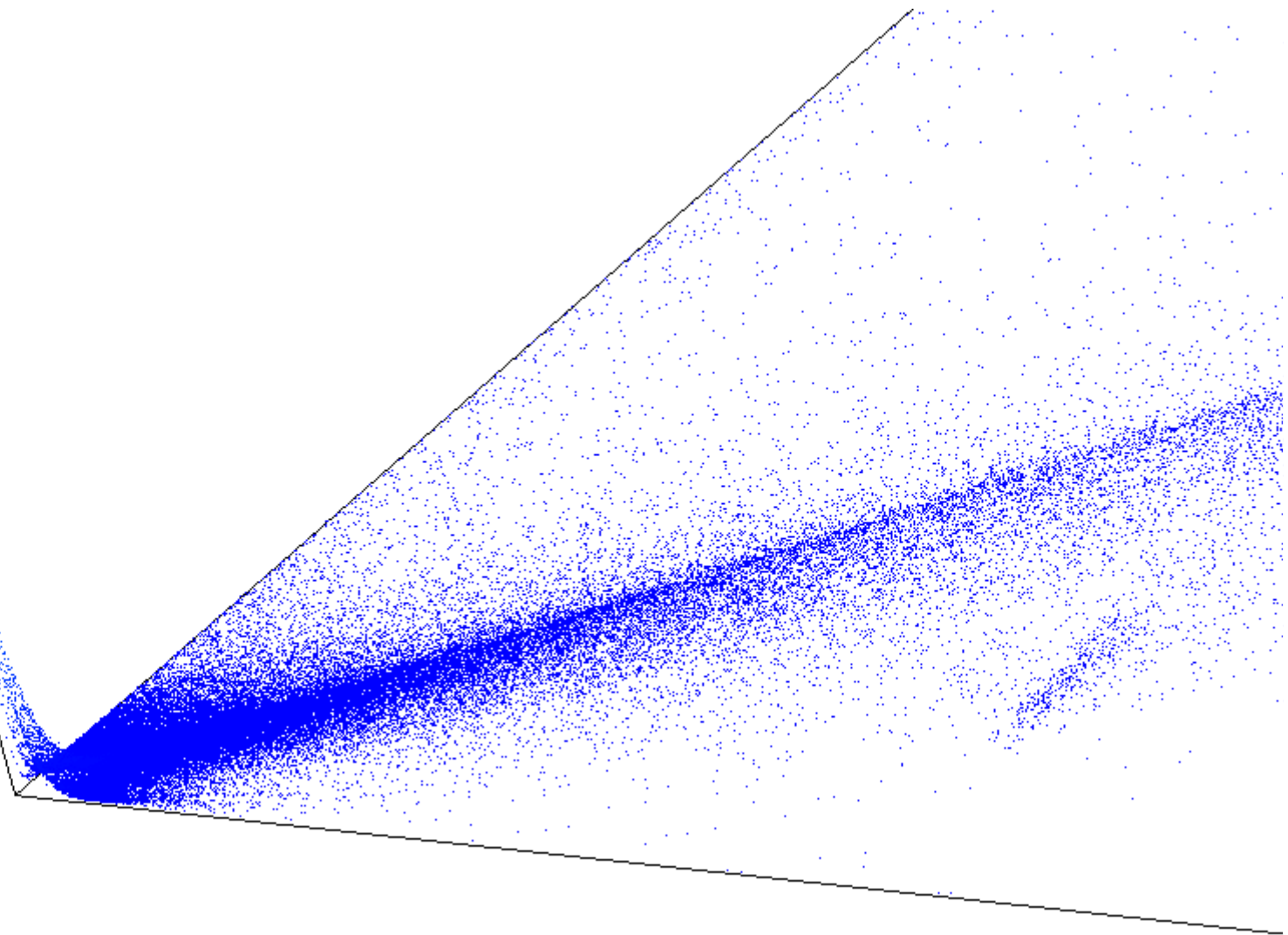
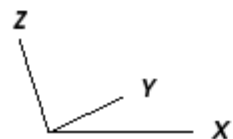
Scatter  
Var: outdegree,  
indegree,  
n,  
n  
4909.  
3682.  
2455.  
1228.  
1.000  
Max: 4909.  
Min: 1.000



