

# Network Analysis in the Age of Large Network Dataset Collections – Challenges, Solutions and Applications

## *Tutorial*



Jérôme KUNEGIS, Renaud LAMBIOTTE  
with acknowledgments to many people

10 Nov 2017  
CIKM'17

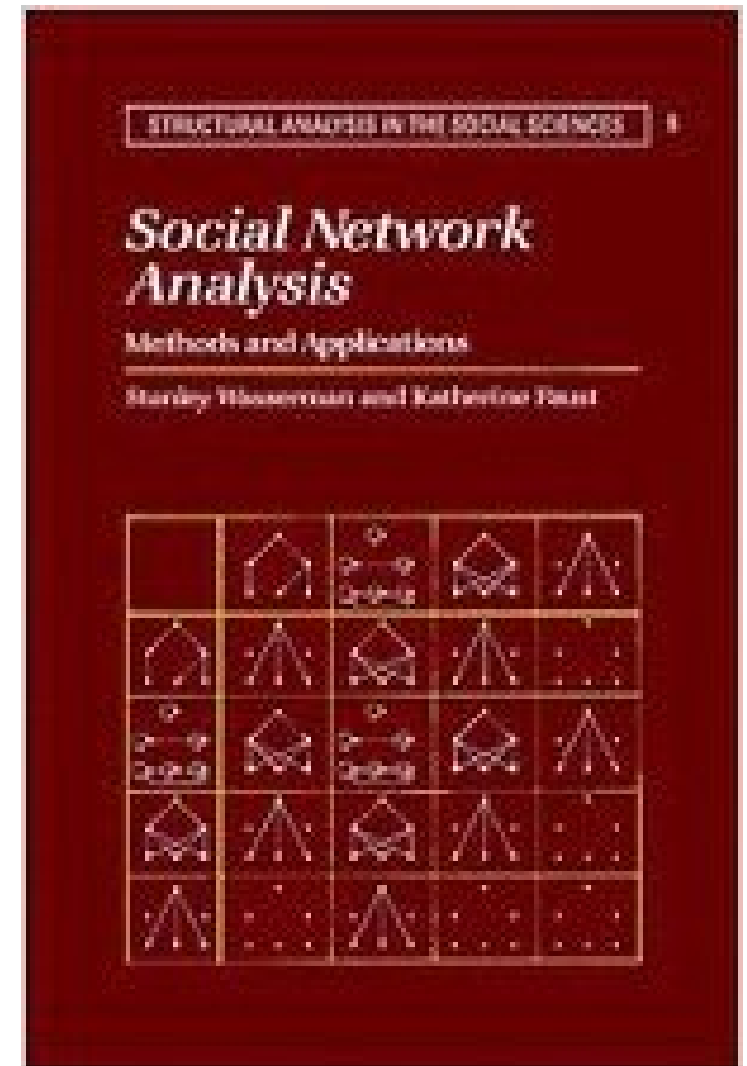
<http://KONECT.cc/>



# Social Network Analysis without the Web

“Social Network Analysis” by  
Wasserman & Faust first edition 1994

Contains 18 datasets (based on 19<sup>th</sup> printing)



# WWW 2014 Best Paper Nominations

- **Community-Based Bayesian Aggregation Models for Crowdsourcing**
- **4** datasets (crowdsourcing)
- **Efficient Estimation for High Similarities using Odd Sketches**
- **5** real-world datasets + synthetic dataset (text documents)
- **Local Collaborative Ranking**
- **3** datasets (rating networks)
- **Engaging with Massive Online Courses**
- **1** dataset (case study)

The best paper award goes to...



# Why Do Researchers Use Multiple Datasets?

- To cover more application areas
- To show that results are generalizable
- To make results more statistically significant

# Showing that Algorithm X is Better Than Y

- Setup: We want to compare two prediction algorithms X and Y
- Experiment: Apply X and Y to dataset A
- Result: X has higher precision than Y
- Conclusion: “Algorithm X performs better than algorithm Y”

Really?

# Let's Make More Experiments

- Add a dataset B to the experiments
- On dataset B, Y performs better than X



No!

# More Datasets

- Add more datasets to the mix
- Algorithm X is better than algorithm Y with 6 out of 10 datasets
- Under Null hypothesis of equal probability of X performing better than Y on any one dataset and results on datasets being independent, the probability that this happens is 17%, i.e., not significant!
- Need about 65 datasets to get a statistically significant result at ( $p \leq 0.05$ ) for a 60% result.



# Who We Are



Jérôme KUNEGIS  
University of Namur  
Belgium



Renaud LAMBIOTTE  
University of Namur  
University of Oxford

# The KONECT.cc Project - Koblenz Network Collection



# On the Spectral Evolution of Large Networks

Jérôme Kunegis

Institute for Web Science and Technologies  
University of Koblenz-Landau  
kunegis@uni-koblenz.de

November 2011

Vom Promotionsausschuss des Fachbereichs 4: Informatik der Universität  
Koblenz-Landau zur Verleihung des akademischen Grades

**Doktor der Naturwissenschaften (Dr. rer. nat.)**

genehmigte Dissertation.

PhD thesis at the University of Koblenz-Landau.

Datum der wissenschaftlichen Aussprache:  
Vorsitz des Promotionsausschusses:  
Berichterstatter:  
Berichterstatter:  
Berichterstatter:

9. November 2011  
Prof. Dr. Karin Harbusch  
Prof. Dr. Steffen Staab  
Prof. Dr. Christian Bauerkhage  
Prof. Dr. Klaus Obermayer



The trick  
is...



Everything  
is a  
NETWORK!

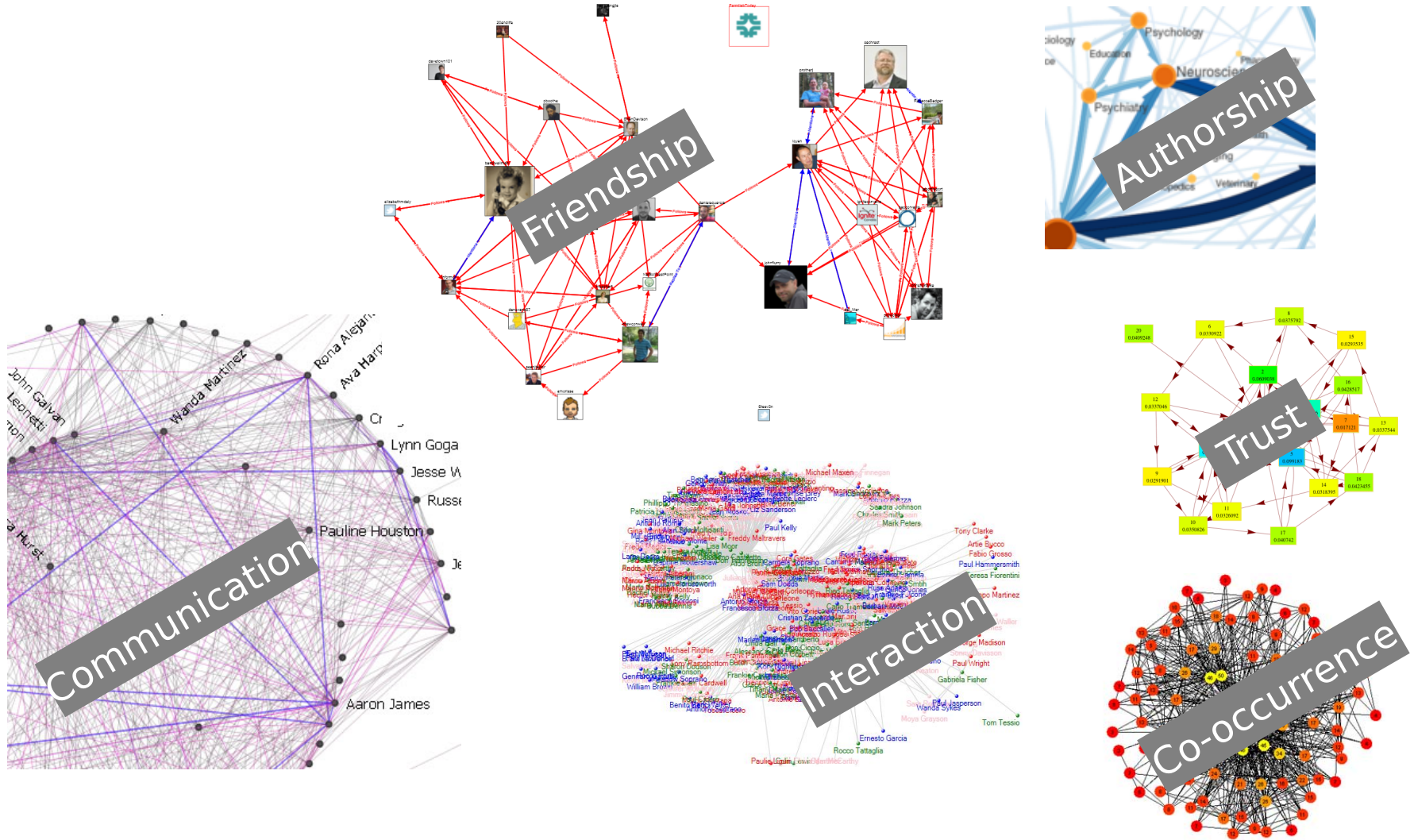
# Explanation : *Why Everything Is a Network*

Argument by considering the structure of large systems :

