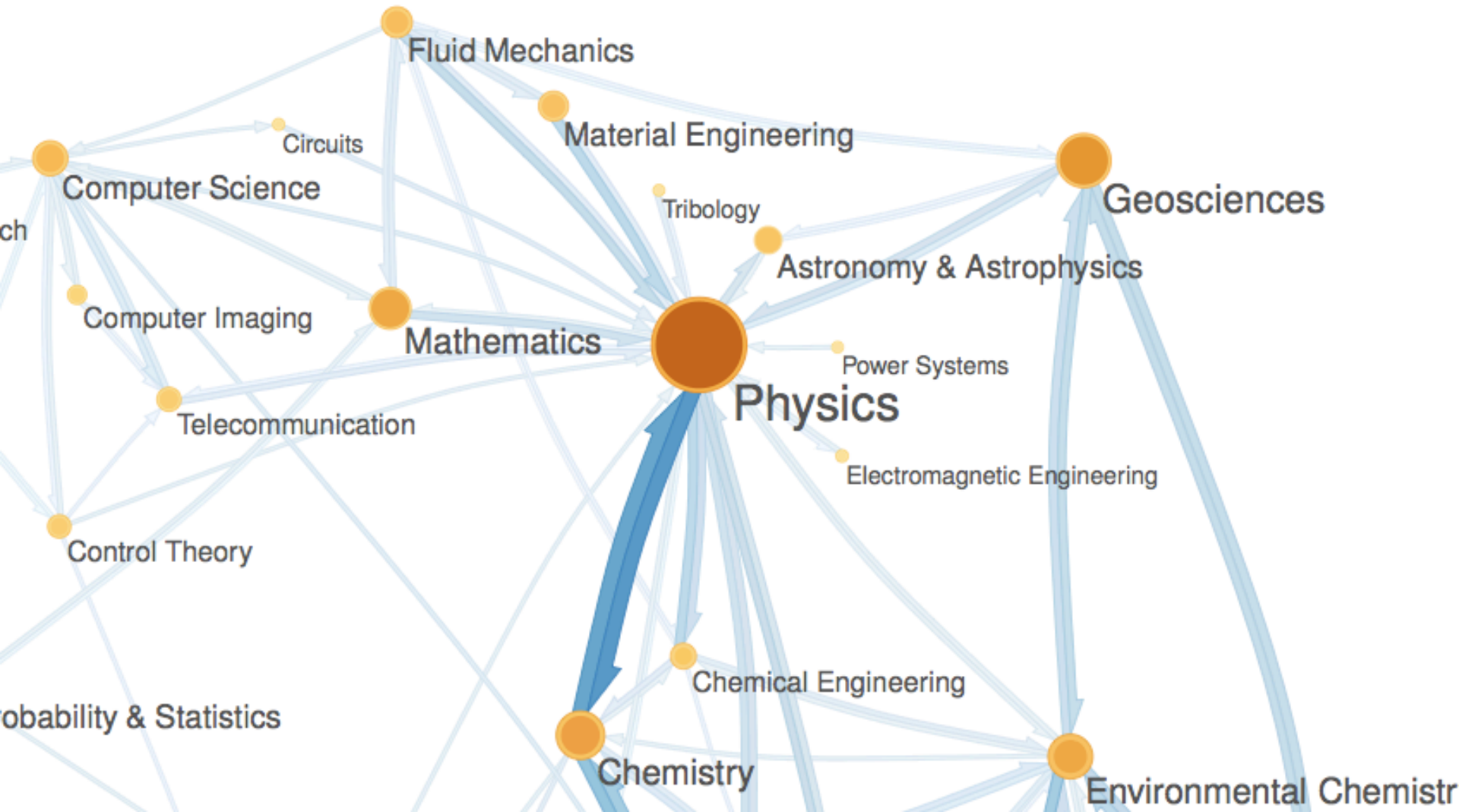
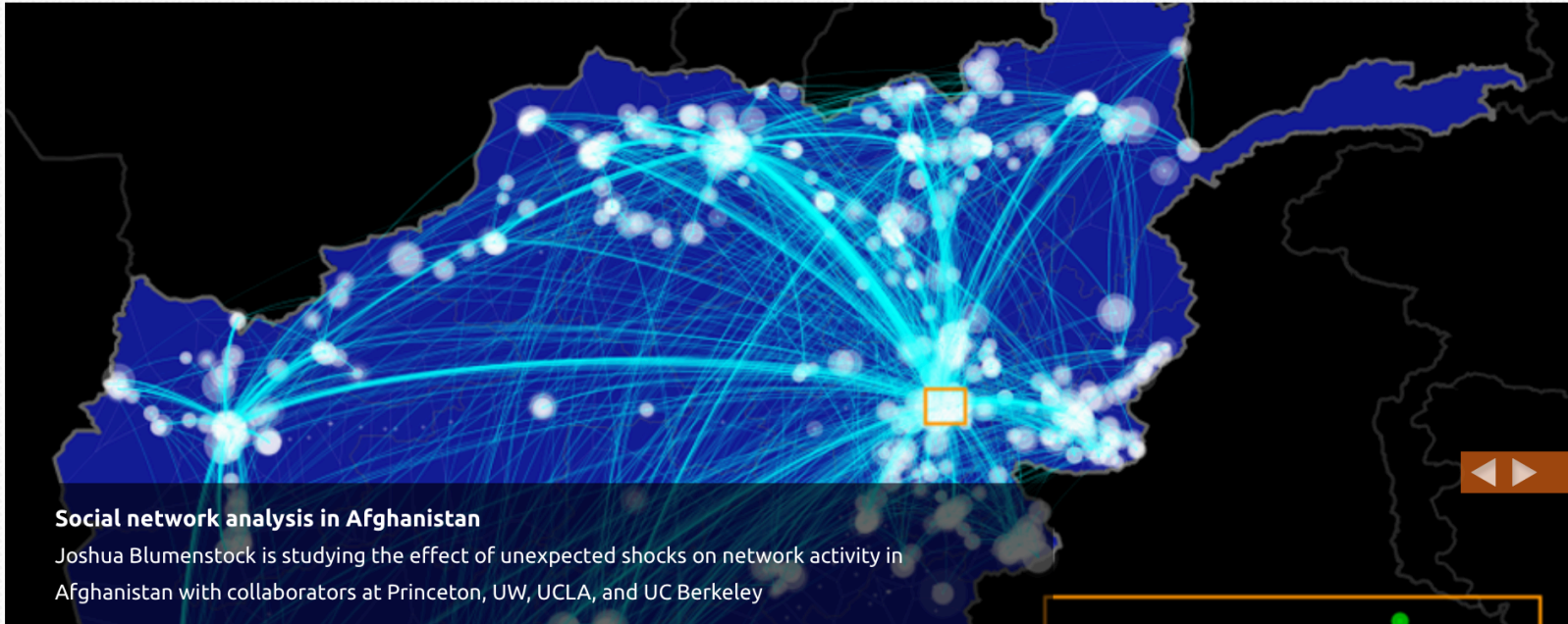


Hierarchically clustering time-directed graphs and the effects of teleportation and memory

Jevin West, Information School, University of Washington

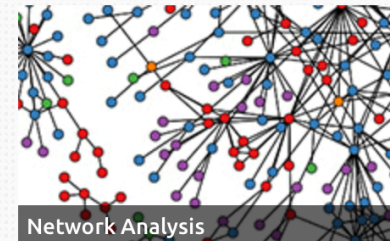




Social network analysis in Afghanistan

Joshua Blumenstock is studying the effect of unexpected shocks on network activity in Afghanistan with collaborators at Princeton, UW, UCLA, and UC Berkeley

Research Focus Areas



News and Updates

28 Blumenstock at Population Association of America

What we do

The DataLab is the nexus for research on Data Science and Analytics at the UW iSchool. We study **large-scale, heterogeneous human data** in an

Network Clustering

Graph Partitioning

Community Detection

Block Models

Module Detection



○ ○ ○ No one size fits all ○ ○

- No canonical solution or one generalizable method for all data and all problems (i.e. there is no method that works best on all networks in all situations)
- Need to know the context for why the user is interested in clustering
- We don't even have a definition of a community
- Umbrella term for many facets



No one size fits all

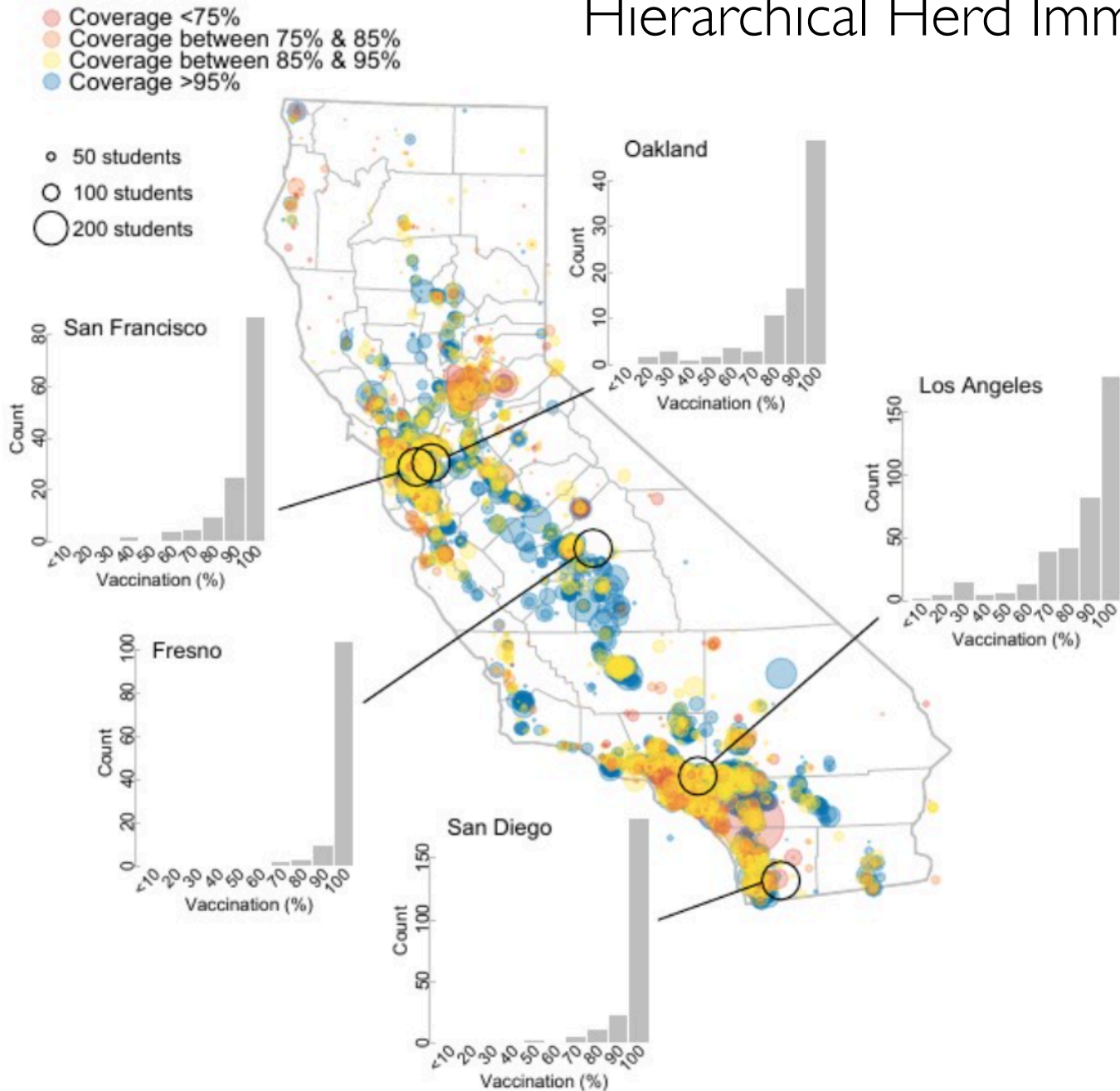
Cut-based: community detection as minimization of some form of constraint violation

Data clustering: community detection framed as a discretized analogue of data clustering, in which densely knit groups of nodes are to be found

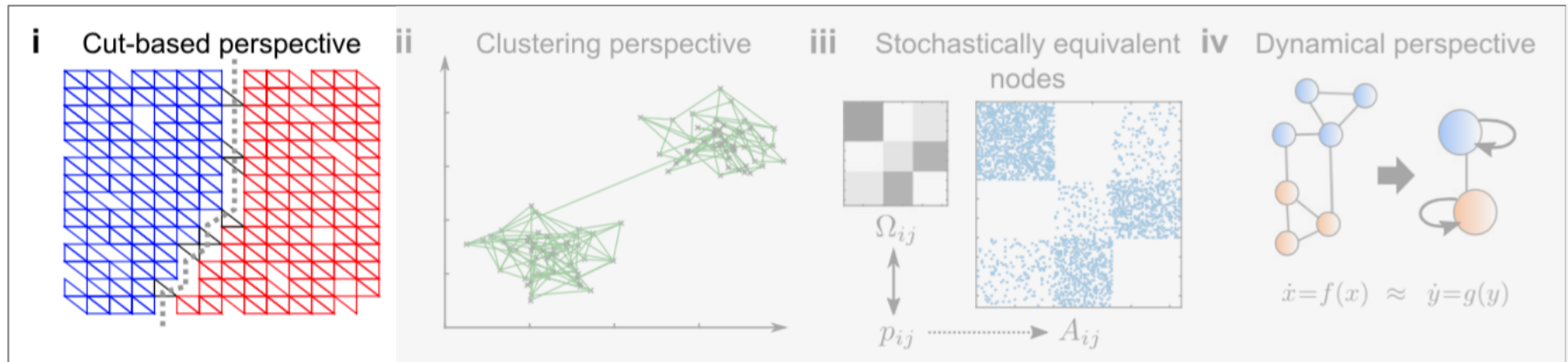
Stochastic equivalence: community detection aiming to identify structurally equivalent nodes in a network, leading to notions such as stochastic block models

Dynamics perspective: community detection looking for simplified descriptions of the dynamical flows occurring on the network, that is, some form of dynamical model reduction

Hierarchical Herd Immunity



Community Detection Perspectives



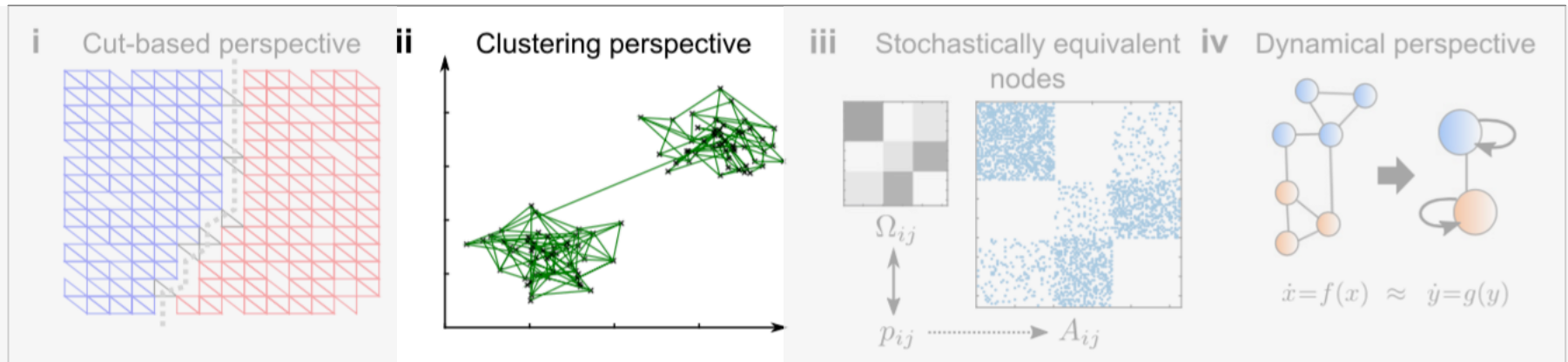
Circuit layout
 Minimizing cuts
 Load balancing
 Eigenvectors
 Spectral methods
 Image segmentation