NORTHWESTERN  
UNIVERSITY

DB: out-in-2000.txt

Scatter  
Var: outdegree,  
indegree,  
n,  
n  
4909.  
3682.  
2455.  
1228.  
1.000  
Max: 4909.  
Min: 1.000



NORTHWESTERN  
UNIVERSITY

# ERGM Model

- Interdependent decision to create and dissolve links
- Statistical “MRI” for structural signatures



# Exponential Random Graph ( $p^*$ ) Model



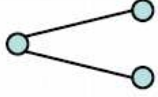
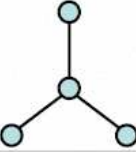
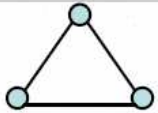
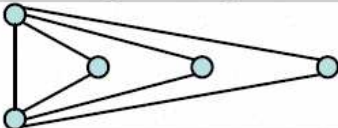
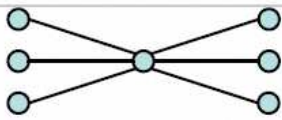
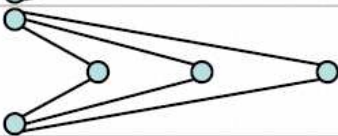
- ERGMs are a class of stochastic models that share the following general form

$$P(Y = y) = \frac{1}{k(\theta)} \exp(\theta^T g(y))$$

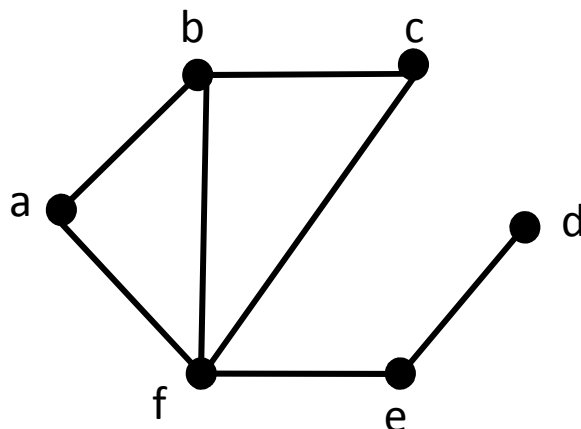
- $Y$  is the network realization, similar to a random variable
- $y$  is the observed network
- $g(y)$  is a vector of network statistics
- $\vartheta^T$  is a vector of coefficient corresponding to  $g(y)$
- $k(\vartheta)$  is a normalizing factor calculated by summing up  $\exp(\vartheta^T g(y))$  over all possible network configurations



# Network Statistics (Undirected Networks)

Edge ( $L$ )		Isolate	
2-Star ( $S_2$ )		3-Star ( $S_3$ )	
Triangle ( $T_1$ )		Alt-Triangle (AT)	
Alt-Star (AS)		Alt-2-Path (A2P)	

Example:



Edge: 7

2-Star:  $1+3+1+0+1+6=12$



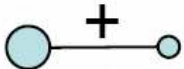

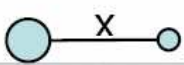


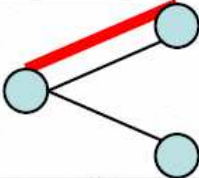
3-Star:  $0+1+0+0+0+4=5$

4-Star: 1

Triangle: 2







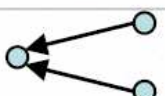

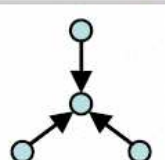
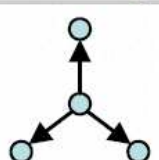


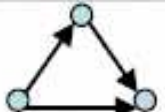
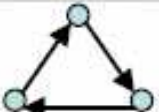
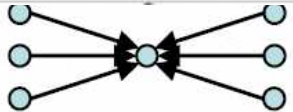
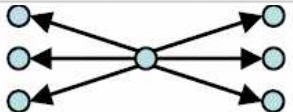
# Attribute Statistics (Undirected Networks)