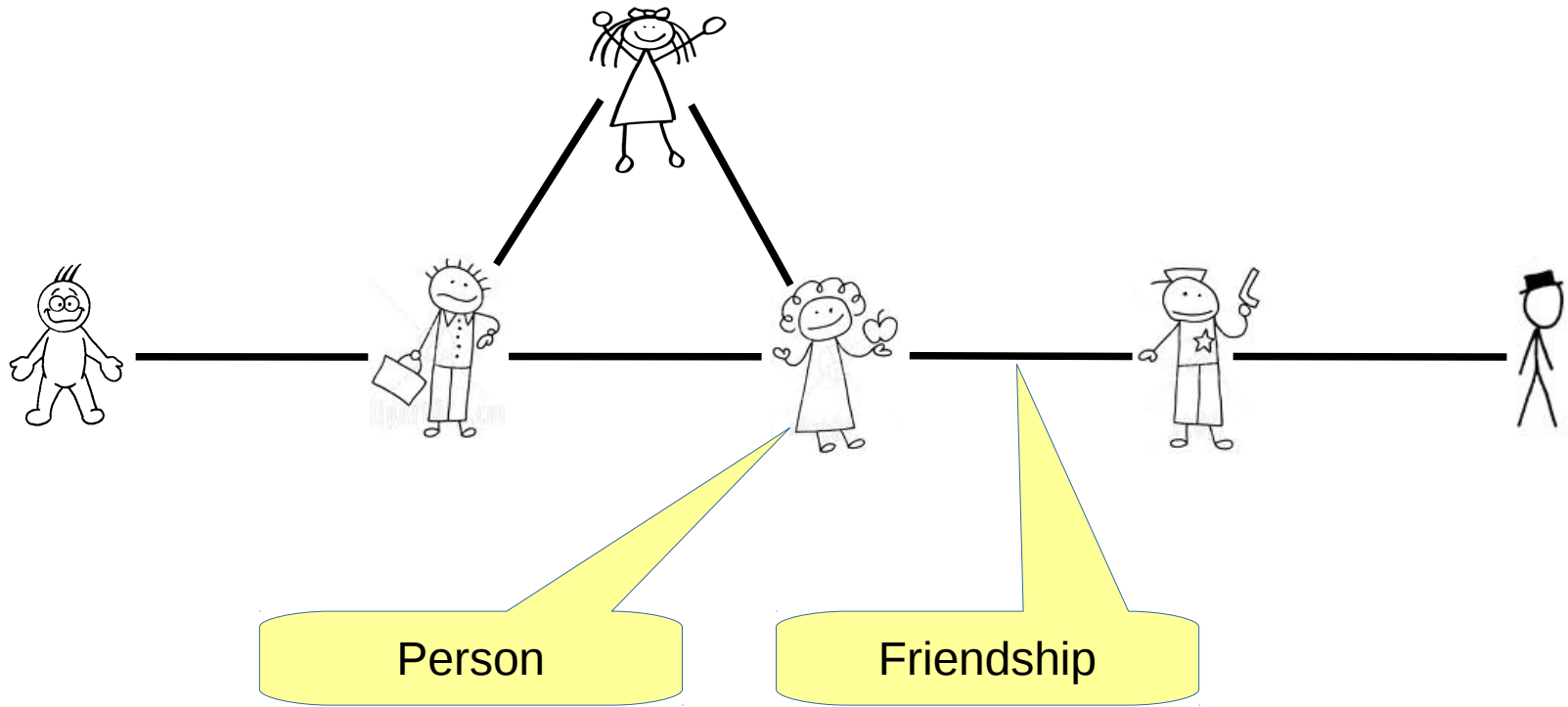
- Large systems consist of many parts
- Any part can only communicate with a finite number of other parts
- This gives rise to a sparse network structure

Caveat : May not be a complex system (e.g. a grid)