Data Clustering
 Maximizing node density
 unknown k , unbalanced
 Conductance
 Local, global
 Modularity

Social Networks
 Connectivity Profiles
 Stochastic equivalence
 SBMs, LFR
 p-values, hypothesis testing
 Bipartite treatment
 Predict missing links

System behavior, processes
 Non-adjacency focused
 Airline network
 Markovian diffusion process
 Undirected, Directed
 InfoMap

Community Detection Perspectives



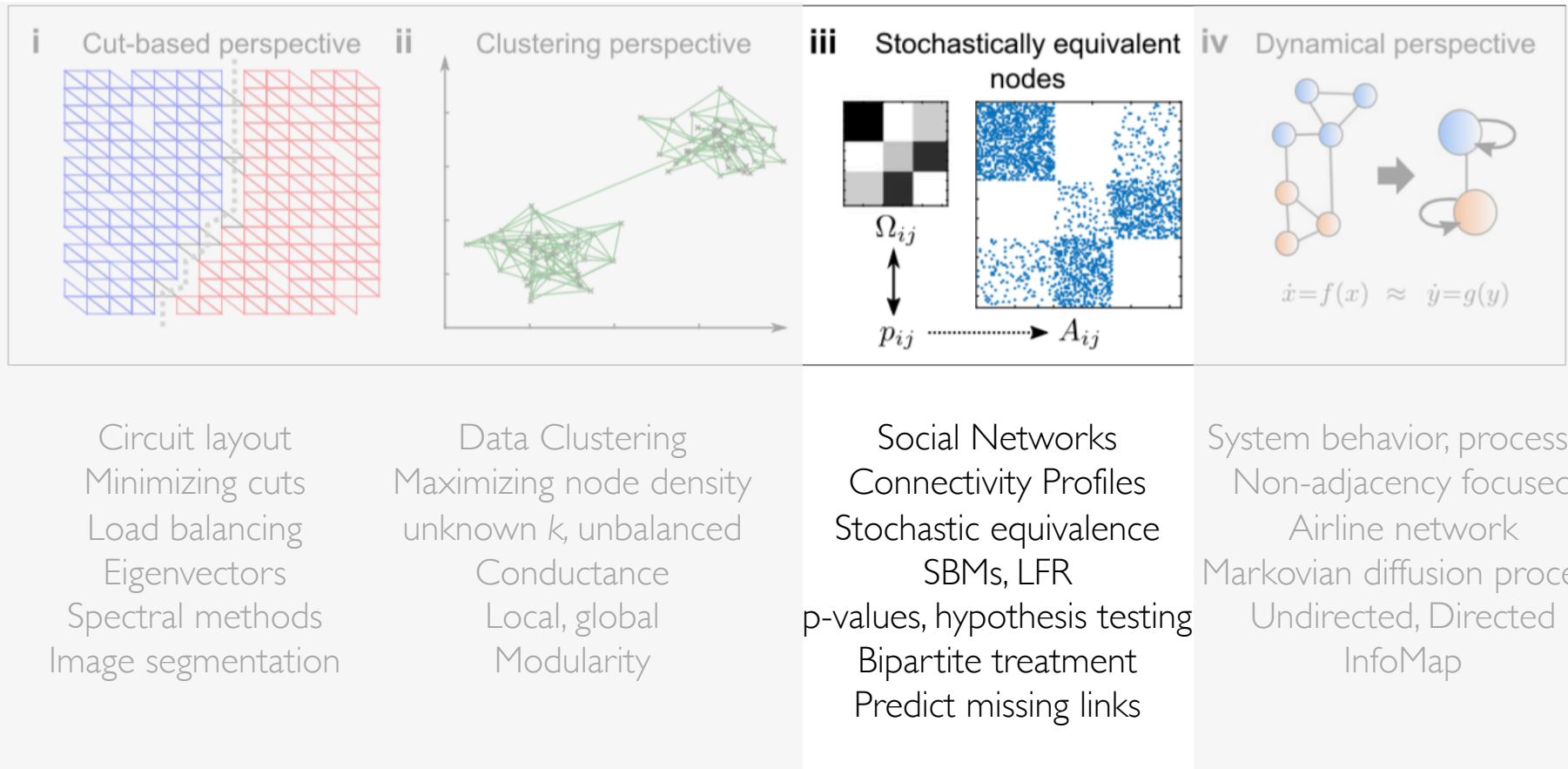
Circuit layout
 Minimizing cuts
 Load balancing
 Eigenvectors
 Spectral methods
 Image segmentation

Data Clustering
 Maximizing node density
 unknown k , unbalanced
 Conductance
 Local, global
 Modularity

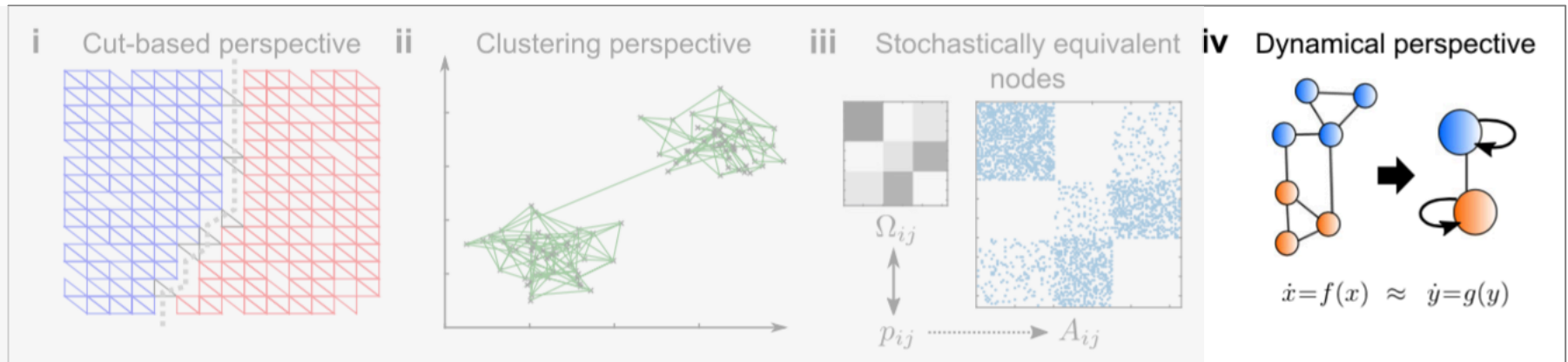
Social Networks
 Connectivity Profiles
 Stochastic equivalence
 SBMs, LFR
 p-values, hypothesis testing
 Bipartite treatment
 Predict missing links

System behavior, processes
 Non-adjacency focused
 Airline network
 Markovian diffusion process
 Undirected, Directed
 InfoMap

Community Detection Perspectives



Community Detection Perspectives



Circuit layout
 Minimizing cuts
 Load balancing
 Eigenvectors
 Spectral methods
 Image segmentation

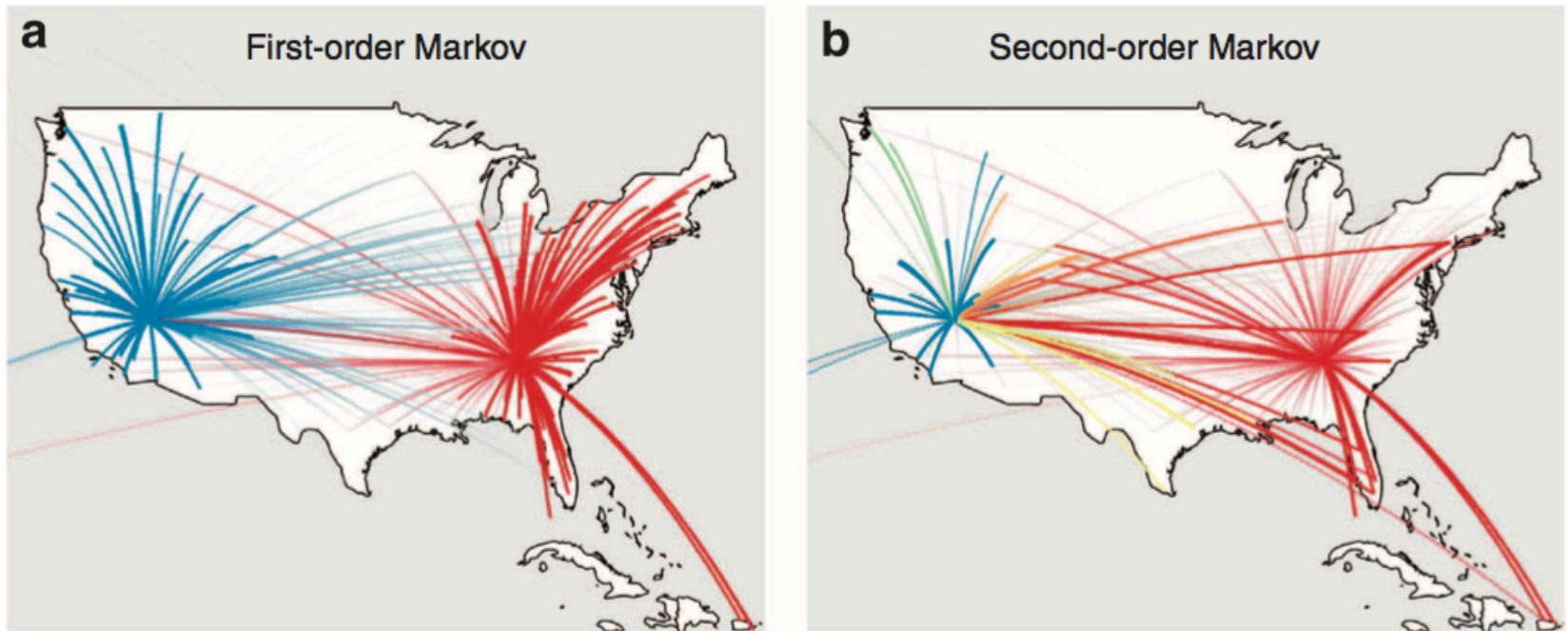
Data Clustering
 Maximizing node density
 unknown k , unbalanced
 Conductance
 Local, global
 Modularity

Social Networks
 Connectivity Profiles
 Stochastic equivalence
 SBMs, LFR
 p-values, hypothesis testing
 Bipartite treatment
 Predict missing links

System behavior, processes
 Non-adjacency focused
 Airline network
 Markovian diffusion process
 Undirected, Directed
 InfoMap



Higher Resolution Maps



Rosvall et al. (2014) Memory in network flows and its effects on spreading dynamics and community detection. *Nature Communications*

In the spirit of clustering context...