<ul style="list-style-type: none"> <li>● – actors with attribute</li> <li>○ – actors with or without attribute</li> <li>[Attr] – attribute name</li> </ul>			
[Attr]-interaction		[Attr]-activity	
[Attr]-Sum		[Attr]-difference <sup>1</sup>	
[Attr]-interaction			
Dyadic covariate 			
[Attr]-Edge		[Attr]-S21	

For more statistics, check Pnet manual:

<http://www.sna.unimelb.edu.au/pnet/download/PNet/PNetManual.pdf>

# Network Statistics (Directed Networks)

Arc		Reciprocity	
sink		source	
In-2-star		Out-2-star	
In-3-star		Out-3-star	
2-path		$T_7$	
Transitive Triad ( $T_9$ )		Cyclic Triad ( $T_{10}$ )	
Alt-in-star (AinS)		Alt-out-star (AoutS)	



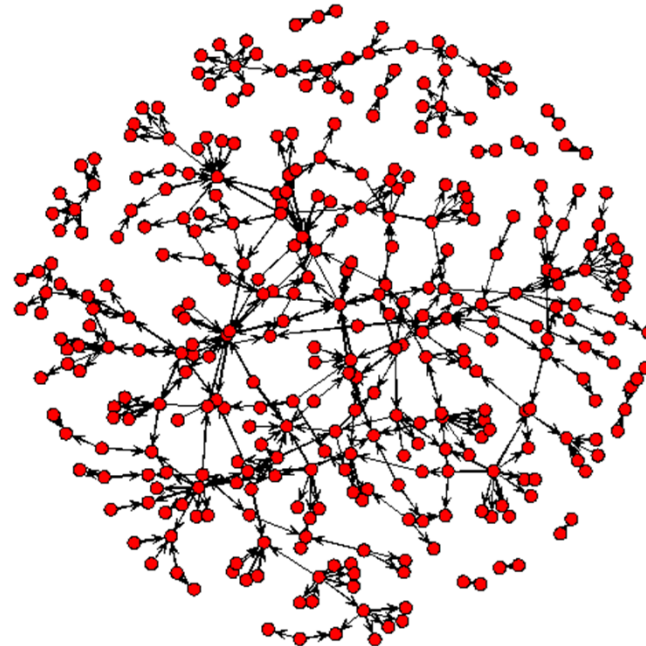
# ERGM Model

- Interdependent decision to create and dissolve links
- Statistical “MRI” for structural signatures
- MCMC simulation
  - Time and space complexity
  - Based on global network statistics (i.e. cannot decompose into local tasks)
  - Inherently sequential