# Well, Only *Almost* Everything Is a Network



# Social Network

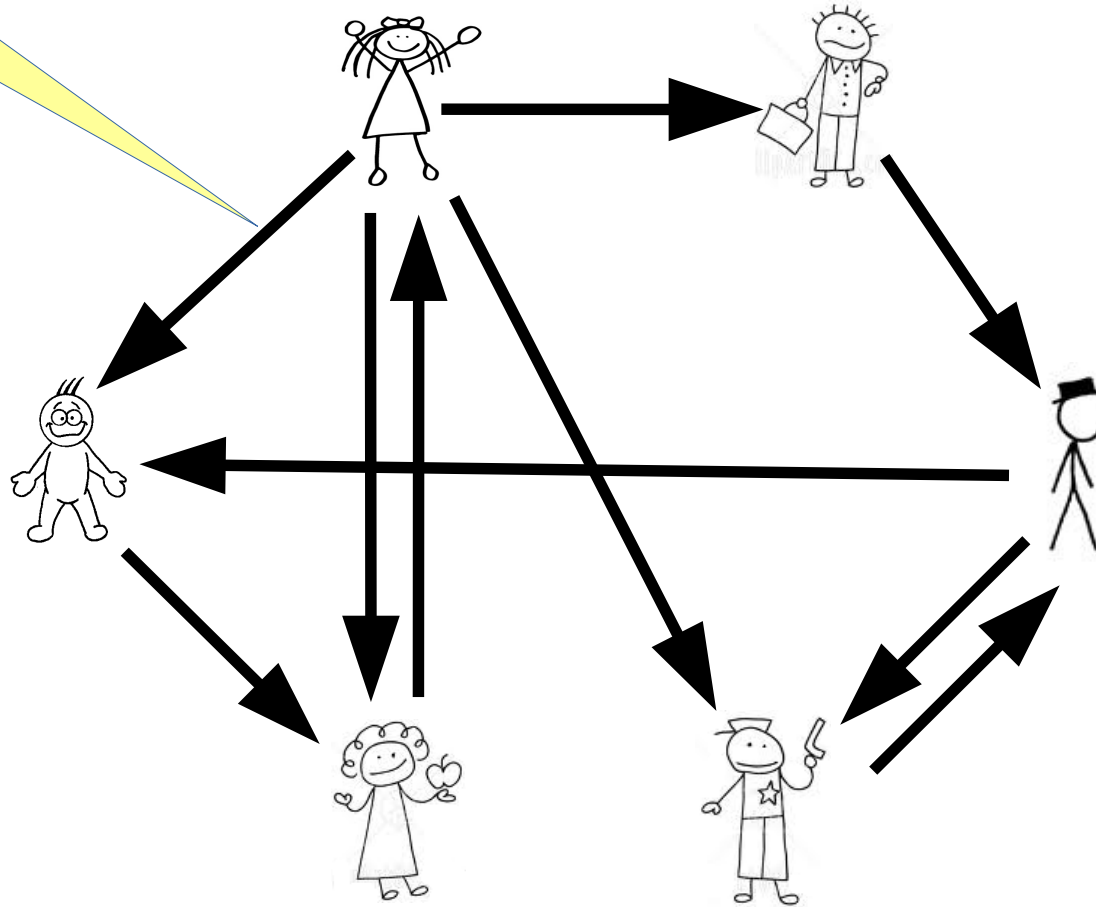




# Trust Network

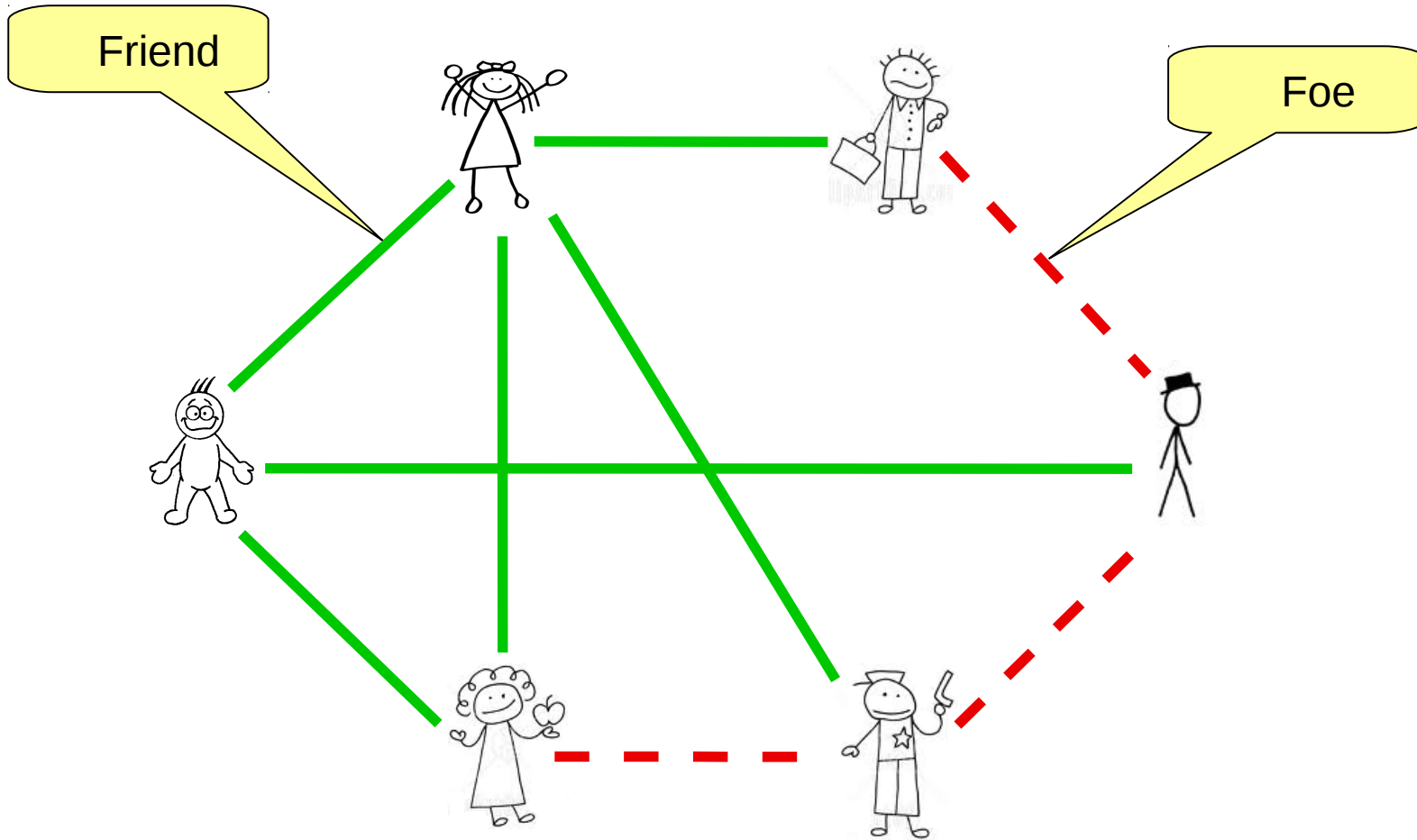


Trust

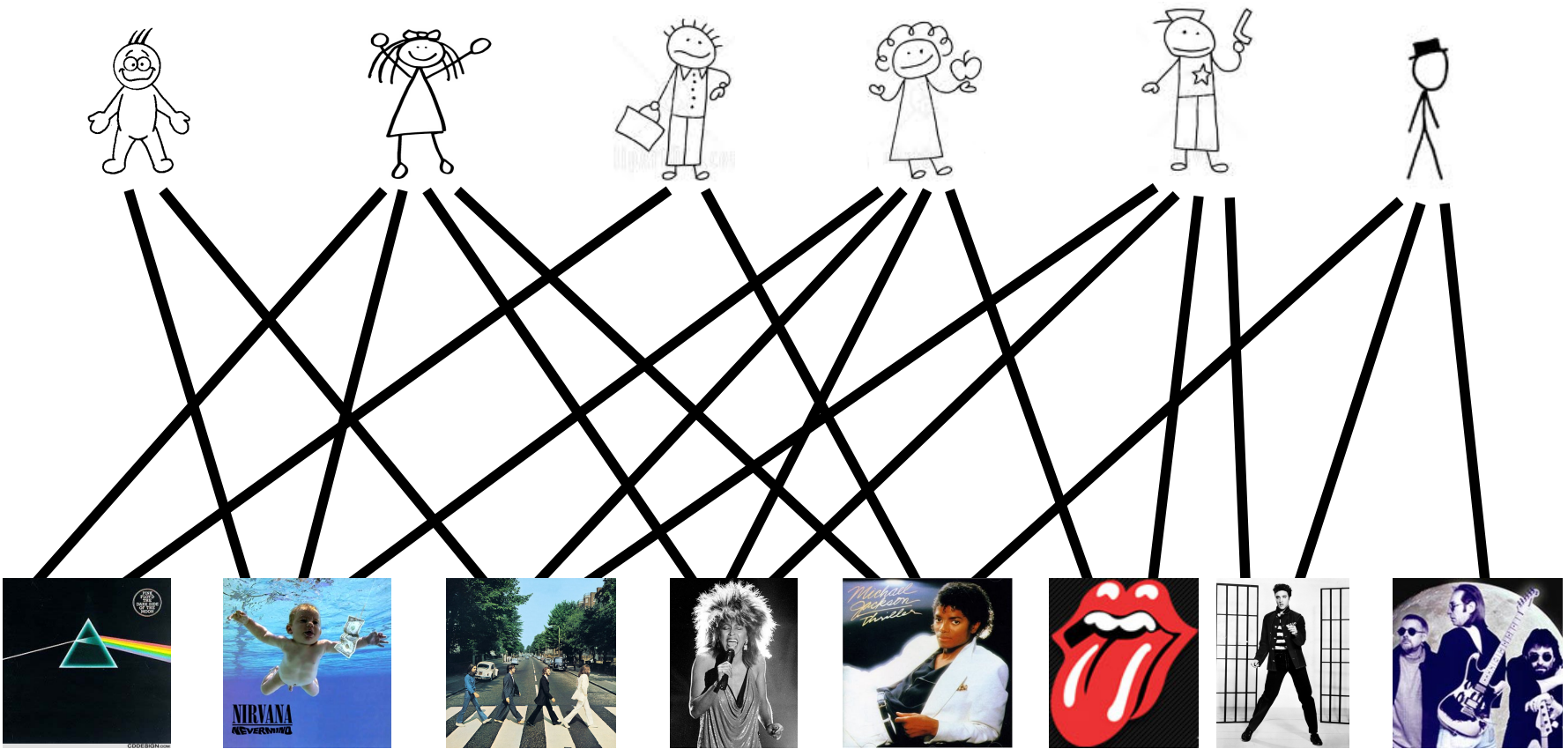


# Signed Network

# Slashdot



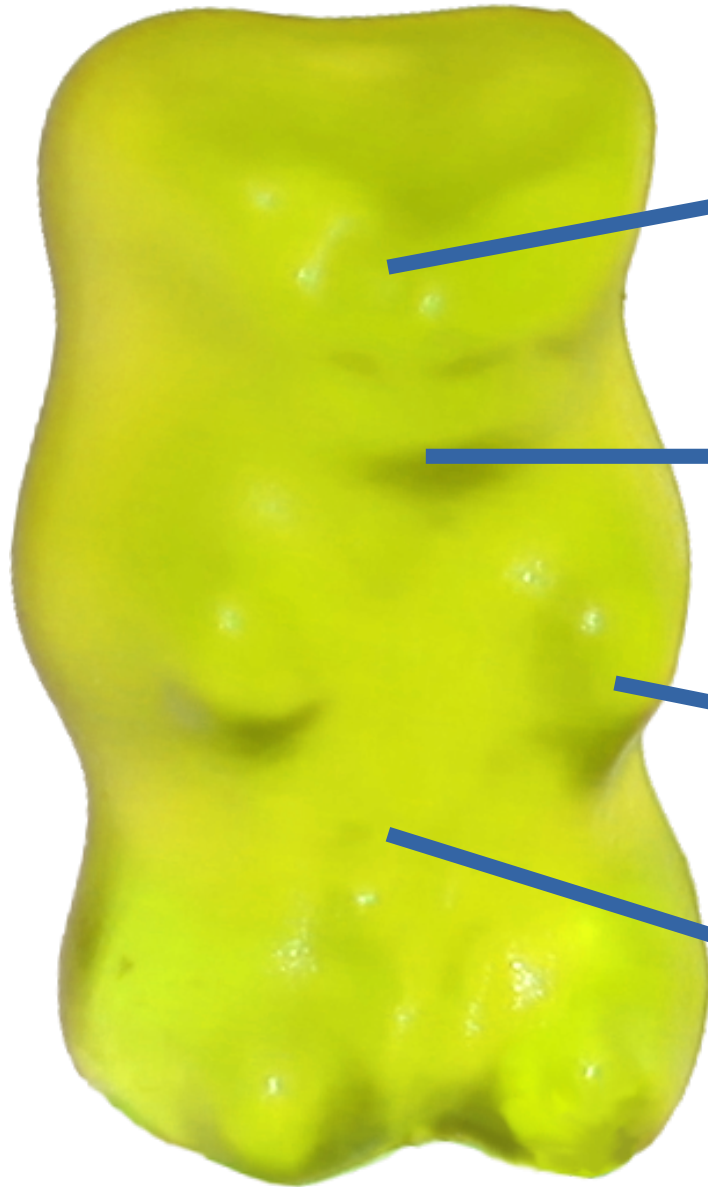
# Bipartite Network



# A Network Dataset Is Like a Gummi Bear



# A Network Dataset Is Like a Gummi Bear



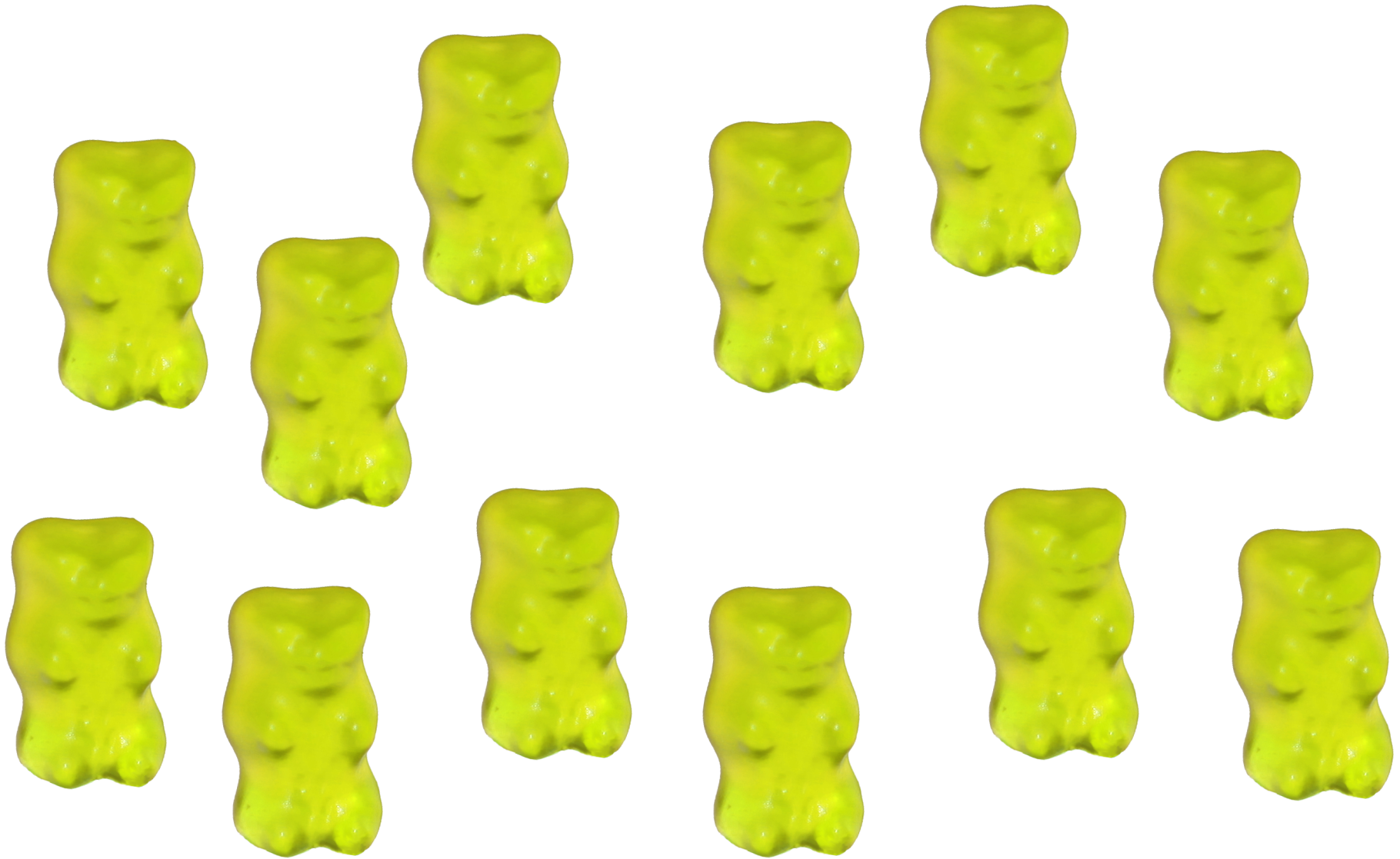
Lots of content  
to analyse

Test network  
models

Evaluate  
prediction  
algorithms

Test scalability of  
algorithms

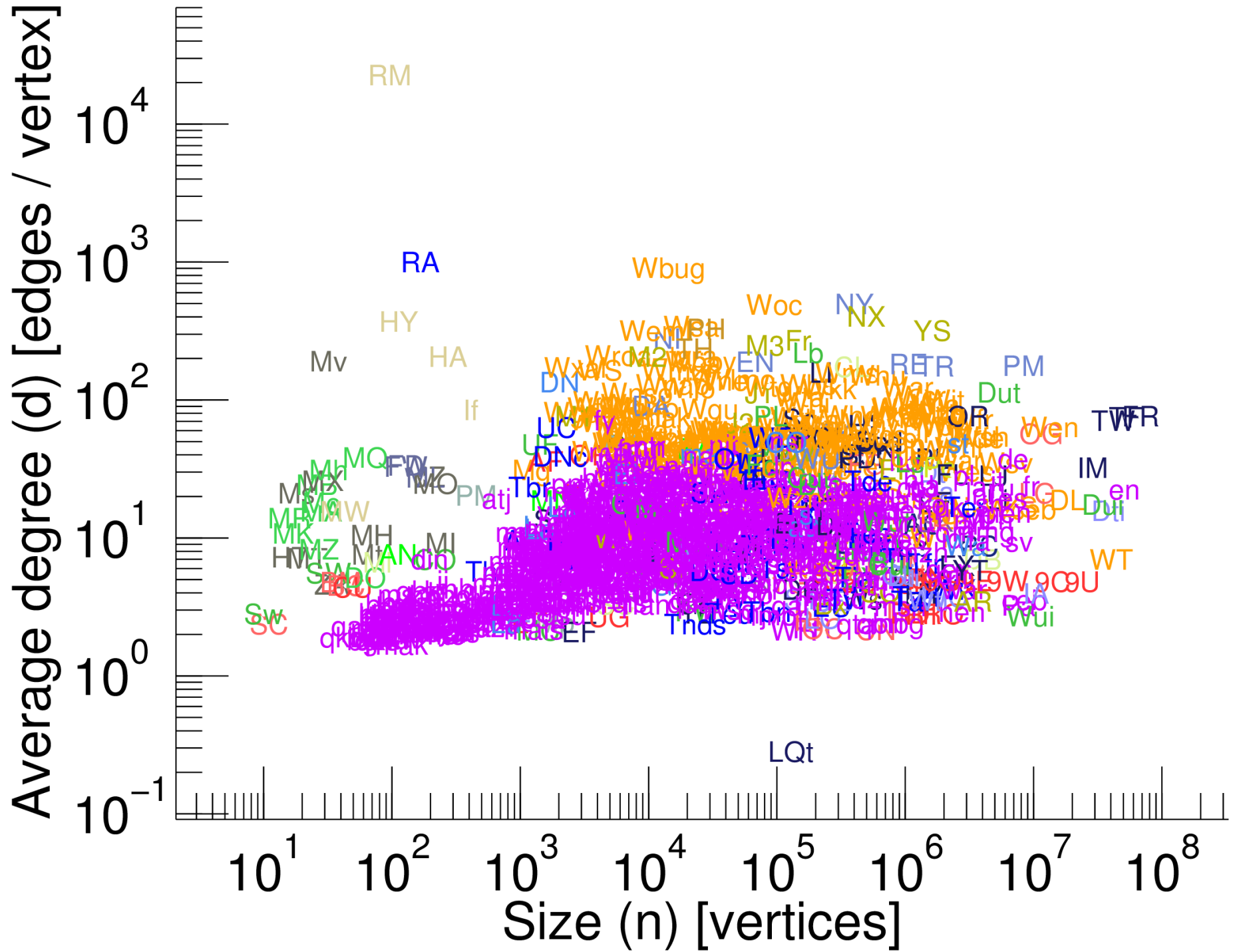
When You Have Tested One, You Have Tested All ?!



Or Do You?



# Diversity of Network Datasets



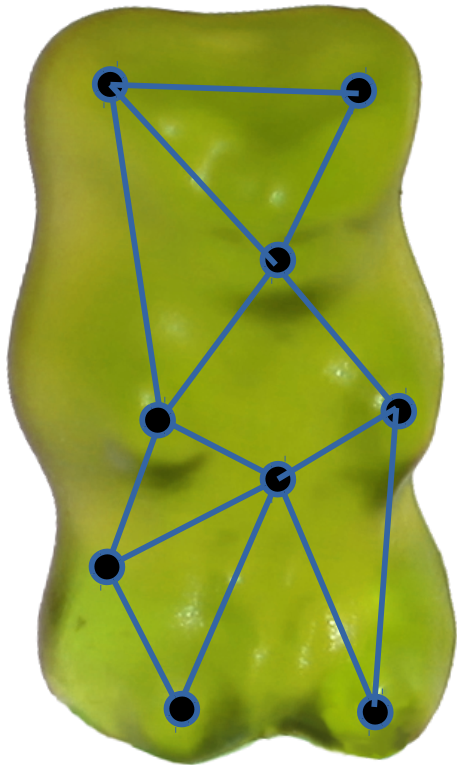


# Network Categories

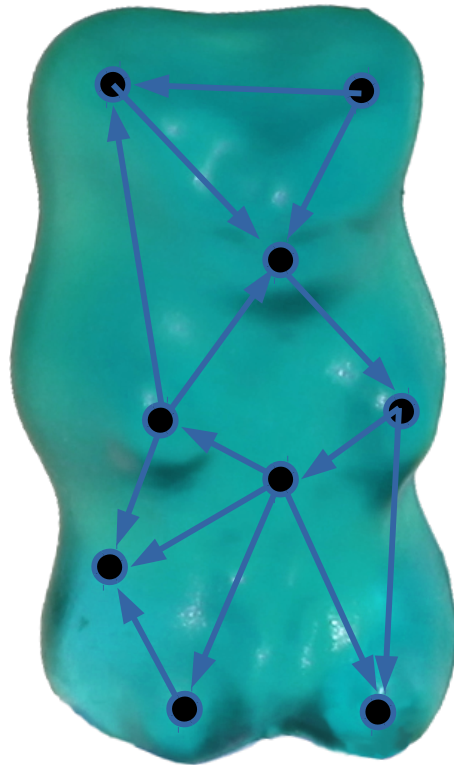
	Internal name	Vertices	Edges	Properties	Count
●	Affiliation	Actors, groups	Membership	B - =	16
●	Animal	Animals	Tie	U D - +	9
●	Authorship	Authors, works	Authorship	B - =	808