The Scholarly Graph



THOMSON REUTERS

PatentVector™



PNAS



dblp
computer science bibliography





The Scholarly Graph



Tens of millions articles, patents, books



Billions of citation links

Years: 1600 – 2016



1. Mapping Knowledge Domains



2. Science of Science

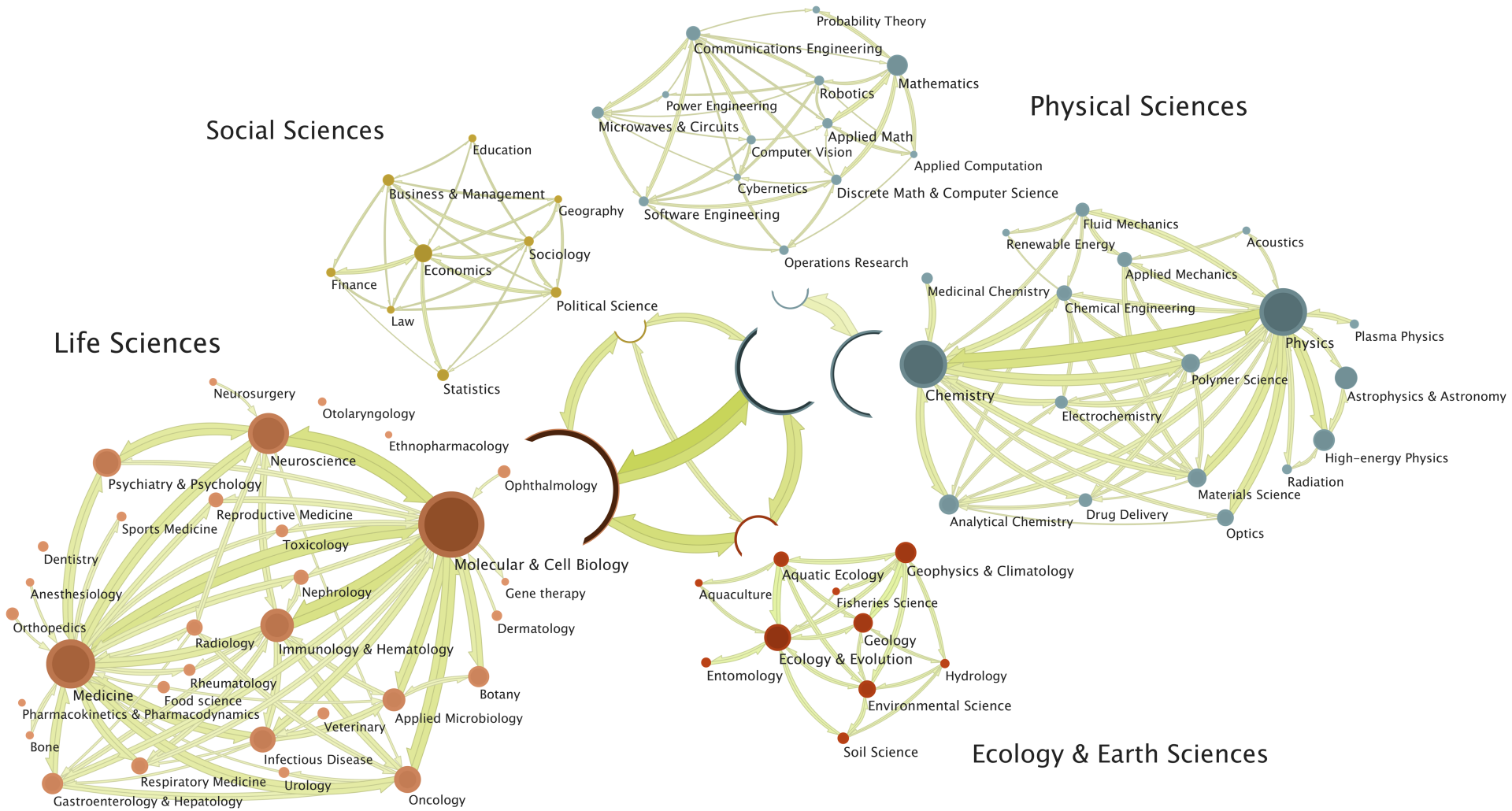
3. Hierarchical Navigation



4. Recommendation

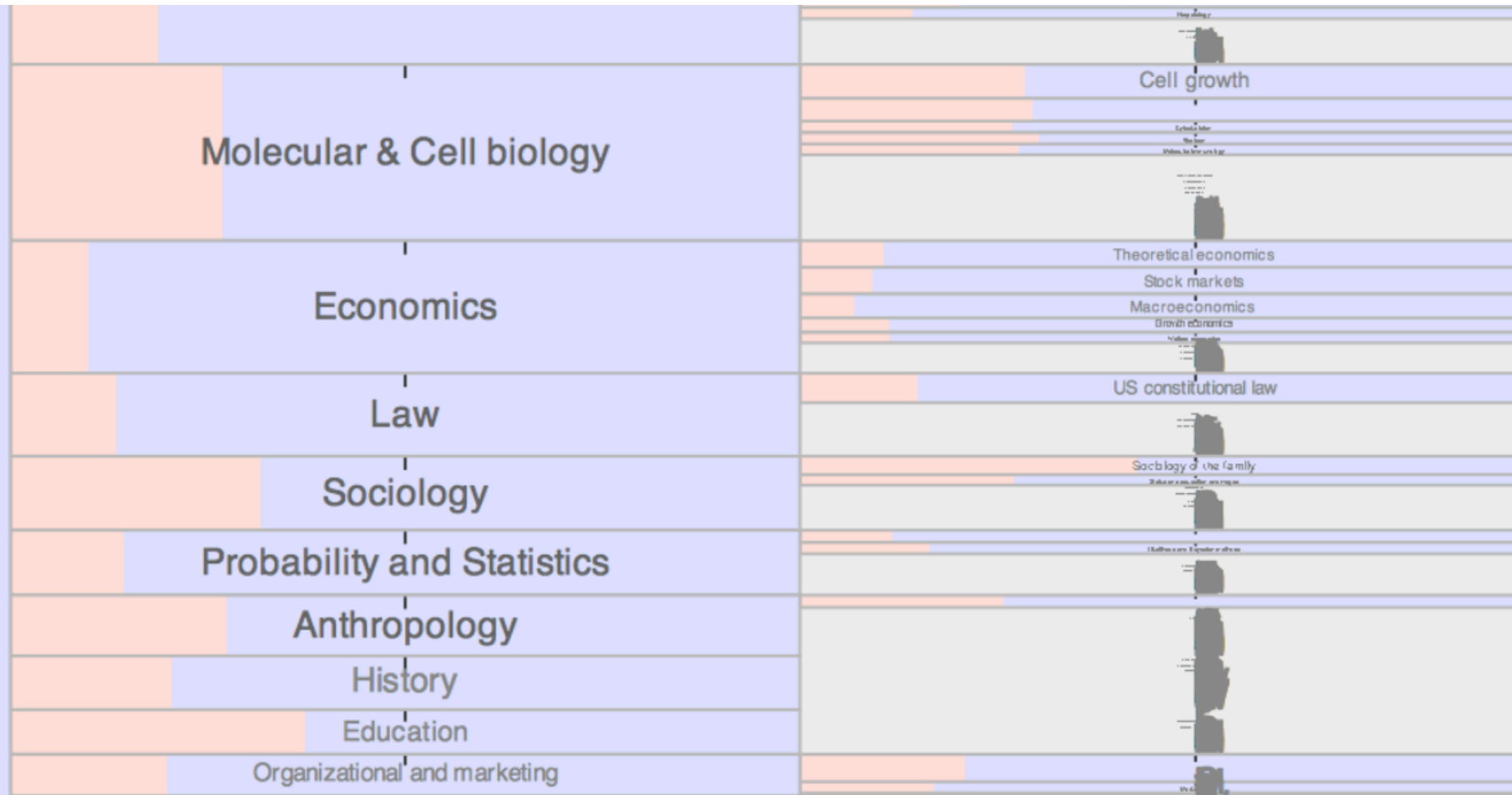


Mapping Knowledge Domains



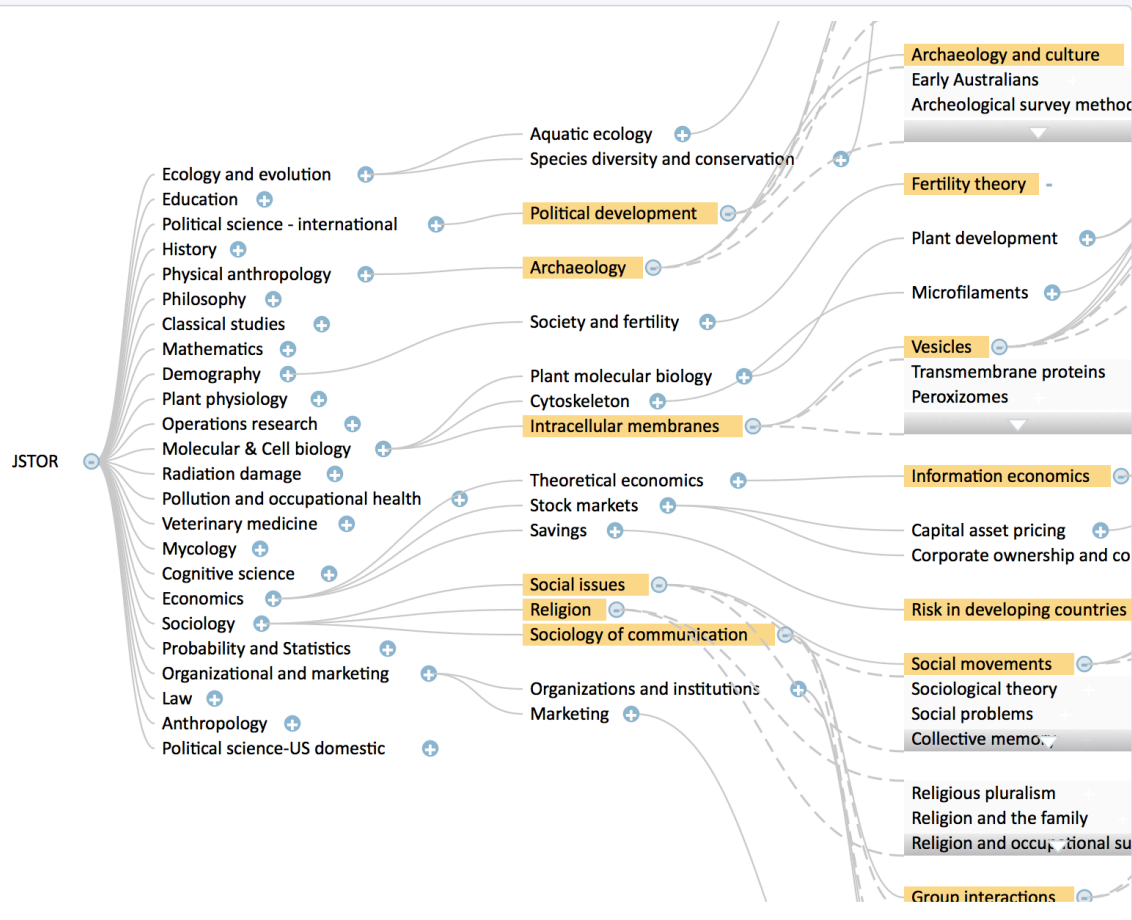
Rosvall, Martin, and Carl T. Bergstrom. "Multilevel compression of random walks on networks reveals hierarchical organization in large integrated systems." *PLoS one* 6.4 (2011): e18209.

The Role of Gender in Science



West, J.D. (2012) The Role of Gender in Scholarly Authorship. *PLoS One*

Hierarchical Navigation



Top Papers

Sort by Year (newest) ▾

[Using Siting Algorithms in the Design of Marine Reserve Networks](#)

Heather Leslie - *Ecological Applications* (2003)

[Mechanism of Filopodia Initiation by Reorganization of a Dendritic Network](#)

Tatyana Svitkina - *The Journal of Cell Biology* (2003)

[Network Structure and Knowledge Transfer: The Effects of Cohesion and Range](#)

Ray Reagans - *Administrative Science Quarterly* (2003)

[A General Model for Designing Networks of Marine Reserves](#)

Eric Sala - *Science* (2002)

[The Density of Social Networks and Fertility Decisions: Evidence from South Nyanza District, Kenya](#)

Hans-Peter Kohler - *Demography* (2001)

[A New Dynamain-Like Protein, ADL6, Is Involved in Trafficking from the trans-Golgi Network to the Central Vacuole in Arabidopsis](#)

Jing Bo Jin - *The Plant Cell* (2001)

[Comparing Sequenced Segments of the Tomato and Arabidopsis Genomes: Large-Scale Duplication Followed by Selective Gene Loss Creates a Network of Synteny](#)

Hsin-Mei Ku - *Proceedings of the National Academy of Sciences of the United States of America* (2000)

[A Noncooperative Model of Network Formation](#)

Venkatesh Bala - *Econometrica* (2000)

Find Papers

by title

by field

by author

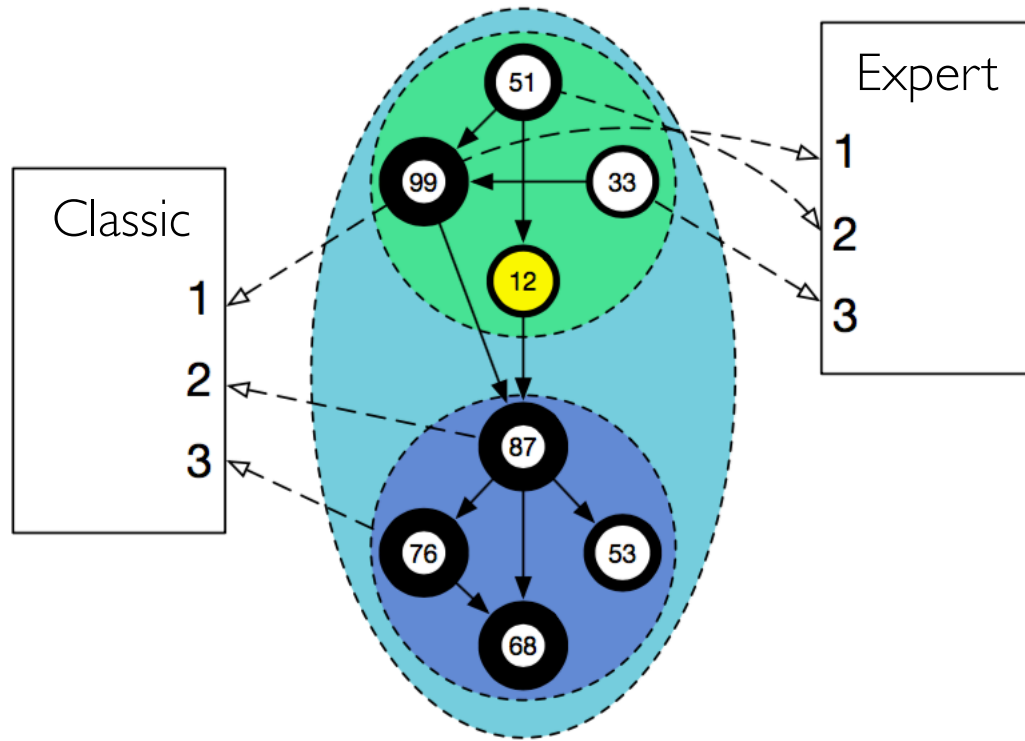
by journal

Active Queries:

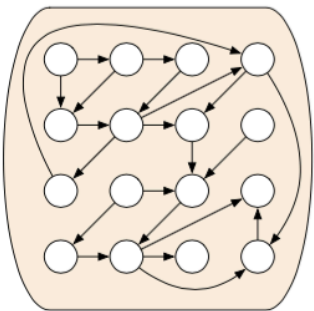
[clear all](#)

keyword: **network**

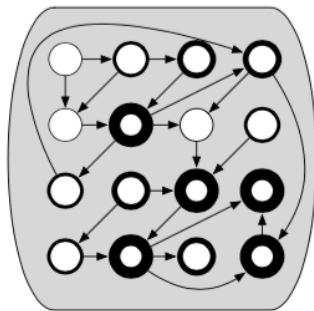
Recommendation



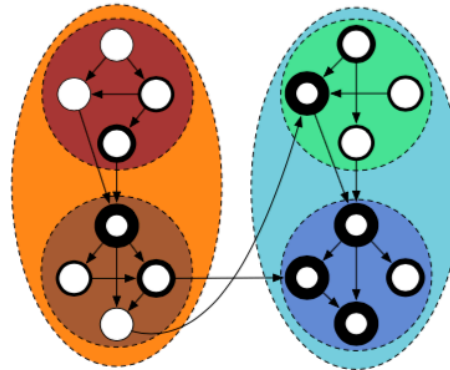
Assemble



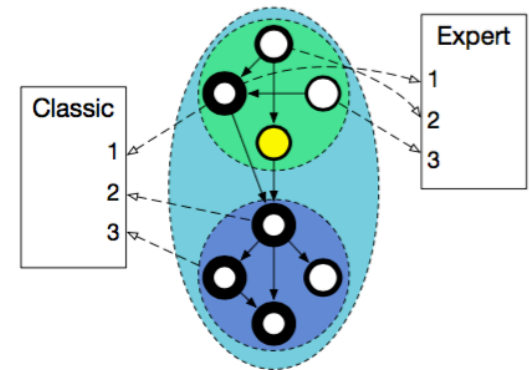
Rank



Cluster

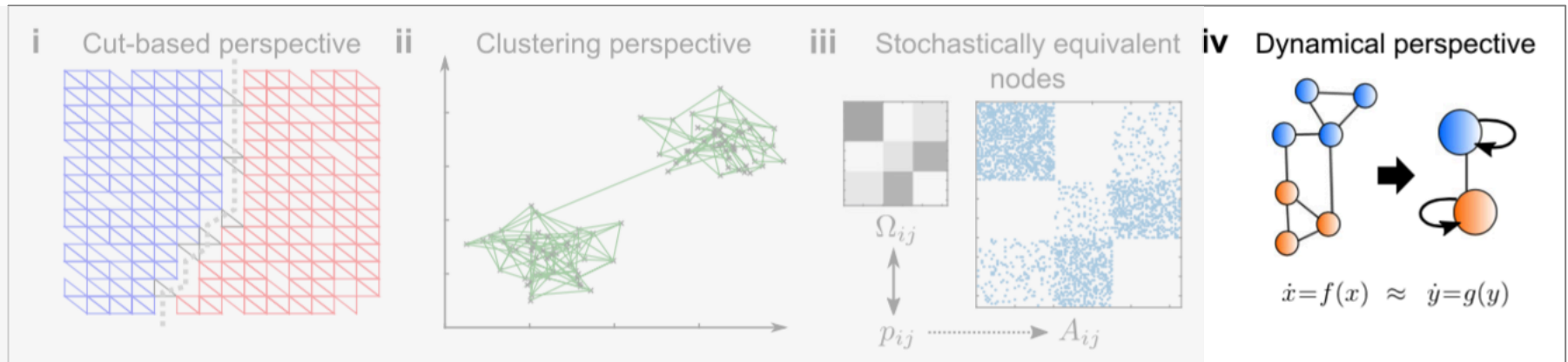


Recommend



West, Wesley-Smith, Bergstrom (2016) A recommendation system based on hierarchical clustering of an article-level citation network. *IEEE, Transactions on Big Data* (in press)