# One wave Snowball Sample

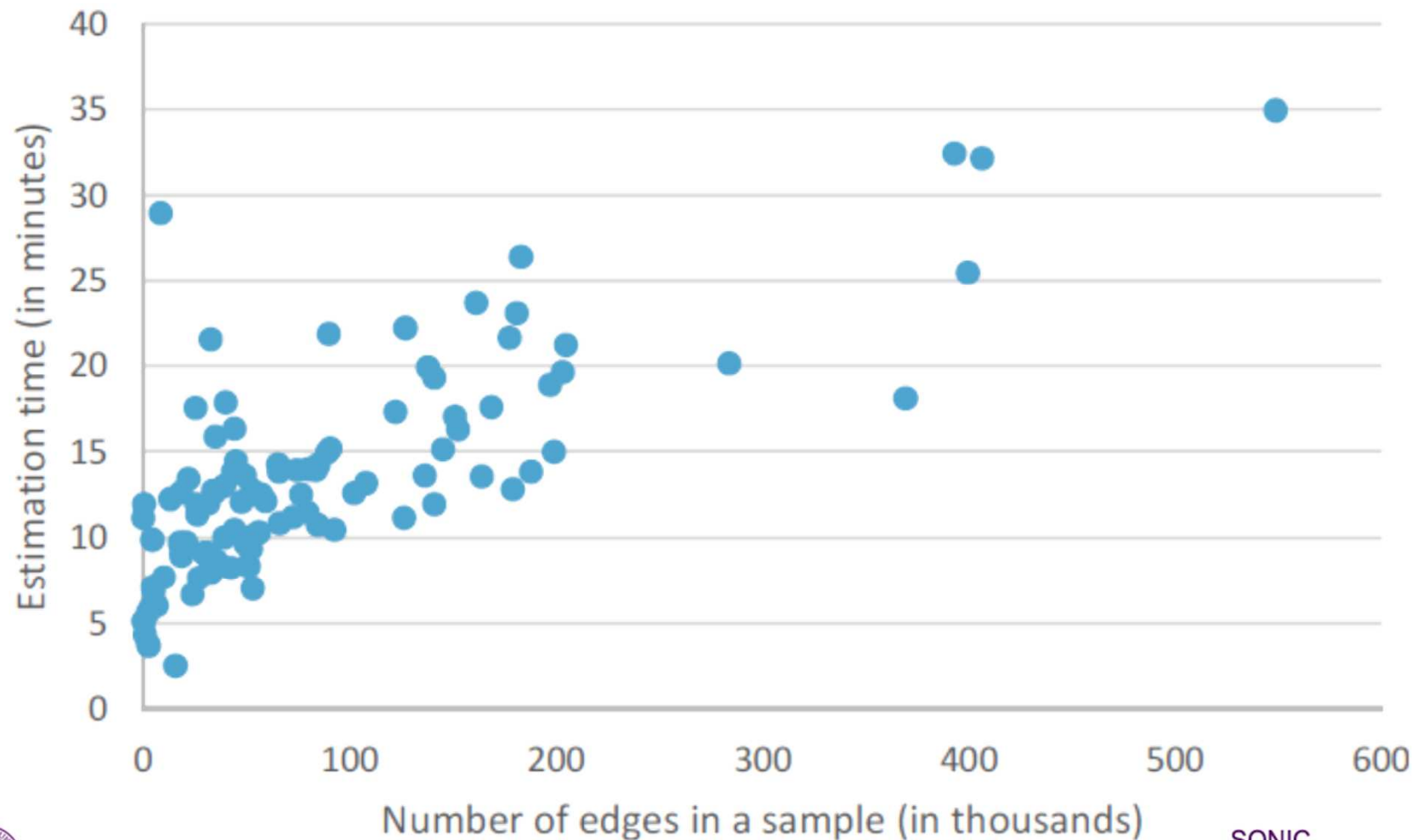
1. Randomly select a user with 1000 to 2000 followers as the seed and find all his/her followers in the first snapshot.
2. Extract unfollow relations among these users



Parameters	Estimate (S.E.)	
Unfollow network structures		
Reciprocal unfollow	<b>3.45*</b>	H2: Supported
Sender attributes of unfollow relations:		
# Followers	<b>0.21*</b>	H3.1:Supported
# Followees	<b>-0.47*</b>	H3.2:Supported
Receiver attributes of unfollow relations:		
# Followers	<b>-0.03*</b>	H3.1:Supported
# Followees	<b>0.07*</b>	H3.2:Supported
Community in follow networks at time 1:		
Mutual following ties	<b>-0.46*(.038)</b>	H1: Supported
# Common followees	<b>-1.83*</b>	H4:Supported
# Common Hashtag	-0.03	H5:Not supported
Sender's interactions to receiver:		
# Replies	0.001	H6: Not supported
# Retweets	-0.009	H6: Not supported
# Mentions	0.321	H6: Not supported
# Favorite	-0.003	H6: Not supported
Unfollow network structures as control variables:		
Edges / density	<b>-1.30*</b>	
Weighted out-stars	<b>-0.57*</b>	
Weighted in-stars	<b>0.19*</b>	