●	Citation	Documents	Citation	D -	7
●	Coauthorship	Authors	Coauthorship	U - =	3
●	Cocitation	Authors	Cocitation	U =	2
●	Communication	Persons	Message	U D - =	42
●	Computer	Computers	Connection	U D - =	13
●	Feature	Items, features	Property	B - = +	17
●	HumanContact	Persons	Real-life contact	U = +	5
●	HumanSocial	Persons	Real-life tie	U D - + ± <sup>+</sup>	12
●	Hyperlink	Web page	Hyperlink	D B - = ⇒	191
●	Infrastructure	Location	Connection	U D - = +	23
●	Interaction	Persons, items	Interaction	D B - = <sup>+</sup>	25
●	Lexical	Words	Lexical relationship	U D - =	5
●	Metabolic	Metabolites	Interaction	U D - =	7
●	Misc	Various	Various	U D - = +	11
●	OnlineContact	Users	Online interaction	U D - = ± <sup>+</sup> ⇒	15
●	Rating	Users, items	Rating	B - = * * *	15
●	Social	Persons	Online tie	U D - + ± * ⇒	46
●	Software	Software Component	Dependency	D - =	3
●	Text	Documents, words	Occurrence	B =	10
●	Trophic	Species	Carbon exchange	D - +	3
<b>Total</b>					<b>1288</b>