Community Detection Perspectives



Circuit layout
 Minimizing cuts
 Load balancing
 Eigenvectors
 Spectral methods
 Image segmentation

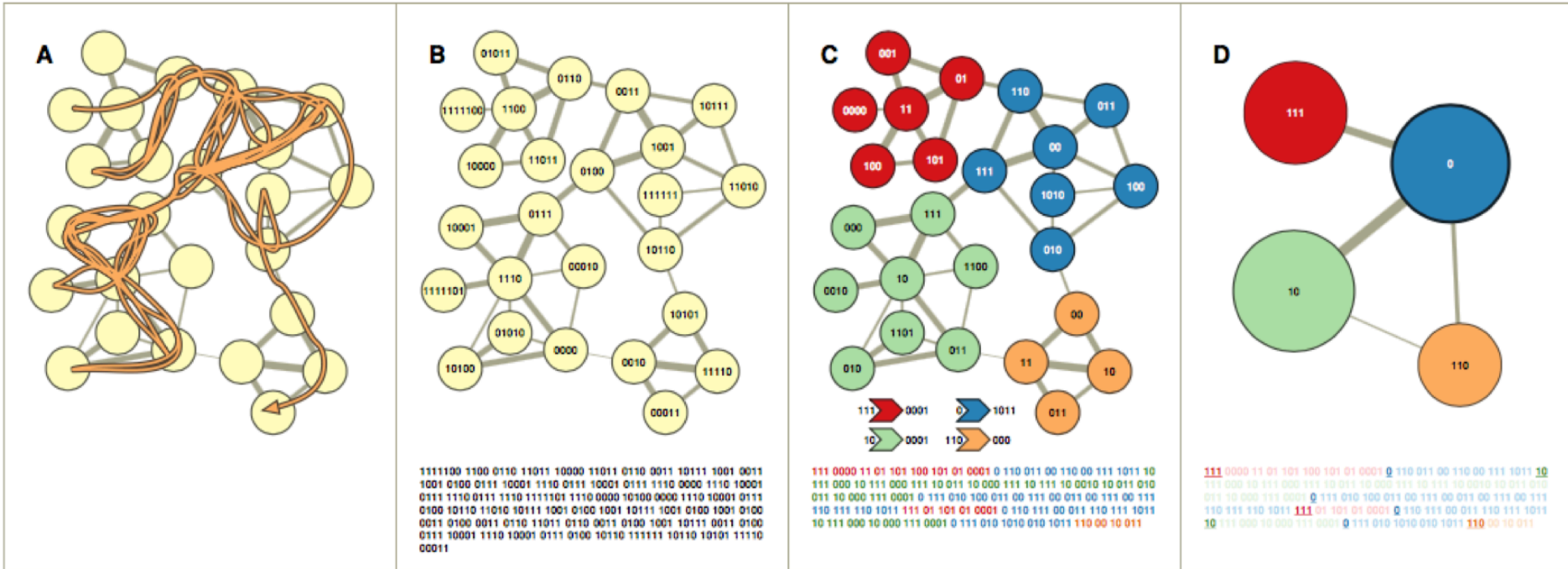
Data Clustering
 Maximizing node density
 unknown k , unbalanced
 Conductance
 Local, global
 Modularity

Social Networks
 Connectivity Profiles
 Stochastic equivalence
 SBMs, LFR
 p-values, hypothesis testing
 Bipartite treatment
 Predict missing links

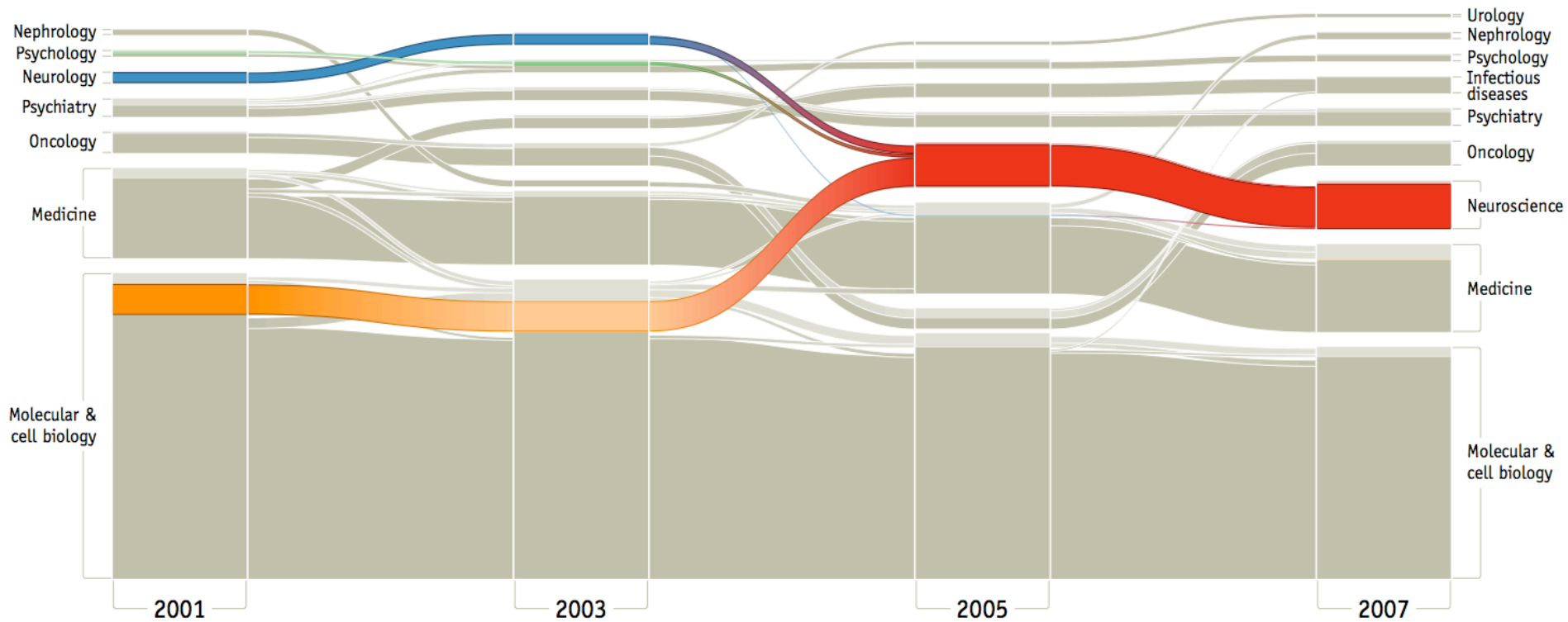
System behavior, processes
 Non-adjacency focused
 Airline network
 Markovian diffusion process
 Undirected, Directed
 InfoMap



Finding regularities in citation networks



The Emergence of Neuroscience



Data

Compressing \longleftrightarrow Finding patterns

If we can find a good code for describing flow on a network, we will have solved the dual problem of finding the important structures with respect to that flow.

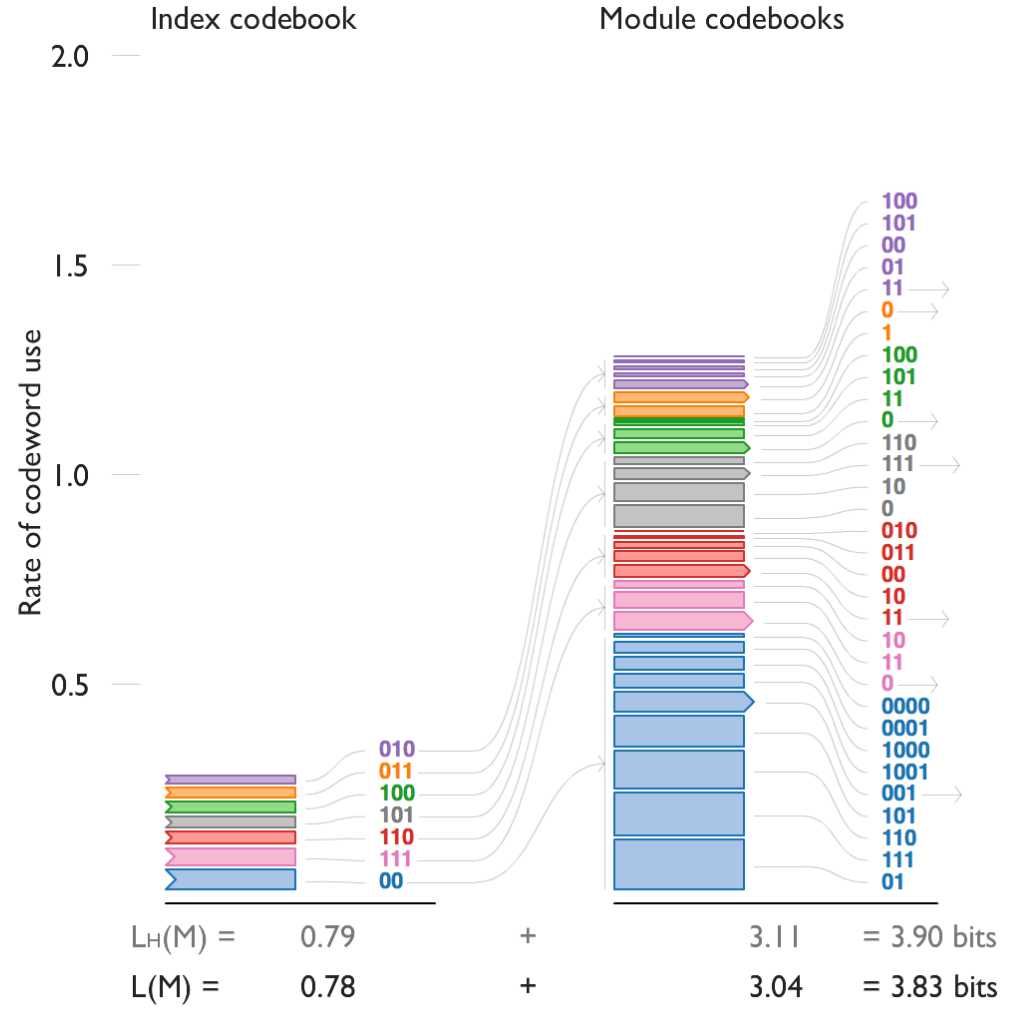
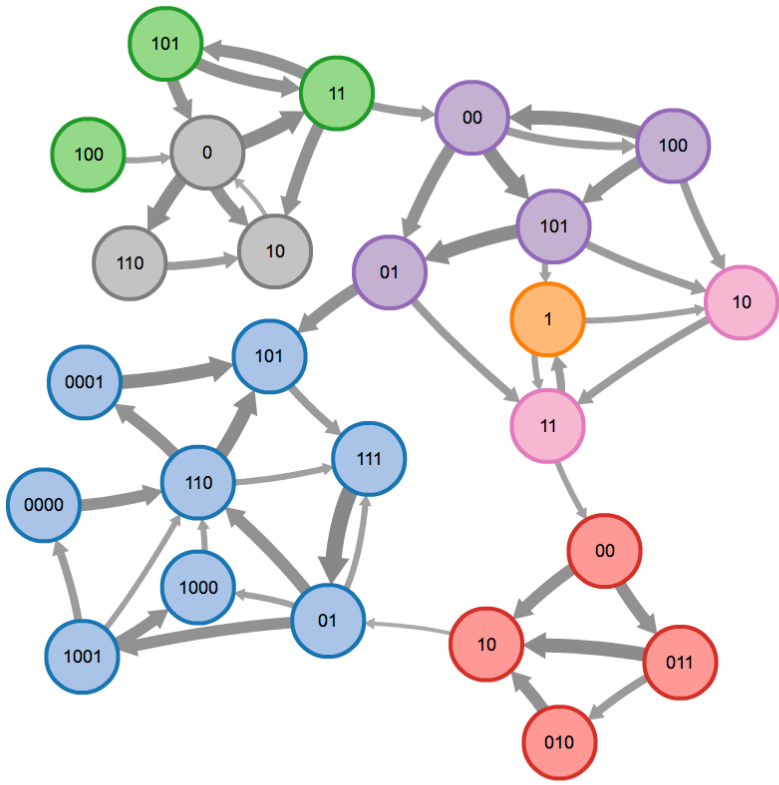
The map equation

frequency of inter-module movements

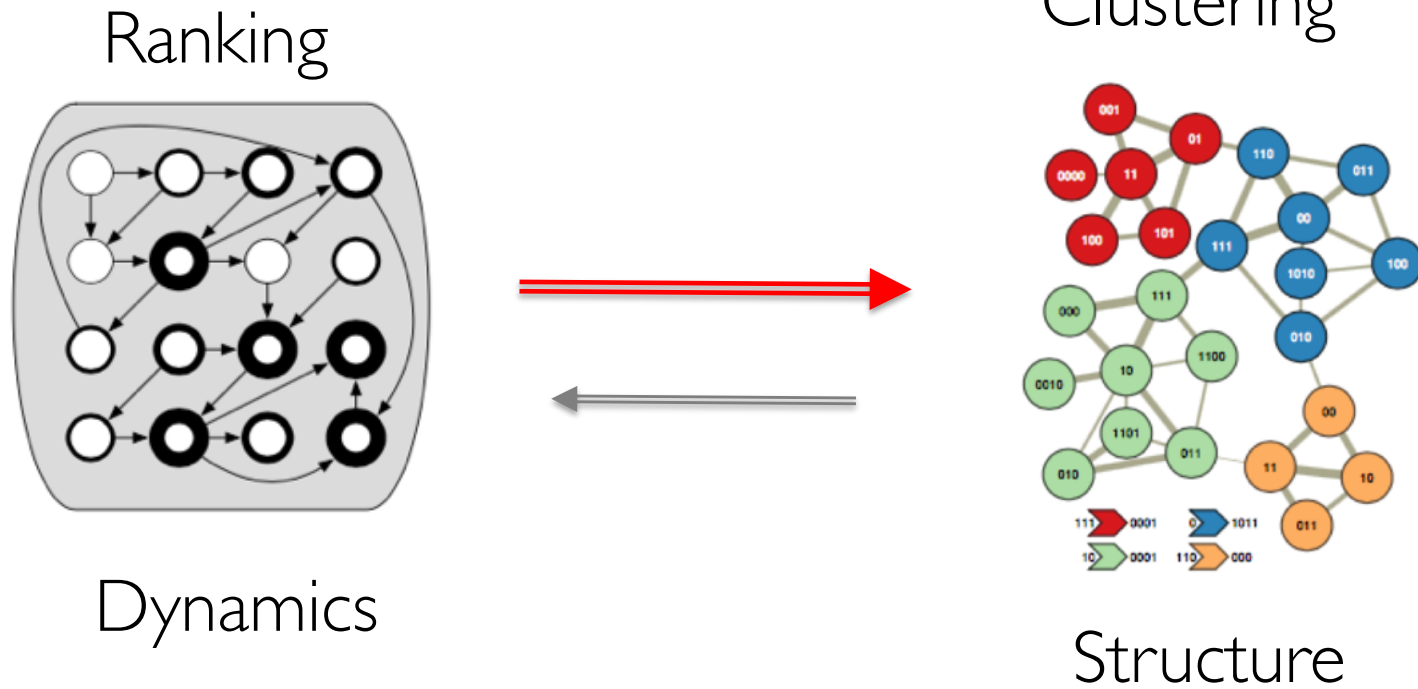
$$L(M) = q_{\curvearrowright} H(Q) + \sum_{i=1}^m p_{\circlearrowleft}^i H(\mathcal{P}^i)$$