# Performance of Estimation in R/Statnet

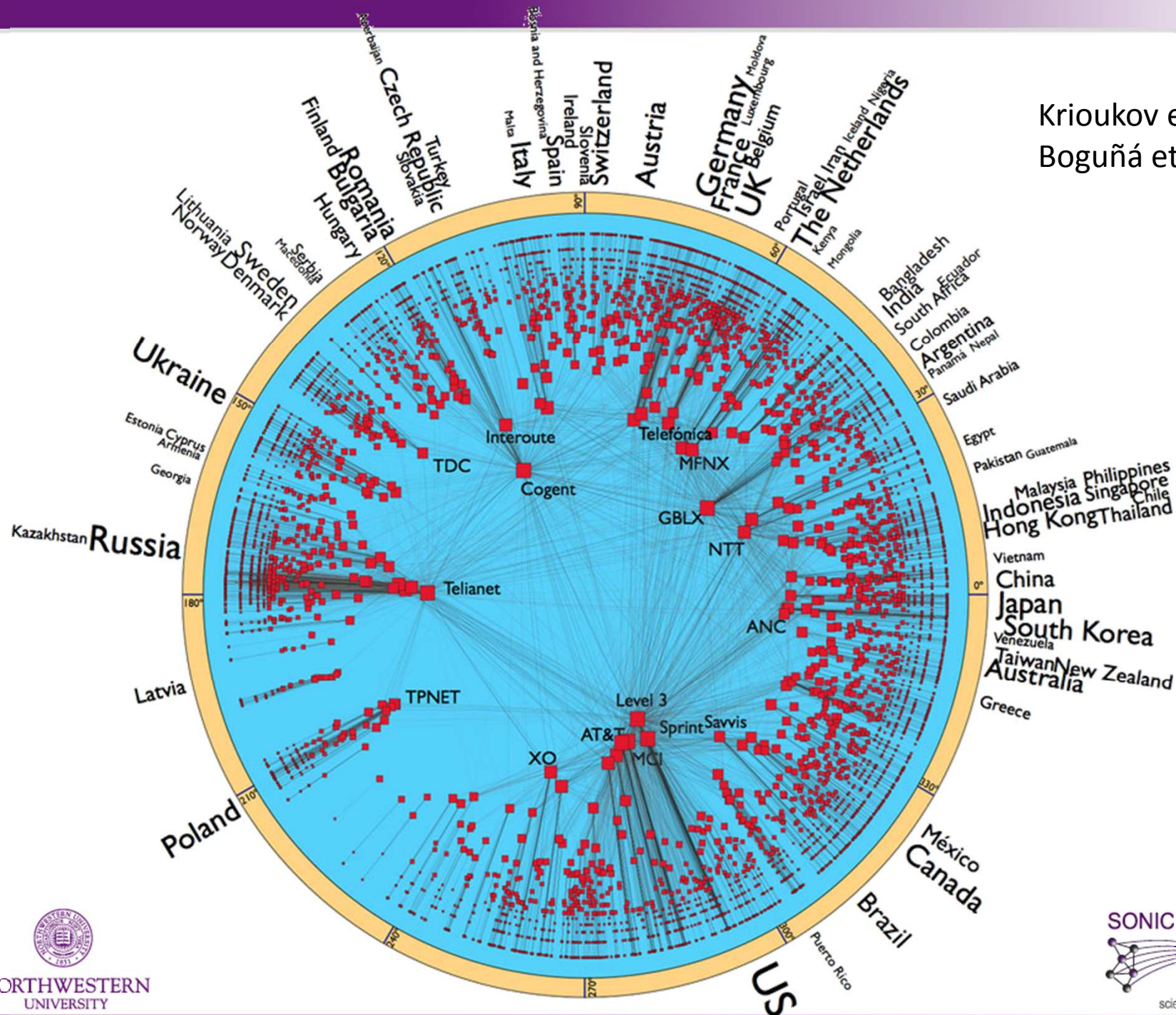


# Performance of ERGM on Gordon

- Difficulty in parallelization
  - Need of global information in the whole network
  - Sequential update in MCMC simulation
- Snowball sampling and meta-analysis
  - Reveal local structures in large networks
  - Utilize the parallel system as a microscope
  - Limitation of snowball sampling
- Future approach
  - Construct locality in a network
  - Map networks to a hyperbolic space



Krioukov et al. 2010  
Boguñá et al. 2010



# Hyperbolic Mapping

- Simplify complex networks
  - From graphs to gyrovector space
  - From edges to distances
  - Foundation of task decomposition
- Incorporate heterogeneous structures
  - Detect and visualize communities and structures in complex networks.
- Simplify statistical methods
  - Convert exponential random graphs to auxiliary fields