# Network Formats

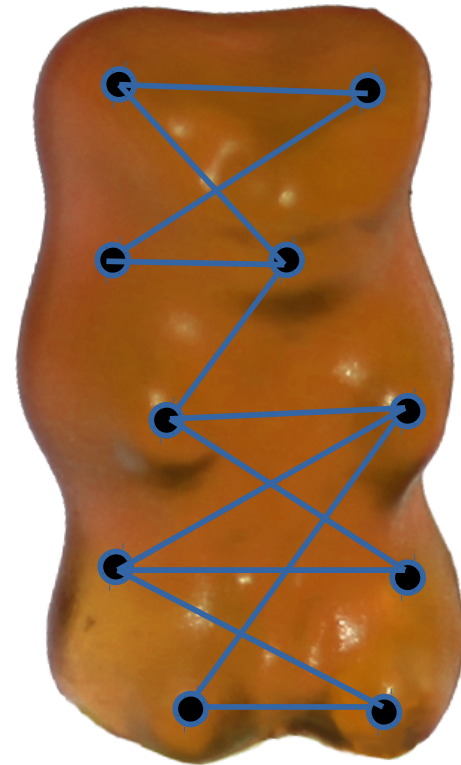
U • Undirected



D • Directed



B • Bipartite

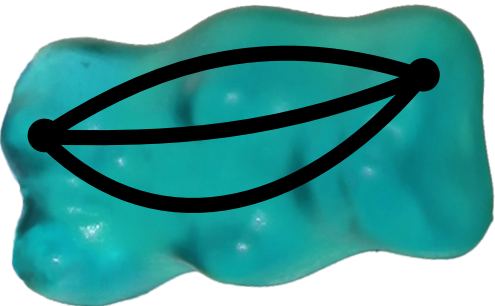


# Edge Weight and Multiplicity Types

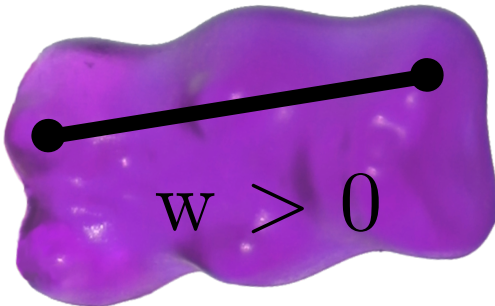
- • Unweighted



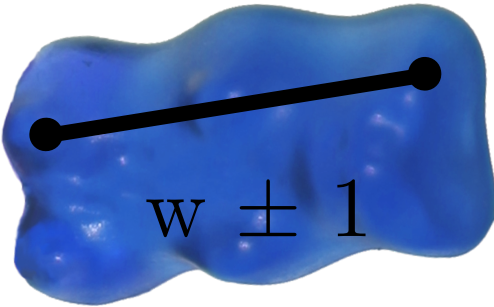
= • Multiple



+ • Positive



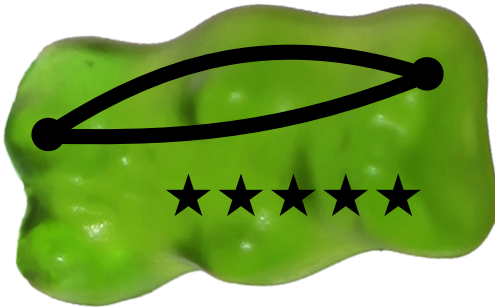
$\pm$  • Signed



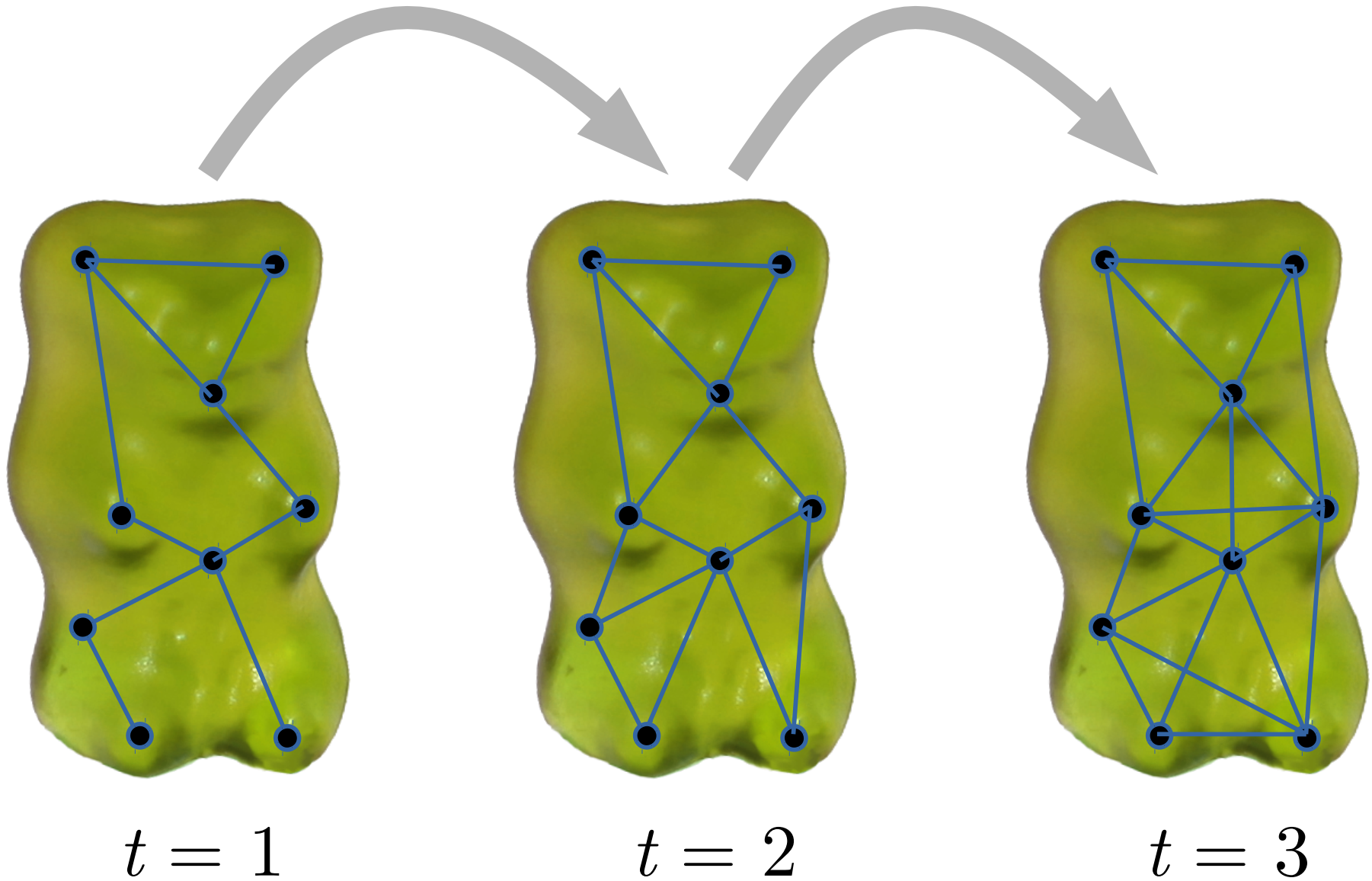
\* • Rating



\*<sup>\*</sup> • Multiple Ratings



# Timestamps



# Collections of Network Datasets

- **SNAP.stanford.edu**

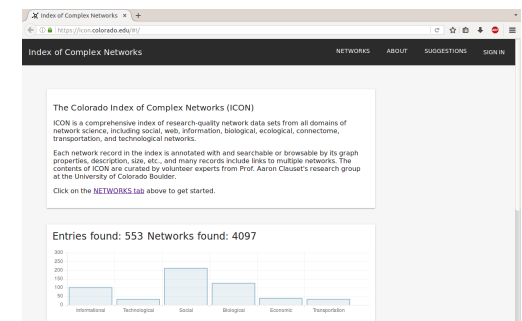
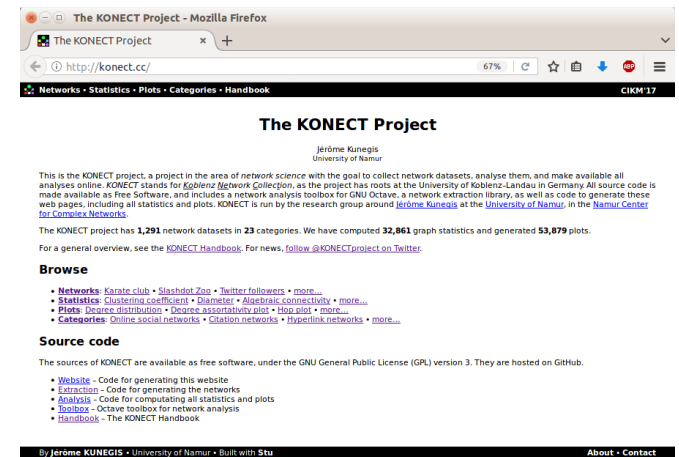
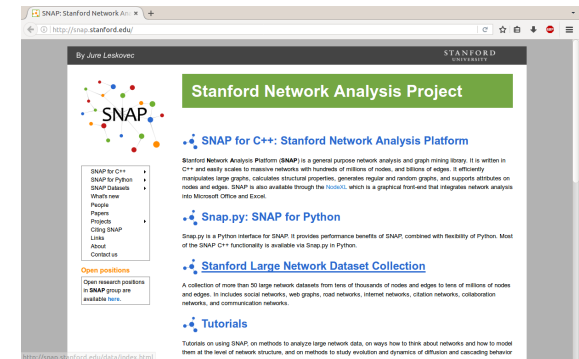
- « Stanford Network Analysis Project »
- by Jure Leskovec, Stanford Univ. (~2009)
- several 100 networks; not systematic
- Available for download
- Some statistics available

- **KONECT.cc**

- « Koblenz Network Collection »
- by Jérôme Kunegis, Univ. of Namur (~2011)
- 1200+ networks
- Most networks available for download
- Many statistics available

- **ICON.colorado.edu**

- « Index of Complex Networks »
- by Aaron Clauset, Univ. of Colorado (~2016)
- 4000+ datasets
- Not available for download (“index”)



# The KONECT Project

Jérôme Kunegis  
University of Namur

This is the KONECT project, a project in the area of *network science* with the goal to collect network datasets, analyse them, and make available all analyses online. *KONECT* stands for *Koblenz Network Collection*, as the project has roots at the University of Koblenz-Landau in Germany. All source code is made available as Free Software, and includes a network analysis toolbox for GNU Octave, a network extraction library, as well as code to generate these web pages, including all statistics and plots. KONECT is run by the research group around [Jérôme Kunegis](#) at the [University of Namur](#), in the [Namur Center for Complex Networks](#).

The KONECT project has **1,291** network datasets in **23** categories. We have computed **32,861** graph statistics and generated **53,879** plots.

For a general overview, see the [KONECT Handbook](#). For news, [follow @KONECTproject on Twitter](#).

## Browse

- **Networks:** [Karate club](#) • [Slashdot Zoo](#) • [Twitter followers](#) • [more...](#)
- **Statistics:** [Clustering coefficient](#) • [Diameter](#) • [Algebraic connectivity](#) • [more...](#)
- **Plots:** [Degree distribution](#) • [Degree assortativity plot](#) • [Hop plot](#) • [more...](#)
- **Categories:** [Online social networks](#) • [Citation networks](#) • [Hyperlink networks](#) • [more...](#)

## Source code

The sources of KONECT are available as free software, under the GNU General Public License (GPL) version 3. They are hosted on GitHub.

- [Website](#) - Code for generating this website
- [Extraction](#) - Code for generating the networks
- [Analysis](#) - Code for computing all statistics and plots
- [Toolbox](#) - Octave toolbox for network analysis
- [Handbook](#) - The KONECT Handbook

<http://konect.cc/networks/ucidata-zachary/>

## Zachary karate club

This is the well-known and much-used Zachary karate club network. The data was collected from the members of a university karate club by Wayne Zachary in 1977. Each node represents a member of the club, and each edge represents a tie between two members of the club. The network is undirected. An often discussed problem using this dataset is to find the two groups of people into which the karate club split after an argument between two teachers.