code length of module names

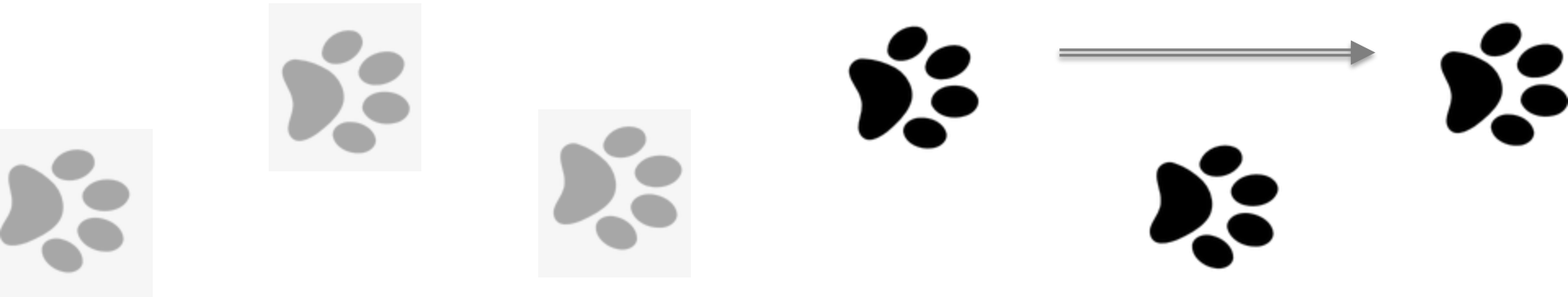
code length of node names in module i



The relationship between ranking and clustering

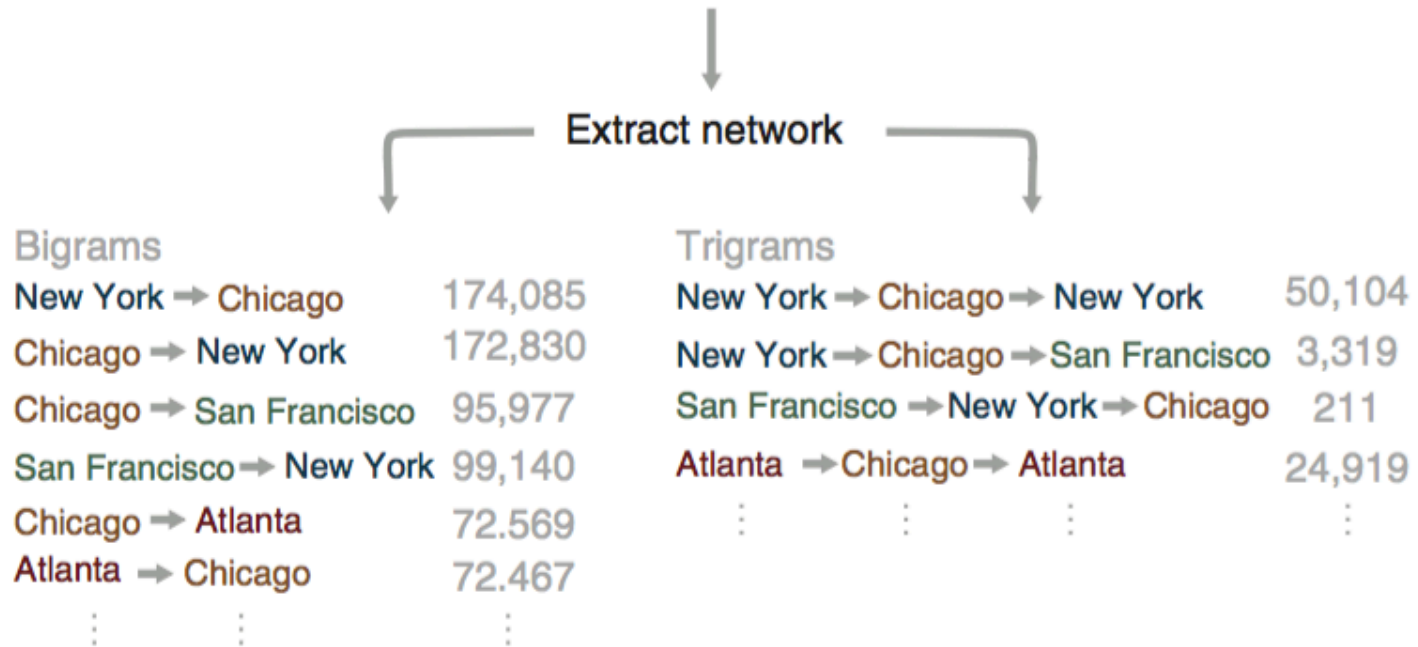


Step Length, Teleportation and Memory



..and their effects on ranking and clustering

Memory: capturing higher order dynamics



Memory: capturing higher order dynamics

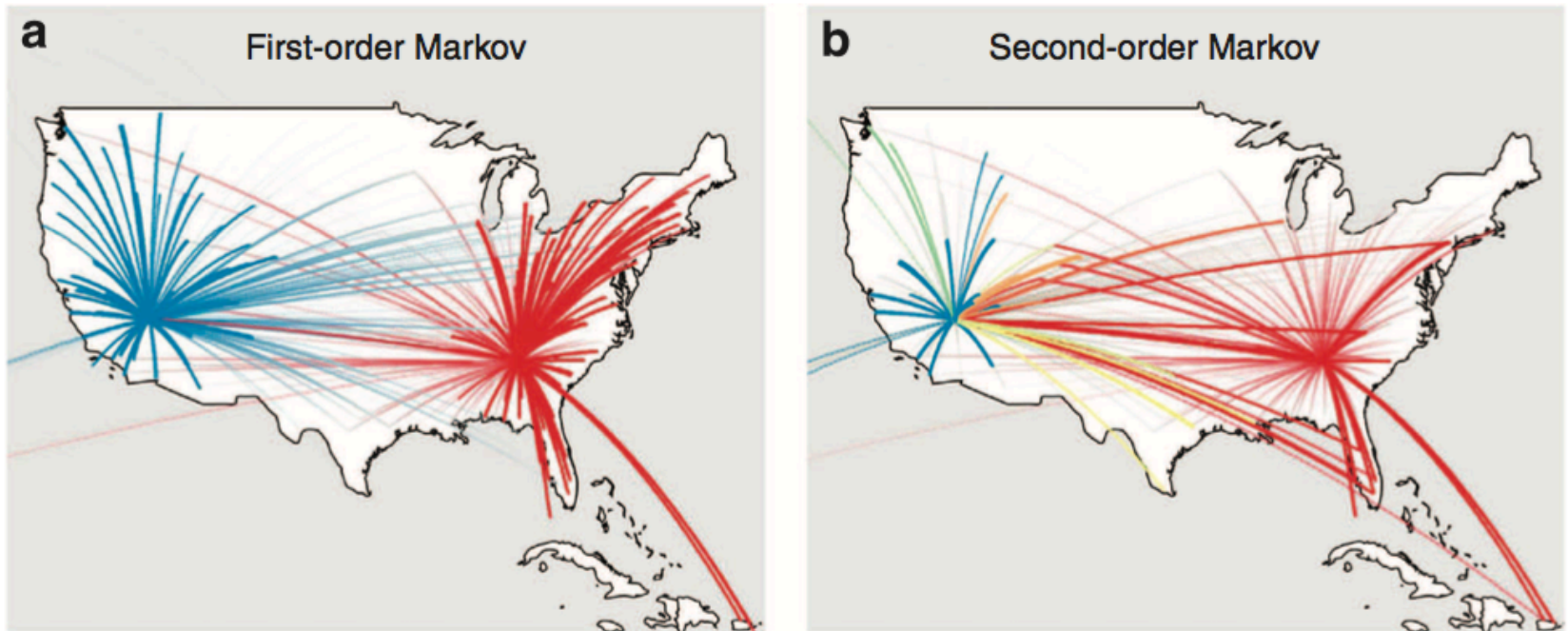
First-order Markov



Second-order Markov

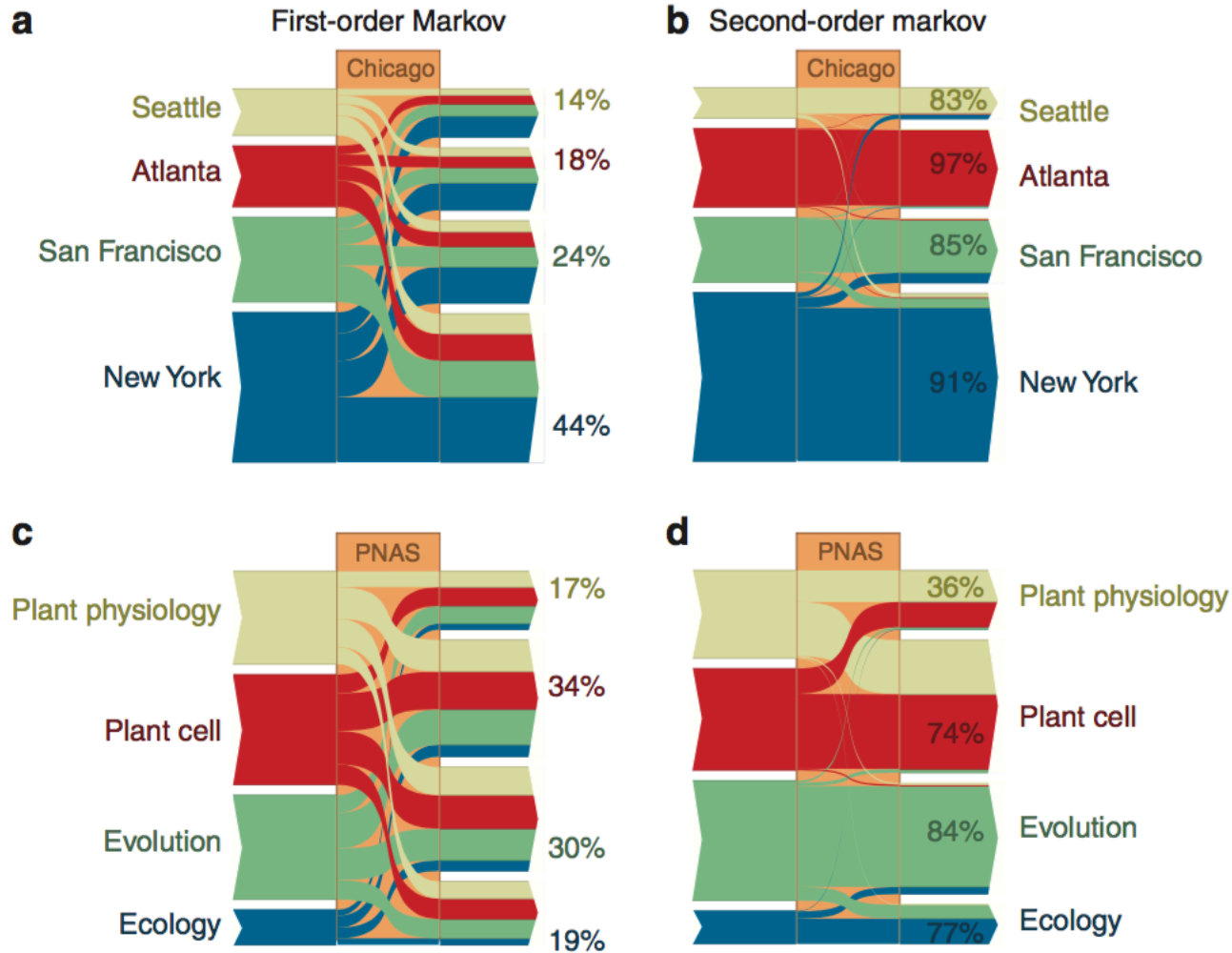


Higher Resolution Maps



Rosvall et al. (2014) Memory in network flows and its effects on spreading dynamics and community detection. *Nature Communications*

Higher Order Dynamics

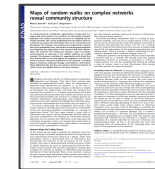
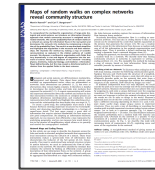


Rosvall et al. (2014) Memory in network flows and its effects on spreading dynamics and community detection. *Nature Communications*

Citation Networks Types



Journal-Level Networks
(Memory)



Article-level Networks

Time-Directed (Acyclic) Graphs

PageRank Variants (EigenFactor)

$$P = \alpha H + (1 - \alpha) a.e^T$$

Matrix representing the random walk over citations

Probability of not teleporting

Cross-citation Matrix dictating the structure of the citation network

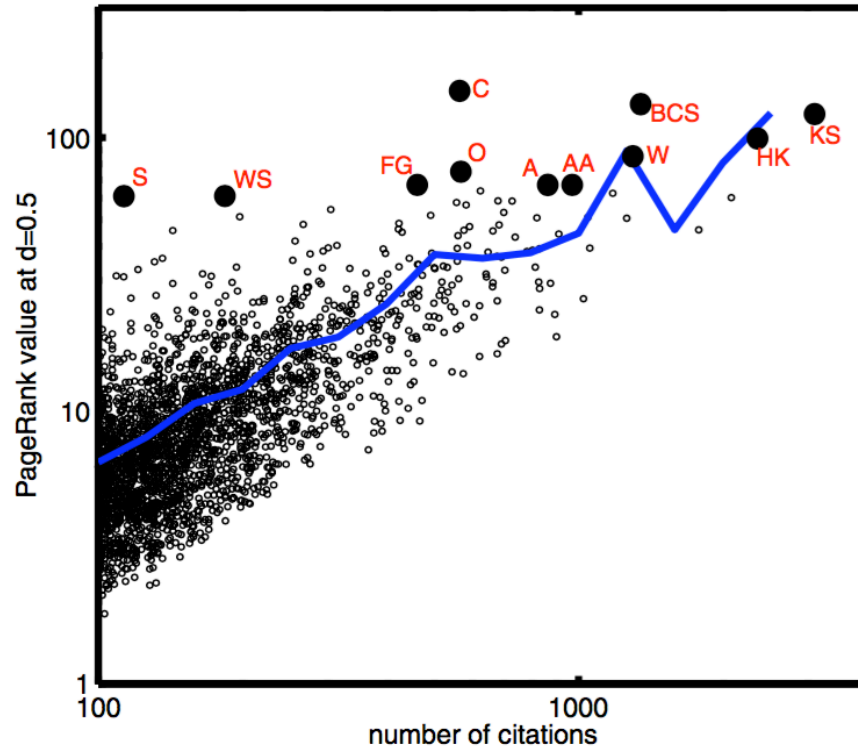
Probability of teleporting to completely new journal weighted by the number of articles in that journal

$$EF = 100 \frac{H \pi}{\sum_i [H \pi]_i}$$

Leading eigenvector of the random walk matrix P .