### Metadata

Code	ZA
Internal name	ucidata-zachary
Name	Zachary karate club
Data source	<a href="http://vlado.fmf.uni-lj.si/pub/networks/data/ucinet/ucidata.htm#zachary">http://vlado.fmf.uni-lj.si/pub/networks/data/ucinet/ucidata.htm#zachary</a>
Consistency check	<input checked="" type="checkbox"/> Dataset passed all tests
<u>Category</u>	<input type="radio"/> <a href="#">Human social network</a>
Dataset timestamp	1977
Node meaning	Member
Edge meaning	Tie
Network format	<input checked="" type="checkbox"/> Unipartite, undirected
Edge type	<input type="checkbox"/> Unweighted, no multiple edges
Loops	<input type="checkbox"/> Does not contain loops

### Statistics

<u>Size</u>	n = 34
<u>Volume</u>	m = 78
<u>Average degree</u>	d = 4.588 24
<u>Maximum degree</u>	d <sub>max</sub> = 17
<u>Fill</u>	p = 0.139 037
<u>Wedge count</u>	s = 528
<u>Claw count</u>	z = 1.764

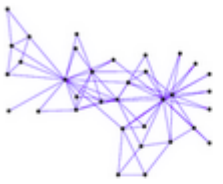
# Statistics

<u>Size</u>	$n = 34$
<u>Volume</u>	$m = 78$
<u>Average degree</u>	$d = 4.588\ 24$
<u>Maximum degree</u>	$d_{\max} = 17$
<u>Fill</u>	$p = 0.139\ 037$
<u>Wedge count</u>	$s = 528$
<u>Claw count</u>	$z = 1,764$
<u>Cross count</u>	$x = 5,082$
<u>Size of LCC</u>	$N = 34$
<u>Relative size of LCC</u>	$N^{\text{rel}} = 1.000\ 00$
<u>Degree assortativity</u>	$\rho = -0.475\ 613$
<u>Degree assortativity p-value</u>	$P_{\rho} = 3.509\ 45 \times 10^{-10}$
<u>Spectral norm</u>	$\ A\ _2 = 6.725\ 70$
<u>Gini coefficient</u>	$G = 0.385\ 370$
<u>Power law exponent</u>	$\gamma = 1.780\ 96$
<u>Tail power law exponent</u>	$\gamma_t = 2.161\ 00$
<u>Relative edge distribution entropy</u>	$H_{\text{er}} = 0.924\ 709$
<u>Clustering coefficient</u>	$c = 0.255\ 682$
<u>Triangle count</u>	$t = 45$
<u>Diameter</u>	$\delta = 5$
<u>50-Percentile effective diameter</u>	$\delta_{0.5} = 1.840\ 54$
<u>90-Percentile effective diameter</u>	$\delta_{0.9} = 3.441\ 98$
<u>Mean distance</u>	$\delta_m = 2.443\ 26$
<u>4-Tour count</u>	$T_4 = 3,500$
<u>Square count</u>	$q = 154$
<u>Algebraic connectivity</u>	$a = 0.468\ 525$

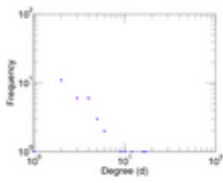


# Plots

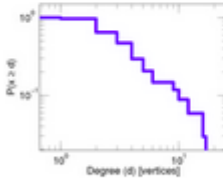
## Fruchterman-Reingold graph drawing



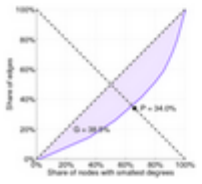
## Degree distribution



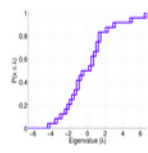
## Cumulative degree distribution



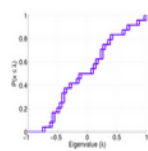
## Lorenz curve



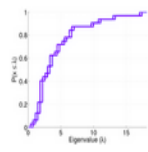
## Spectral distribution of the adjacency matrix



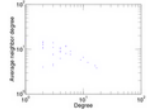
## Spectral distribution of the normalized adjacency matrix



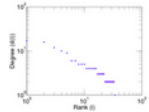
## Spectral distribution of the Laplacian



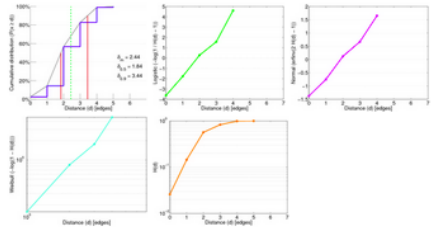
## Degree assortativity



## Zipf plot



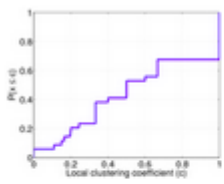
## Hop distribution



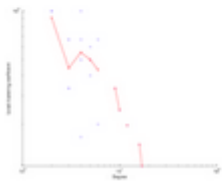
## Double Laplacian graph drawing



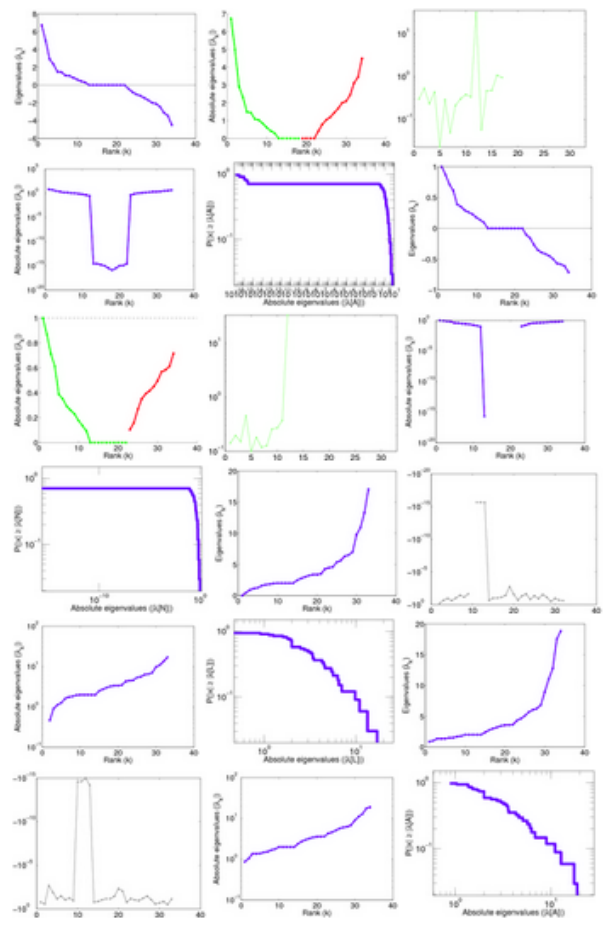
## Clustering coefficient distribution



## Average neighbor degree distribution



## Matrix decompositions plots



## Downloads

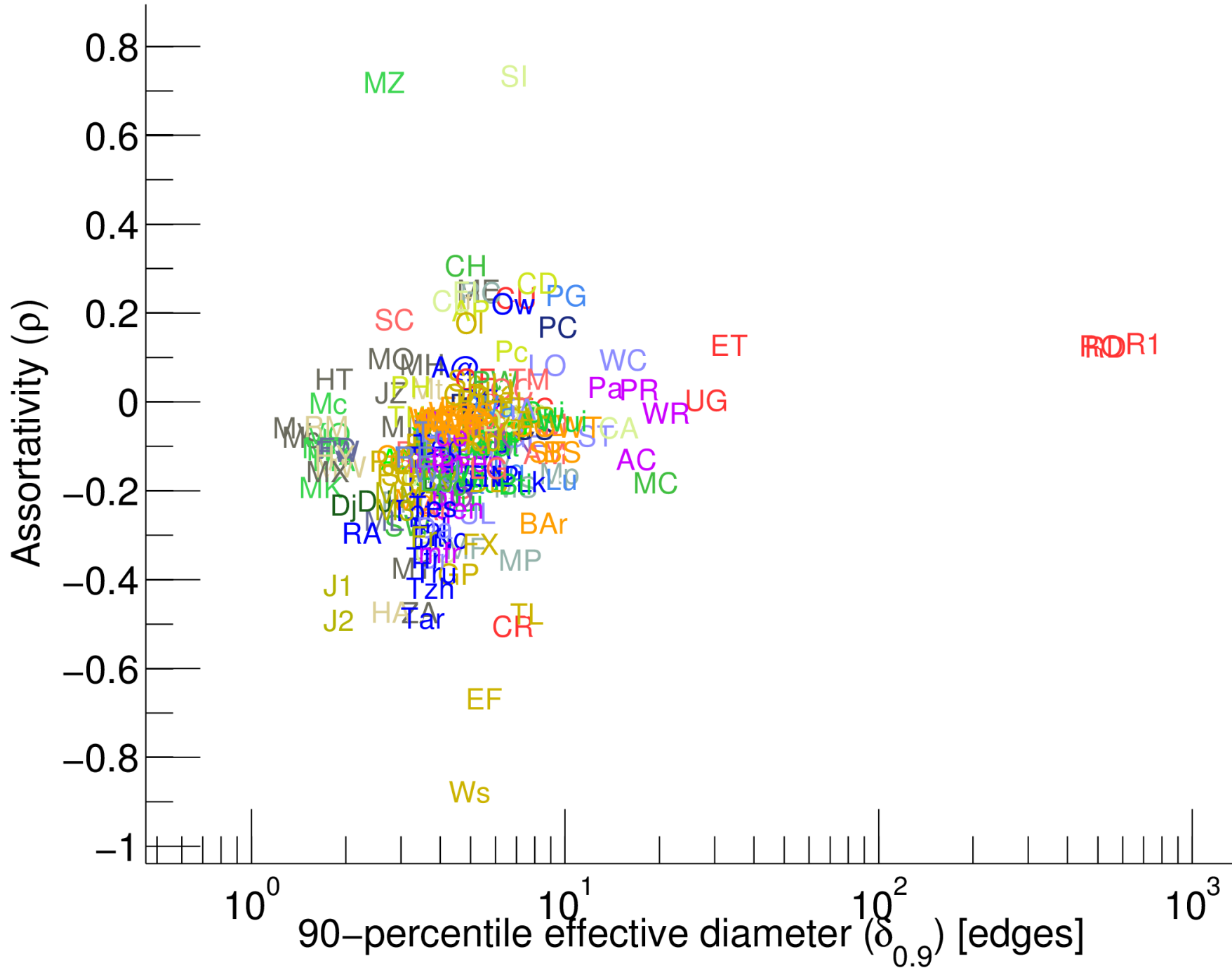
- [Data as TSV](#) (1,856 bytes)
- [Extraction code](#) (GitHub)

## References

- [1] Jérôme Kunegis. KONECT – The Koblenz Network Collection. In *Proc. Int. Conf. on World Wide Web Companion*, pages 1343–1350, 2013. [ [http](#) ]
- [2] Wayne Zachary. An information flow model for conflict and fission in small groups. *J. of Anthropol. Res.*, 33:452–473, 1977.

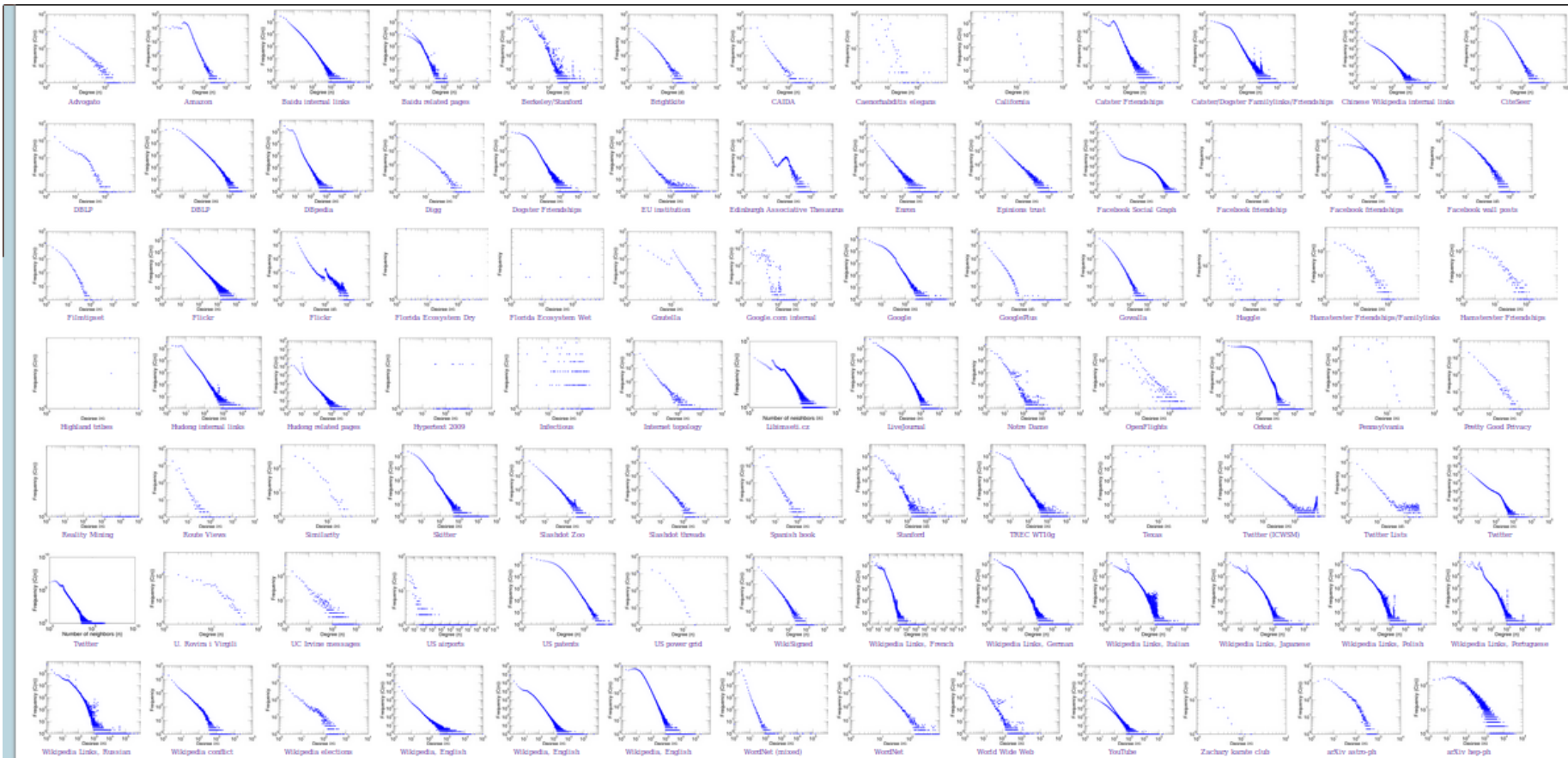
# Network Comparison <http://konect.cc/statistics/>

Command: `stu @scatter.diameff90.assortativity`



# Network Comparison: Plots

<http://konect.cc/plots/>

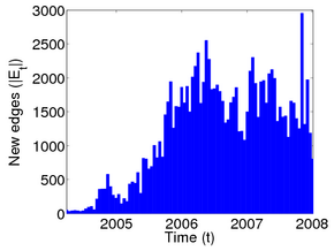


Example: Degree distribution

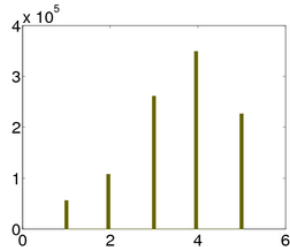
Command: `stu @degree`

# More Plots

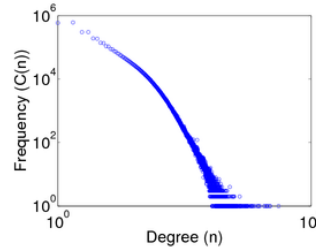
<http://konect.cc/plots/>



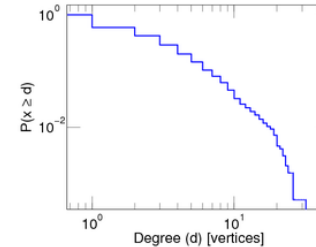
Temporal distribution



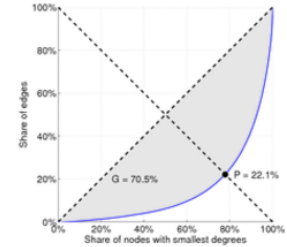
Edge weight and multiplicity distribution



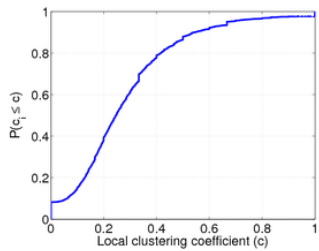
Degree distribution



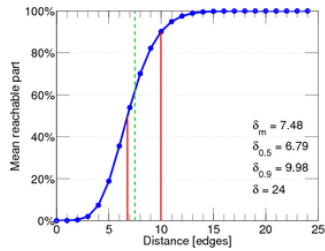
Cumulative degree distribution



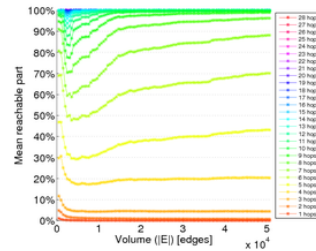
Lorenz curve



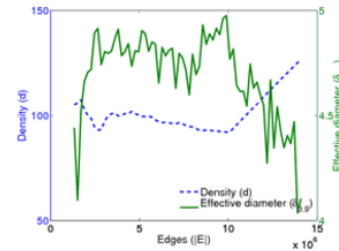
Clustering coefficient distribution



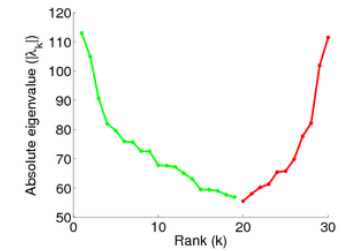
Hop plot



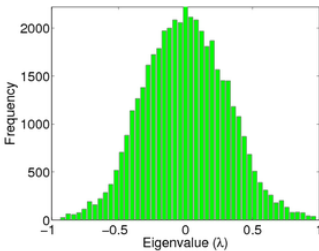
Temporal hop plot



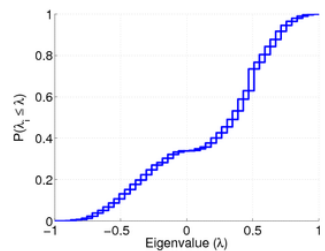
Evolution of average degree and diameter



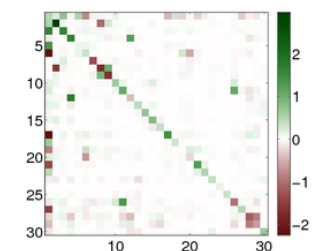
Top-k eigenvalues



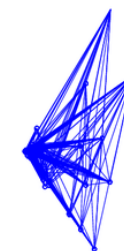
Spectral distribution



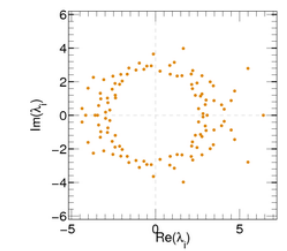
Cumulative spectral distribution



Spectral diagonality test



Network visualization



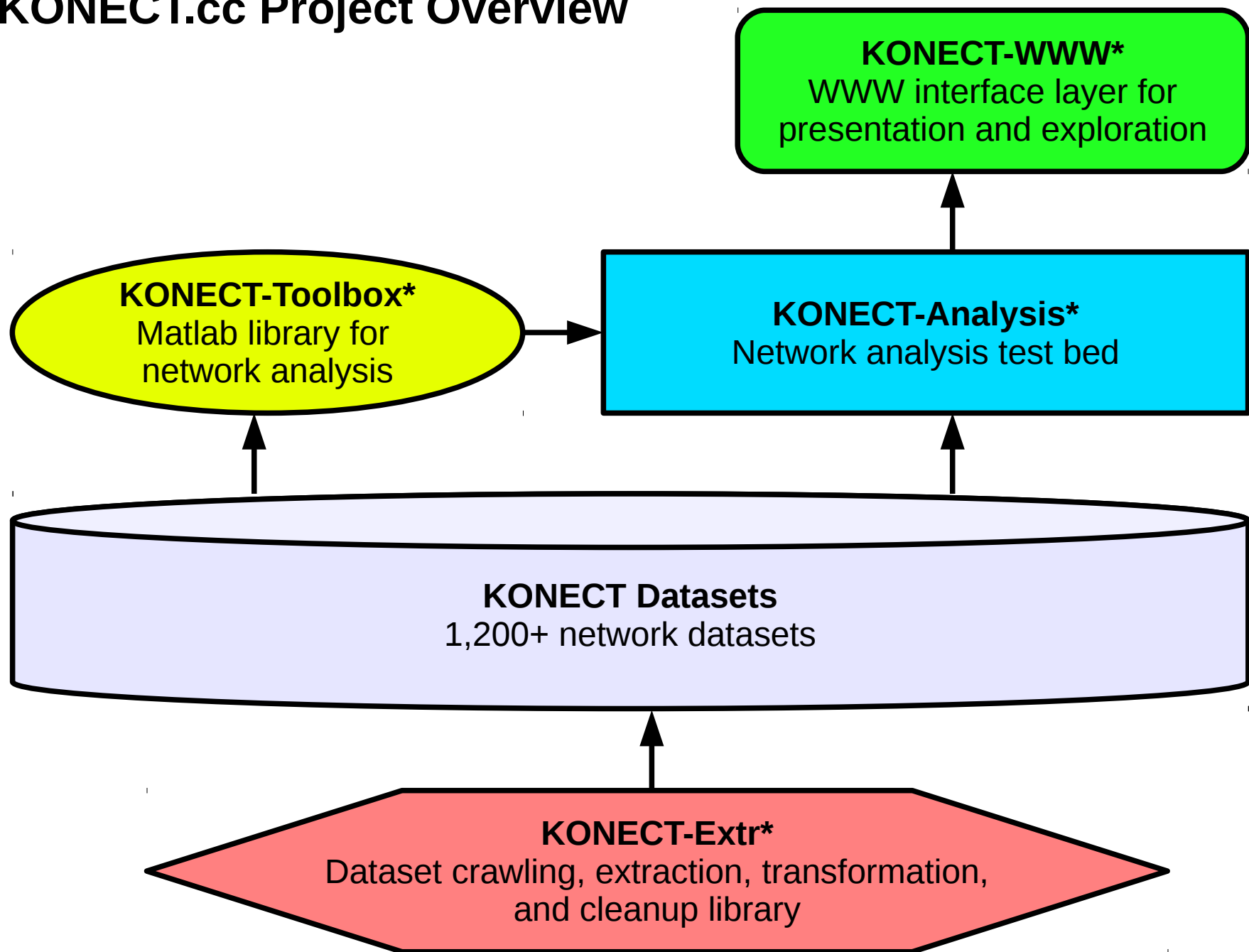
Complex eigenvalues of the asymmetric adjacency matrix