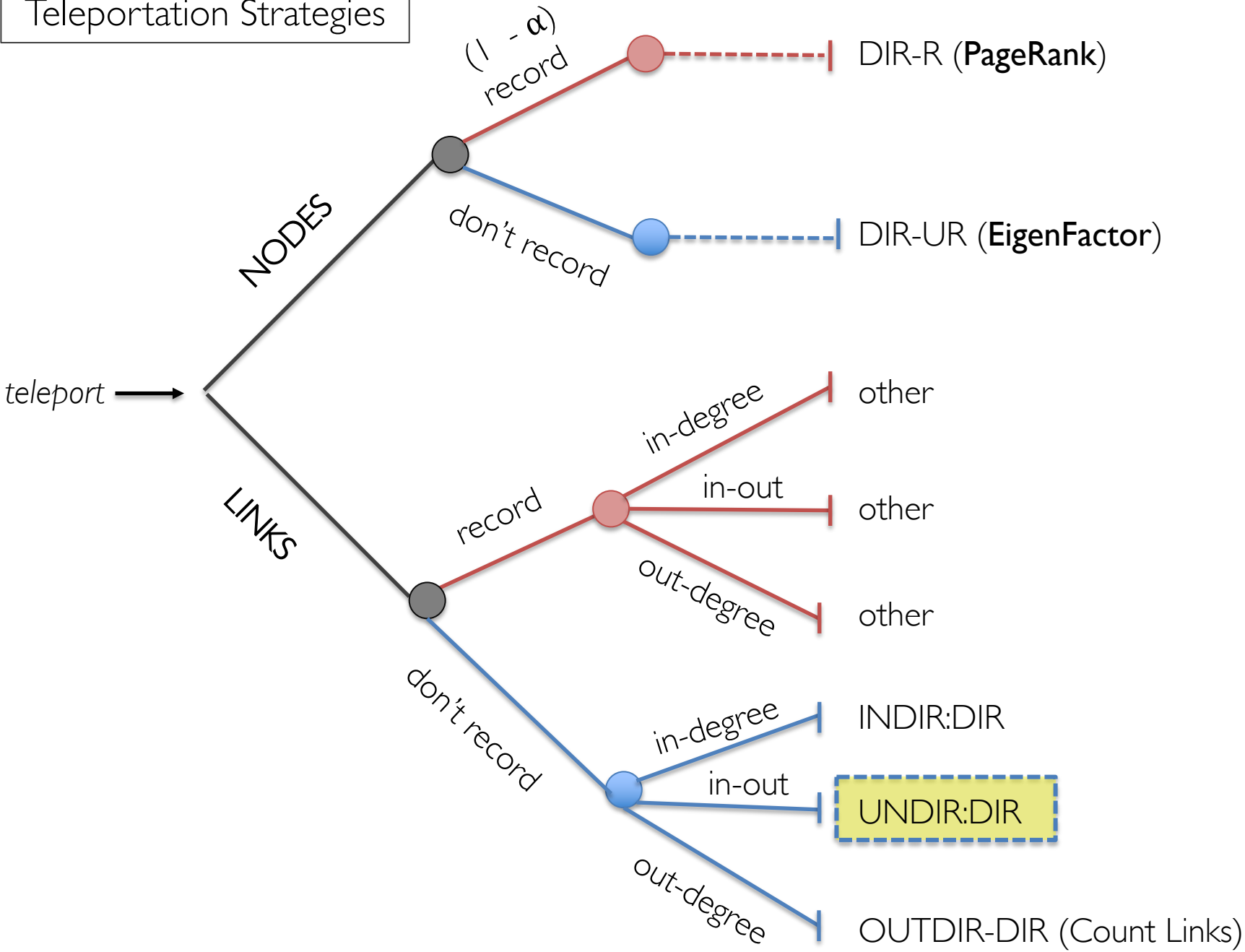
Normalization

PageRank Pitfalls



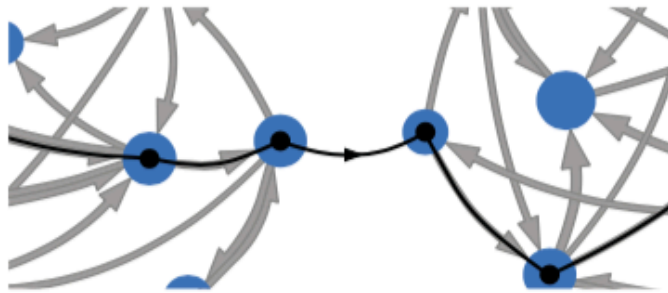
Maslov, S. & Redner, S. (2008) Promise and Pitfalls of Extending Google's PageRank Algorithm to Citation Networks. *The Journal of Neuroscience*

Teleportation Strategies

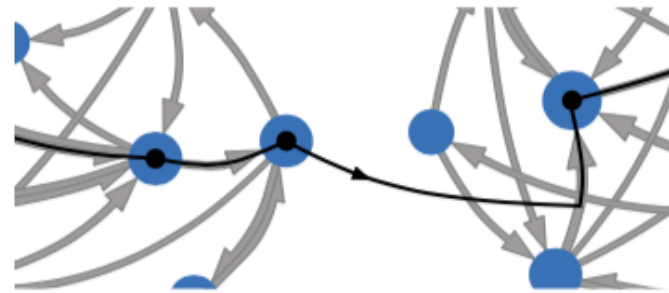


Smart Teleportation

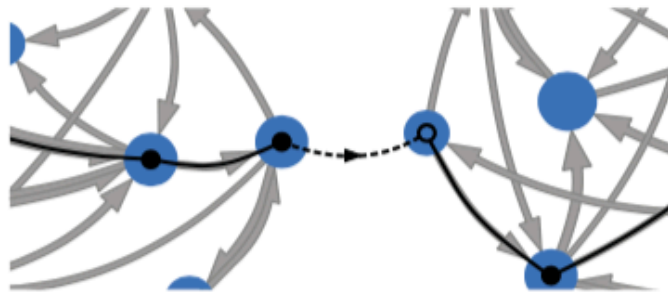
(a) Recorded node teleportation



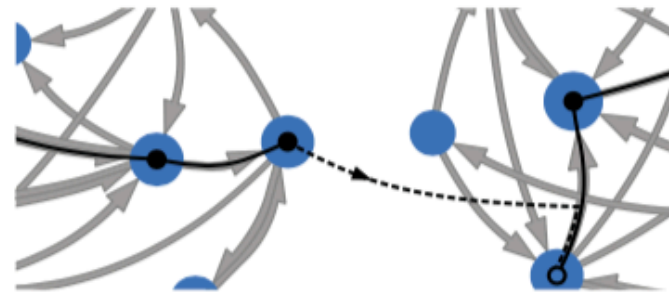
(b) Recorded link teleportation



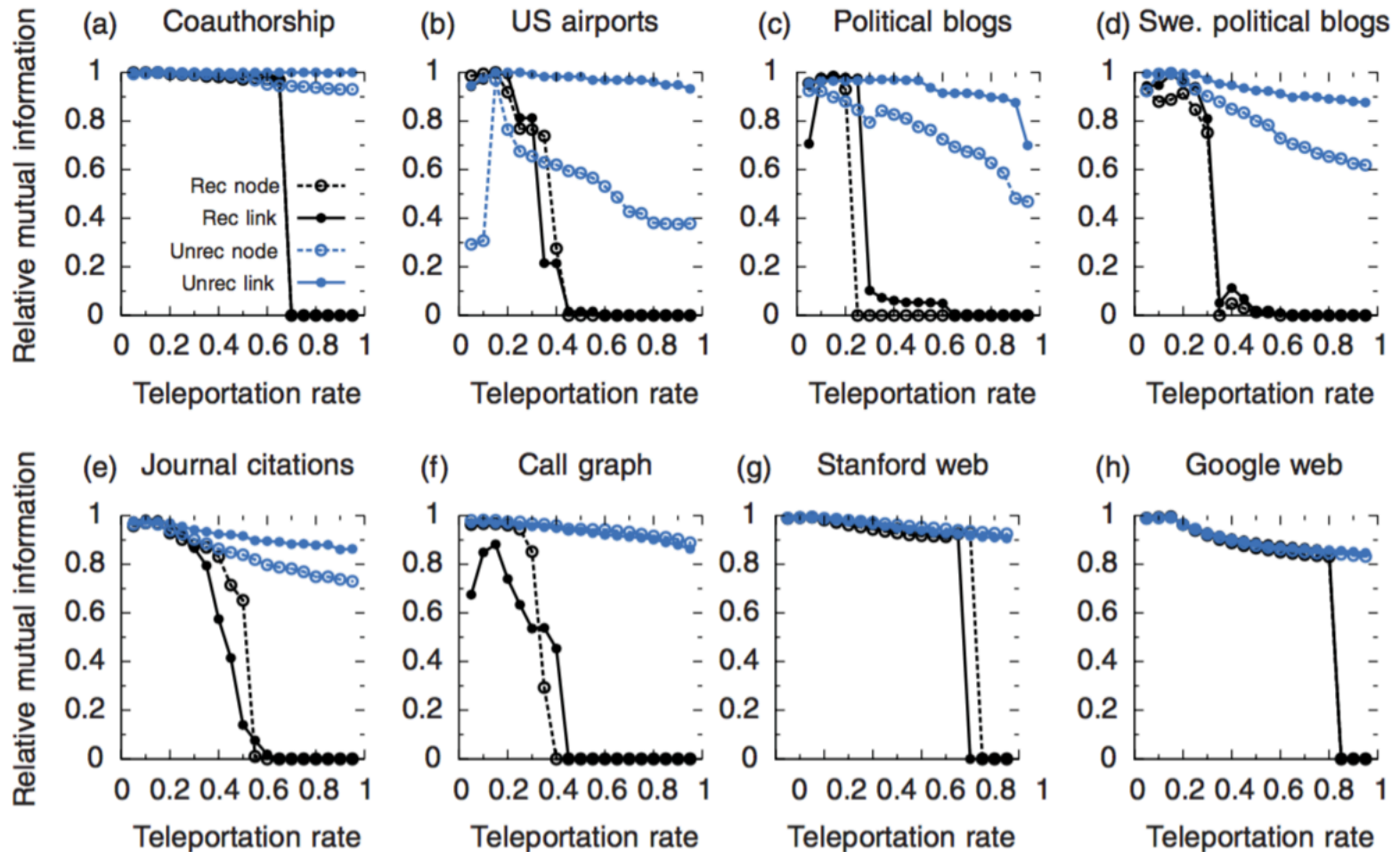
(c) Unrecorded node teleportation



(d) Unrecorded link teleportation

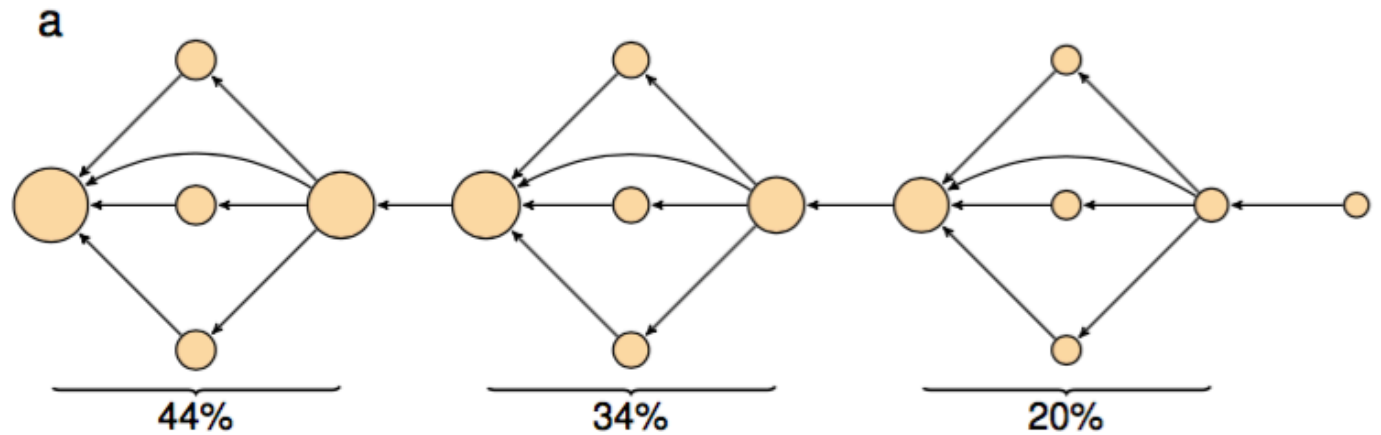


Smart Teleportation and Clustering



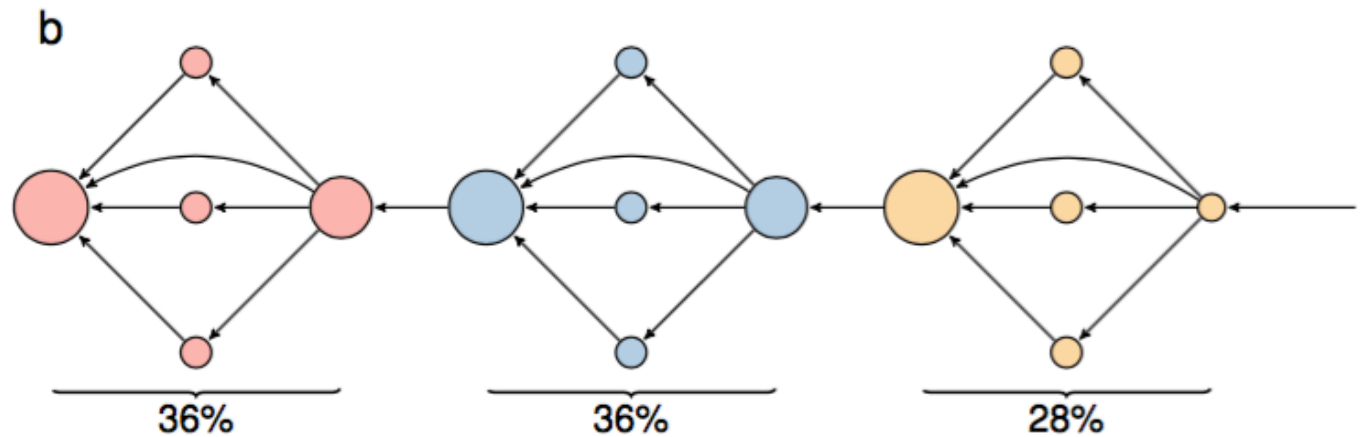
Article-level Ranking and Mapping

DIR-R (PageRank)