# Download

<http://konect.cc/networks/>

Code	Name	Category	F.	W.	M.	Size	Volume	Avg. degree	Download
<b>Ds</b>	discogs_1style	Features	<b>B</b>	■		244,147	5,255,950	21.56	 
<b>AN</b>	Adjective-noun relationships	Lexical	<b>B</b>	■		194	425	3.94	 
<b>AD</b>	Adrogato	Social	<b>D</b>	+		6,551	51,332	7.64	 
<b>AM</b>	Amazon	Contact	<b>U</b>	■		605,751	3,367,368	4.20	 
<b>AR</b>	Amazon ratings	Ratings	<b>B</b>	+	Ⓢ	3,376,972	5,838,041	2.72	 
<b>AP</b>	arXiv astro-ph	Contact	<b>U</b>	■		37,544	396,160	10.55	 
<b>AC</b>	arXiv cond-mat	Authorship	<b>B</b>	■		38,741	58,595	3.50	 
<b>PH</b>	arXiv hep-ph	Contact	<b>U</b>	■		28,093	12,730,098	453.14	 
<b>PHc</b>	arXiv hep-ph	Reference	<b>D</b>	■		60,388	421,578	6.98	 
<b>THc</b>	arXiv hep-th	Reference	<b>D</b>	■		46,239	352,607	7.31	 
<b>TH</b>	arXiv hep-th	Contact	<b>U</b>	■		22,908	11,209,368	489.32	 
<b>th</b>	arXiv hep-th (KDD Cup)	Reference	<b>D</b>	■		27,770	352,607	12.70	
<b>BAI</b>	Baidu	Reference	<b>D</b>	■		2,141,300	17,794,839	8.31	
<b>BAr</b>	Baidu	Reference	<b>D</b>	■		415,641	3,284,387	7.90	
<b>BS</b>	Berkeley/Stanford	Reference	<b>D</b>	■		1,297,580	7,600,595	5.86	 
<b>Bti</b>	BibSonomy ti	Folksonomy	<b>B</b>	■	Ⓢ	975,963	2,555,060	12.48	
<b>Bui</b>	BibSonomy ui	Folksonomy	<b>B</b>	■	Ⓢ	777,084	2,555,060	440.99	
<b>But</b>	BibSonomy ut	Folksonomy	<b>B</b>	■	Ⓢ	210,467	2,555,060	440.99	
<b>BK</b>	Brightkite	Social	<b>U</b>	■		58,228	214,078	3.68	 
<b>PM</b>	Caenorhabditis elegans	Contact	<b>U</b>	■		453	4,596	10.15	 
<b>IN</b>	CAIDA	Physical	<b>U</b>	■		26,475	106,762	4.03	 
<b>RO</b>	California	Physical	<b>U</b>	■		3,930,412	5,533,214	1.41	 
<b>Sc</b>	Catster	Social	<b>U</b>	■		149,700	5,449,275	36.40	  
<b>Scd</b>	Catster/Dogster	Social	<b>U</b>	■		624,127	15,705,337	25.16	  
<b>CS</b>	CiteSeer	Reference	<b>D</b>	■		723,131	1,764,929	2.44	
<b>Cti</b>	CiteULike ti	Folksonomy	<b>B</b>	■	Ⓢ	685,046	2,411,619	15.74	
<b>Cui</b>	CiteULike ui	Folksonomy	<b>B</b>	■	Ⓢ	754,484	2,411,619	106.18	
<b>Cut</b>	CiteULike ut	Folksonomy	<b>B</b>	■	Ⓢ	175,992	2,411,619	106.18	
<b>CN</b>	Countries	Affiliation	<b>B</b>	■		512,761	557,587	1.09	 
<b>PI</b>	DBLP	Reference	<b>D</b>	■		12,591	49,793	3.95	 

# KONECT.cc Project Overview



\* GitHub package

konect\_dentropy.m  
konect\_diameff.m  
konect\_diammean.m  
konect\_effective\_diameter.m  
konect\_eigl.m  
konect\_eign.m  
konect\_eigskew.m  
konect\_first\_index.m  
konect\_fromto.m  
konect\_gini\_direct.m  
konect\_gini.m  
konect\_hopdistr\_ex.m  
konect\_hopdistr.m  
konect\_imageubu\_complex.m  
konect\_imageubu.m  
konect\_jain.m  
konect\_join.m  
konect\_label\_statistic.m  
konect\_map.m  
konect\_matrix.m  
konect\_mauc.m  
konect\_network\_rank\_abs.m  
konect\_normalize\_additively.m  
konect\_normalized\_entropy.m  
konect\_normalize\_matrix.m  
konect\_normalize\_rows.m  
konect\_order\_dedicom.m  
konect\_own.m





## Handbook of Network Analysis KONECT project

Jérôme Kunegis  
University of Namur, Belgium  
naXys – Namur Center for Complex Systems  
[konect.cc](http://konect.cc)

November 4, 2017

### Abstract

This is the handbook for the KONECT project, a scientific project to archive network datasets, compute systematic network theoretic statistics about them, visualize their properties, and provide corresponding data and Free Software tools to researchers.

<https://github.com/kunegis/konect-handbook>

<http://KONECT.cc/>

# KONECT Datasets in the Wild

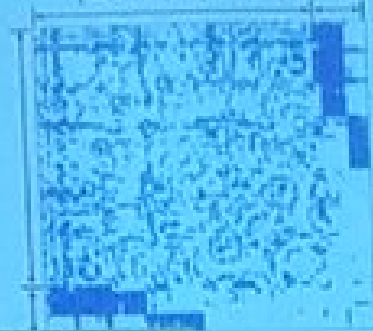
## Structural-Core Pattern: Pattern



### Pattern 3: Structural-Core Pattern

Degeneracy-cores have structural patterns such as core-periphery and communities

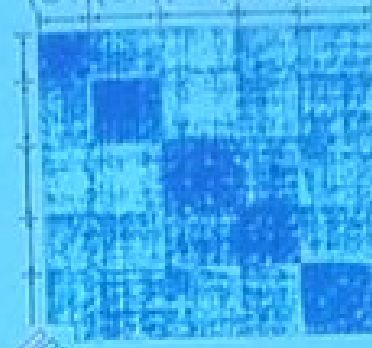
(Core) (Periphery)



Handwritten text: *Handwritten text*

Core-periphery

(C1)(C2)(C3)(C4)(C5)



Communities

Nodes in degeneracy-cores are not homogeneous

Scholar

About 203 results (0.02 sec)



All citations

## Konect: the koblenz network collection

Articles

Search within citing articles

<https://scholar.google.com/scholar?cites=7174338004474749050>

Case law

## Recent advances in graph partitioning

My library

[A Buluç](#), [H Meyerhenke](#), [I Safro](#), [P Sanders...](#) - *Algorithm ...*, 2016 - Springer

Cited by 76 [Related articles](#) [All 13 versions](#) [Cite](#) [Save](#) [More](#)

Any time

## Estimating clustering coefficients and size of social networks via random walk

Since 2017

[SJ Hardiman](#), [L Katzir](#) - ... of the 22nd international conference on World ..., 2013 - dl.acm.org

Since 2016

Abstract Online social networks have become a major force in today's society and economy.

Since 2013

The largest of today's social networks may have hundreds of millions to more than a billion users. Such networks are too large to be downloaded or stored locally, even if terms of use

Custom range...

Cited by 44 [Related articles](#) [All 10 versions](#) [Cite](#) [Save](#) [More](#)

Sort by relevance

## BFS and coloring-based parallel algorithms for strongly connected components and related problems

Sort by date

[GM Słota](#), [S Rajamanickam...](#) - *Parallel and Distributed ...*, 2014 - ieeexplore.ieee.org

include citations

Abstract: Finding the strongly connected components (SCCs) of a directed graph is a fundamental graph-theoretic problem. Tarjan's algorithm is an efficient serial algorithm to find SCCs, but relies on the hard-to-parallelize depth-first search (DFS). We observe that  
Cited by 39 [Related articles](#) [All 13 versions](#) [Cite](#) [Save](#) [More](#)

Create alert

## Pulp: Scalable multi-objective multi-constraint partitioning for small-world networks

[GM Słota](#), [K Madduri...](#) - *Big Data (Big Data)*, 2014 ..., 2014 - ieeexplore.ieee.org

Abstract: We present PuLP, a parallel and memory-efficient graph partitioning method specifically designed to partition low-diameter networks with skewed degree distributions. Graph partitioning is an important Big Data problem because it impacts the execution time  
Cited by 19 [Related articles](#) [All 11 versions](#) [Cite](#) [Save](#) [More](#)

# Learning Goals

- Everything is a network
- Take datasets
- Give datasets
- Increase the number of datasets you use
- Compare datasets
- Transfer knowledge across datasets
- Automate everything

<http://konect.cc/>

<https://github.com/kunegis/konect-handbook>

<https://github.com/kunegis/stu>

We are on Twitter: [@KONECTproject](https://twitter.com/KONECTproject)



# Outline

- 1 Introduction
- 2 Statistics
- 3 Plots
- 4 Studies
  - (a) Diversity
  - (b) Graph generator
  - (c) Power laws
  - (d) Preferential attachment
  - (e) Conflict
- 5 Automate everything : Stu



# One More Note

This is a tutorial, not a paper talk – questions, comments and direct participation are welcome :)