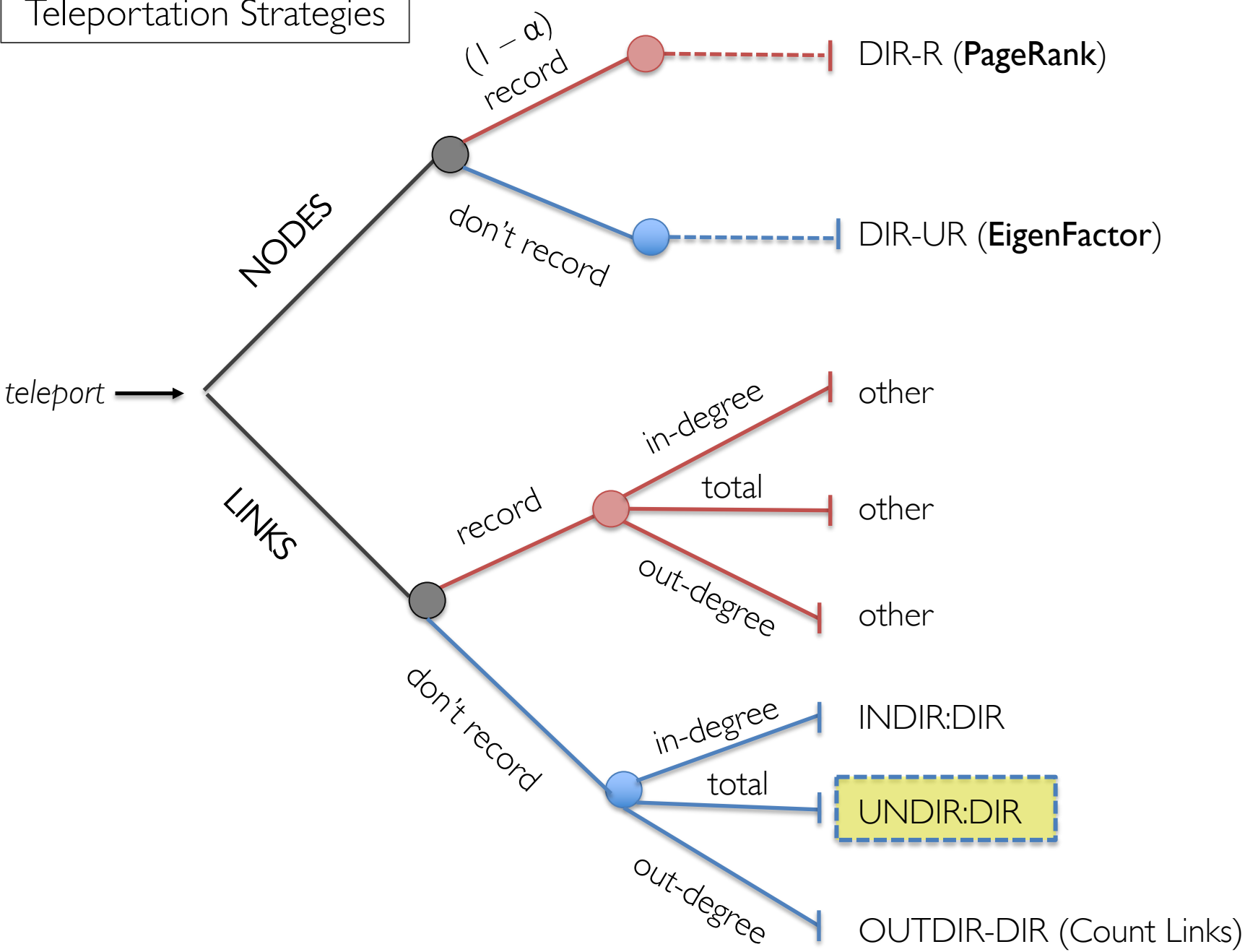


UNDIR:DIR

Smooths ranking
~
better clustering



Teleportation Strategies



Article-level Eigenfactor

$$w_i = \sum_j^n (Z_{ij} + Z_{ij}^T)$$

$$\mathbf{H}_{ij} = \frac{\mathbf{Z}_{ij}}{\mathbf{Z}_i}$$

$$\text{ALEF} = n \frac{\mathbf{H}_{ij}^T \cdot w_i}{\sum_i [\mathbf{H}_{ij}^T \cdot w_i]_i}$$

Static Ranking of Scholarly Papers using Article-Level Eigenfactor (ALEF)

Ian Wesley-Smith
The Information School
University of Washington
Seattle, WA 98195 USA
iwsmith@uw.edu

Carl T. Bergstrom
Department of Biology
University of Washington
Seattle, WA 98195 USA
cbergst@uw.edu

Jevin D. West
The Information School
University of Washington
Seattle, WA 98195 USA
jevinw@uw.edu

ABSTRACT

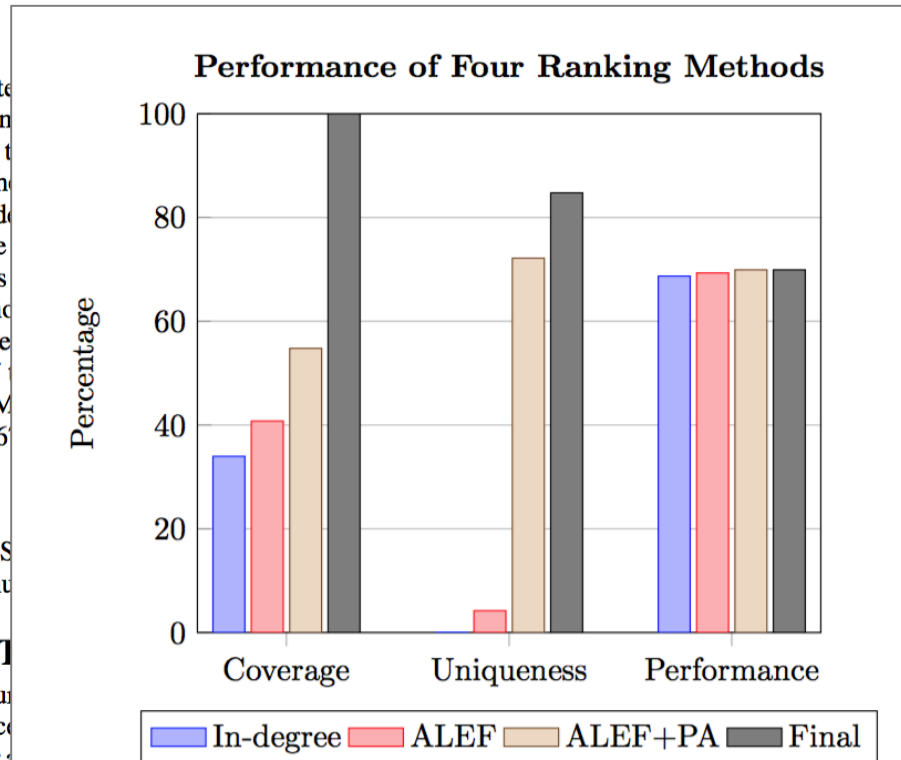
Microsoft Research hosted the 2016 WSDM Cup Challenge on the Microsoft Academic Graph (MAG) rankings for the articles to be evaluated against the Academic Graph provided by Microsoft. To evaluate its performance upon multiple facets of the contest (122M papers and 757M hyperlinks), the contest was scored at 0.6

Keywords

Information Retrieval; Scholarly Works; Scholarly Communication

1. INTRODUCTION

The scholarly literature is growing rapidly and some estimates place the number of scientific advances occurring



material in the vast academic world. Identifying documents that are both sufficiently high importance to rank is a problem of matching; the problem of ranking. Once documents are connected by hyperlinks among them, the graph provides extensive contextual information for ranking. Google's PageRank algorithm is an example. Before Google, most search results were returned without considering the quality of the results returned. Finding the most relevant web directories such as the Open Directory Project. Google showed us that latent structure in the world wide web lies the information. The problem of determining relevance and

graph structure, with papers and the practice of scholarly citation. The process of scholarly discovery has been applying network-based methods to determine the importance of entities [1]. Although our methods were not able to determine the importance of journals, we have developed a methodology which operates

a wealth of new approaches to investigation, and scholars publish their work in a rapidly diversifying variety of venues. Despite all this,

on individual articles: Article-Level Eigenfactor (ALEF) [27].

The 2016 WSDM Cup Challenge presented an excellent opportunity

WSDM CUP CHALLENGE

SIGN-UPS FOR THE WSDM CUP CHALLENGE ARE NOW CLOSED

The Graph

The Microsoft Academic Graph is a heterogeneous graph containing scientific publication records, citation relationships between publications, as well as authors, institutions, journal and conference "venues," and fields of study.

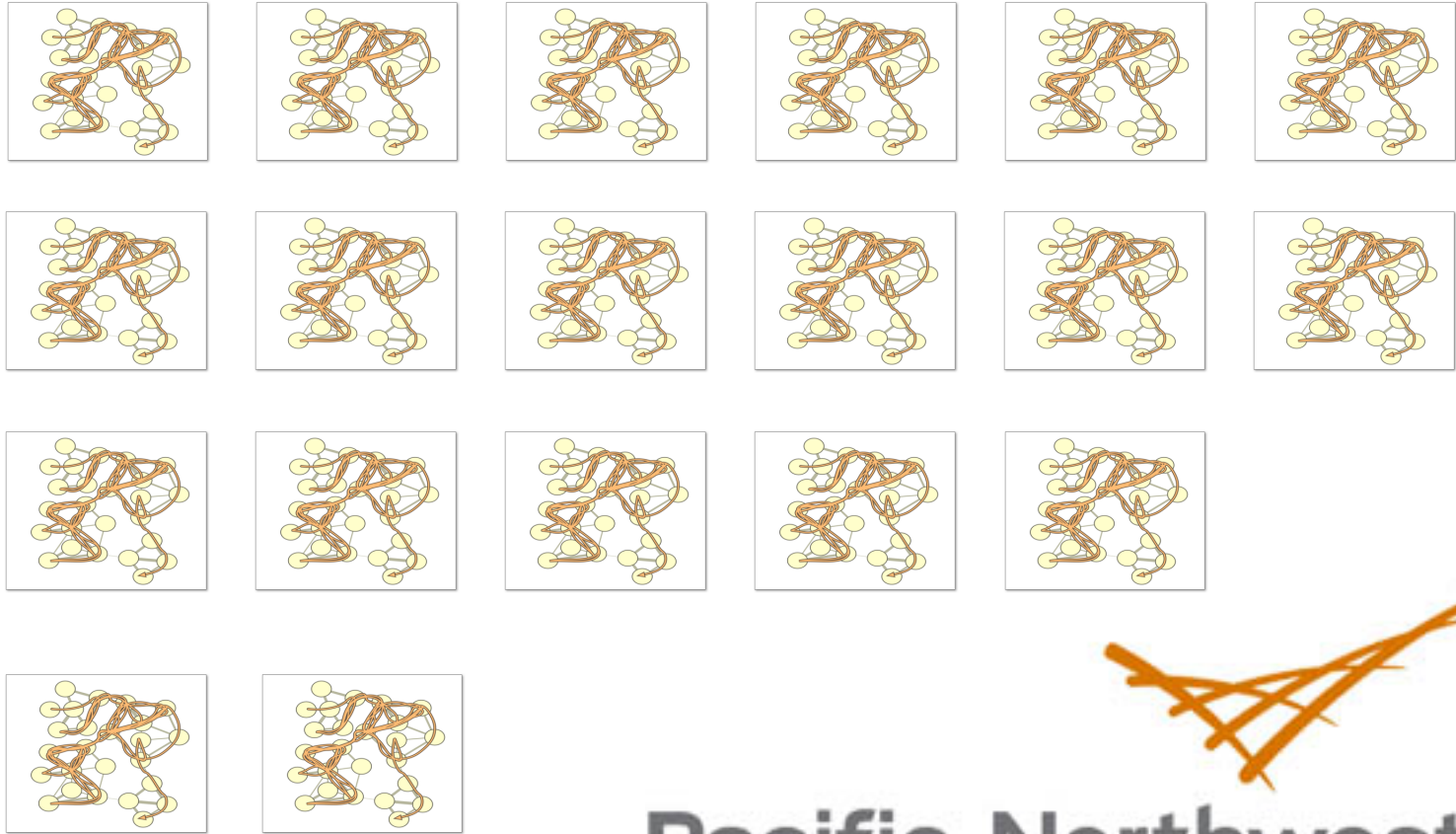
The Data

This data is available as a set of zipped text files stored in Microsoft Azure blob storage and available via HTTP. The file size (zipped) is ~30GB and may be downloaded [here](#).

The Challenge

The goal of the Ranker Challenge is to assess the query-independent importance of scholarly articles, using data from the Microsoft Academic Graph.

Running Experiments

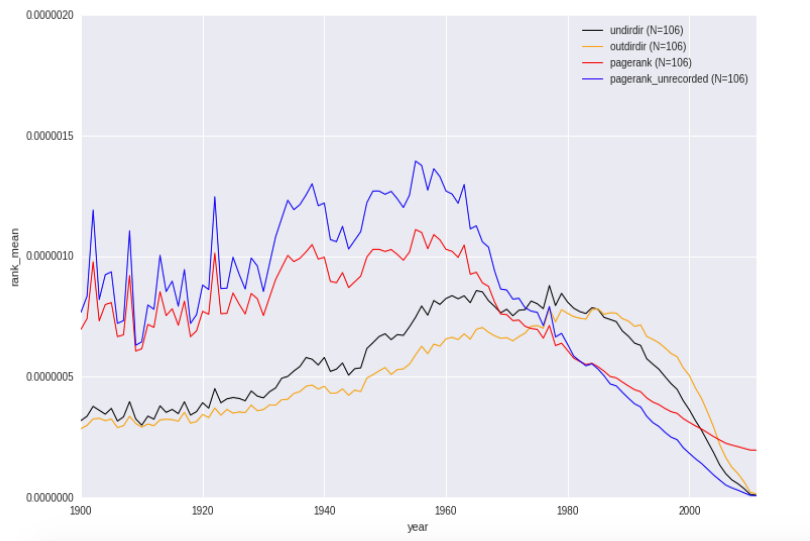


Pacific Northwest
NATIONAL LABORATORY

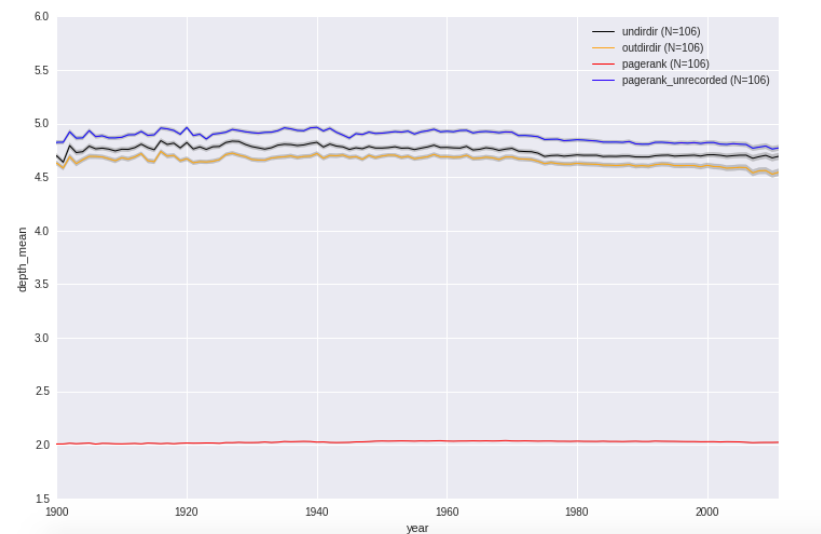
Clustering on time-directed networks

- Empirical exploration of hierarchical partitions with varying dynamics
- The effects of changing recorded teleportation ranking and clustering

Ranking Effects

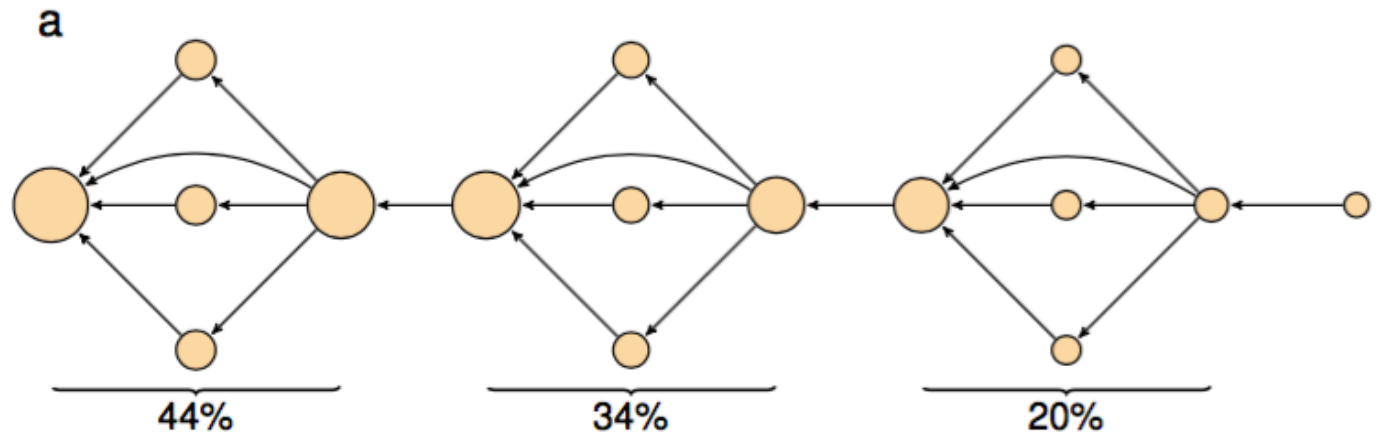


Clustering Effects



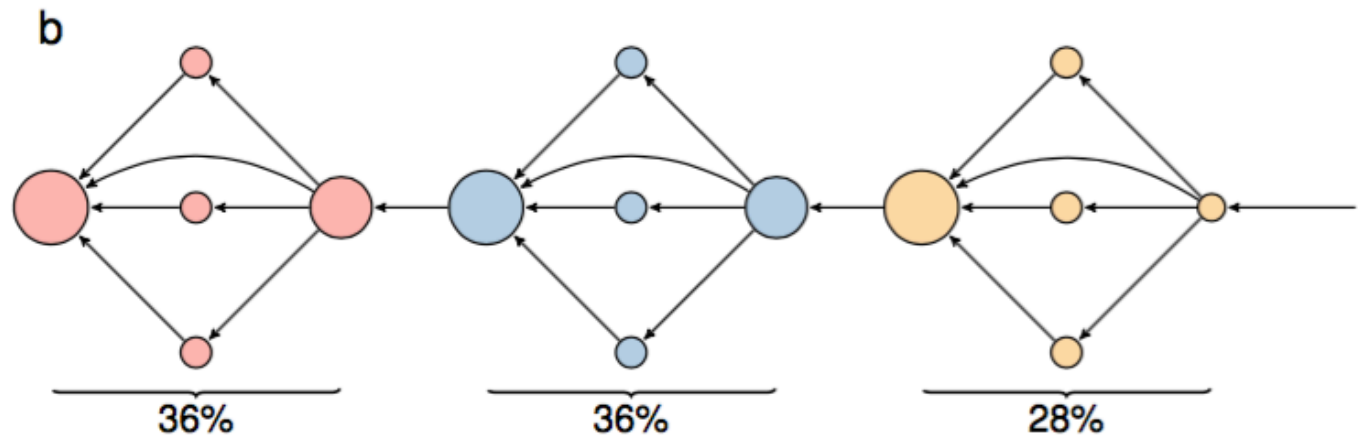
Article-level Ranking and Mapping

DIR-R (PageRank)

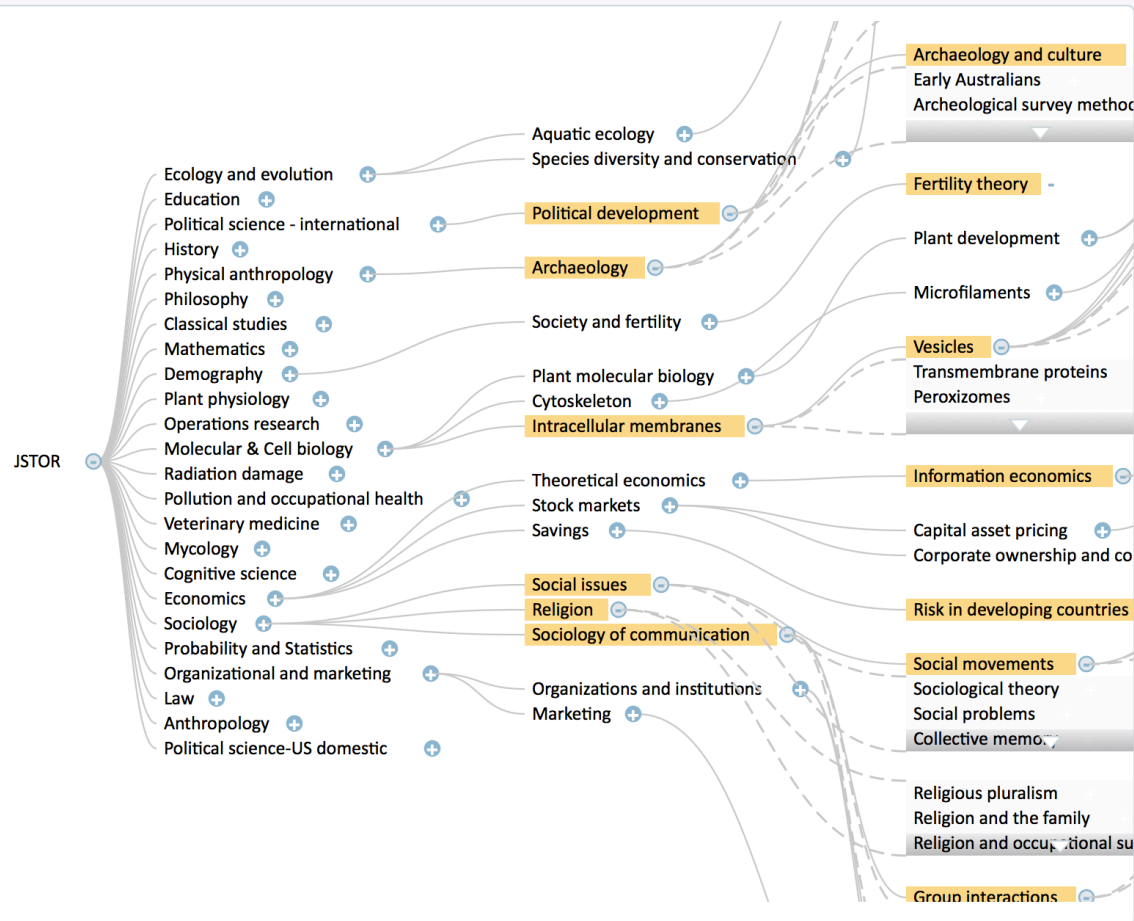


UNDIR:DIR

Smooths ranking
~
better clustering



Revealing Hierarchical Structure



Top Papers Sort by Year (newest) ▾

- [Using Siting Algorithms in the Design of Marine Reserve Networks](#)

Heather Leslie - *Ecological Applications* (2003)
- [Mechanism of Filopodia Initiation by Reorganization of a Dendritic Network](#)

Tatyana Svitkina - *The Journal of Cell Biology* (2003)
- [Network Structure and Knowledge Transfer: The Effects of Cohesion and Range](#)

Ray Reagans - *Administrative Science Quarterly* (2003)
- [A General Model for Designing Networks of Marine Reserves](#)

Eric Sala - *Science* (2002)
- [The Density of Social Networks and Fertility Decisions: Evidence from South Nyanza District, Kenya](#)

Hans-Peter Kohler - *Demography* (2001)
- [A New Dynammin-Like Protein, ADL6, Is Involved in Trafficking from the trans-Golgi Network to the Central Vacuole in Arabidopsis](#)

Jing Bo Jin - *The Plant Cell* (2001)
- [Comparing Sequenced Segments of the Tomato and Arabidopsis Genomes: Large-Scale Duplication Followed by Selective Gene Loss Creates a Network of Synteny](#)

Hsin-Mei Ku - *Proceedings of the National Academy of Sciences of the United States of America* (2000)
- [A Noncooperative Model of Network Formation](#)

Venkatesh Bala - *Econometrica* (2000)

Find Papers

by title

by field

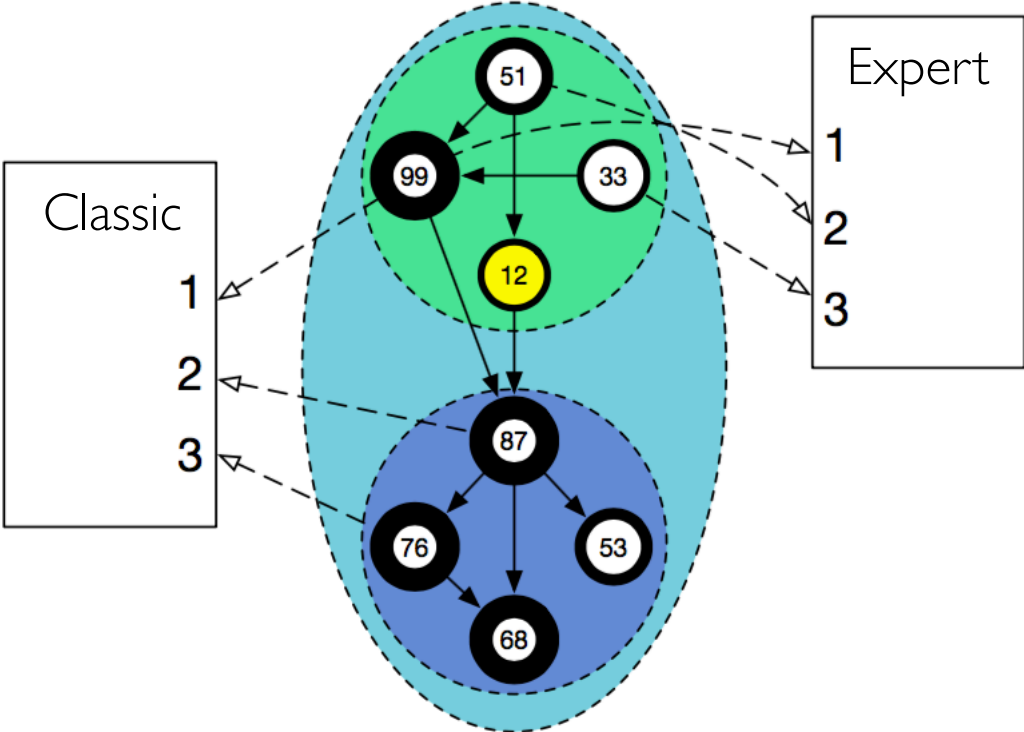
by author

by journal

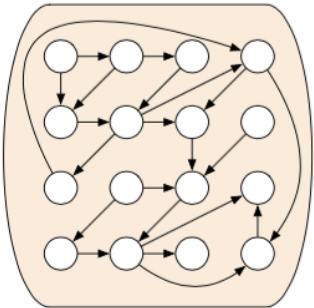
Active Queries: [clear all](#)

✕ keyword: network

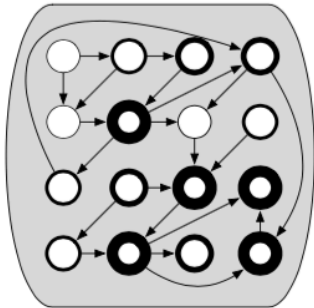
Recommend



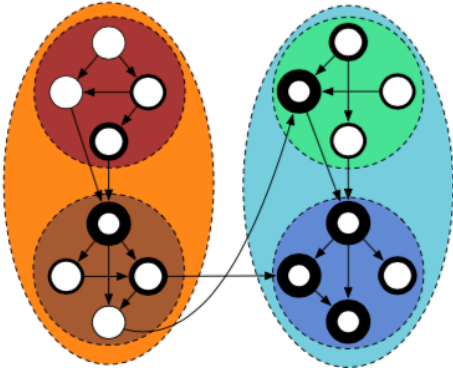
Assemble



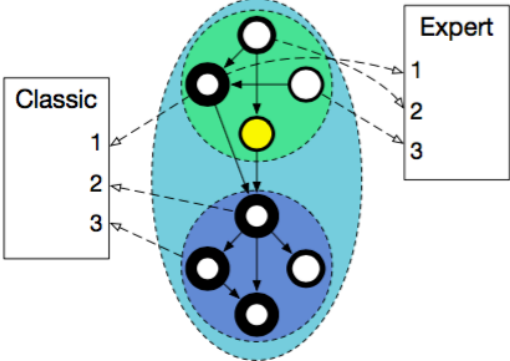
Rank



Cluster

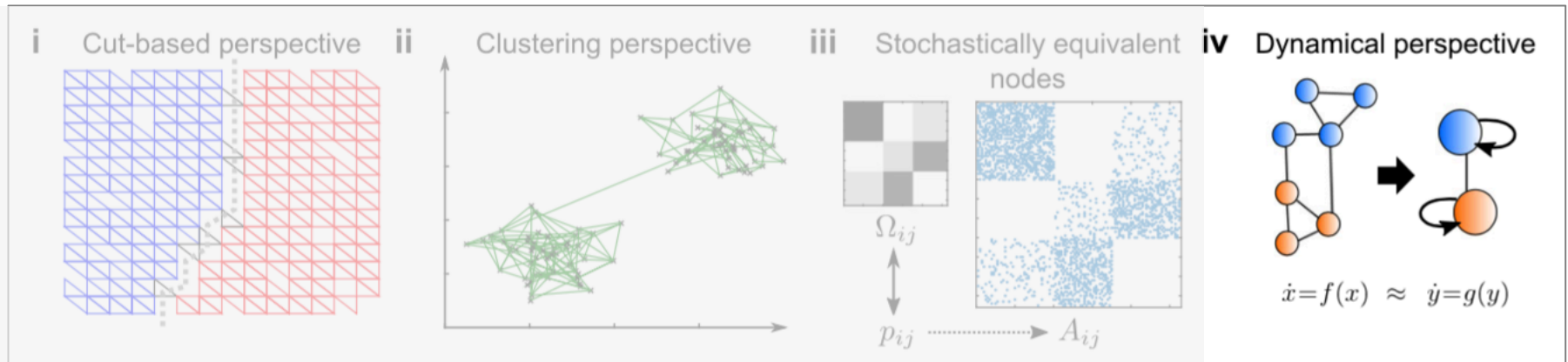


Recommend



West, Wesley-Smith, Bergstrom (2016) A recommendation system based on hierarchical clustering of an article-level citation network. *IEEE, Transactions on Big Data* (in press)

Community Detection Perspectives



Circuit layout
 Minimizing cuts
 Load balancing
 Eigenvectors
 Spectral methods
 Image segmentation

Data Clustering
 Maximizing node density
 unknown k , unbalanced
 Conductance
 Local, global
 Modularity

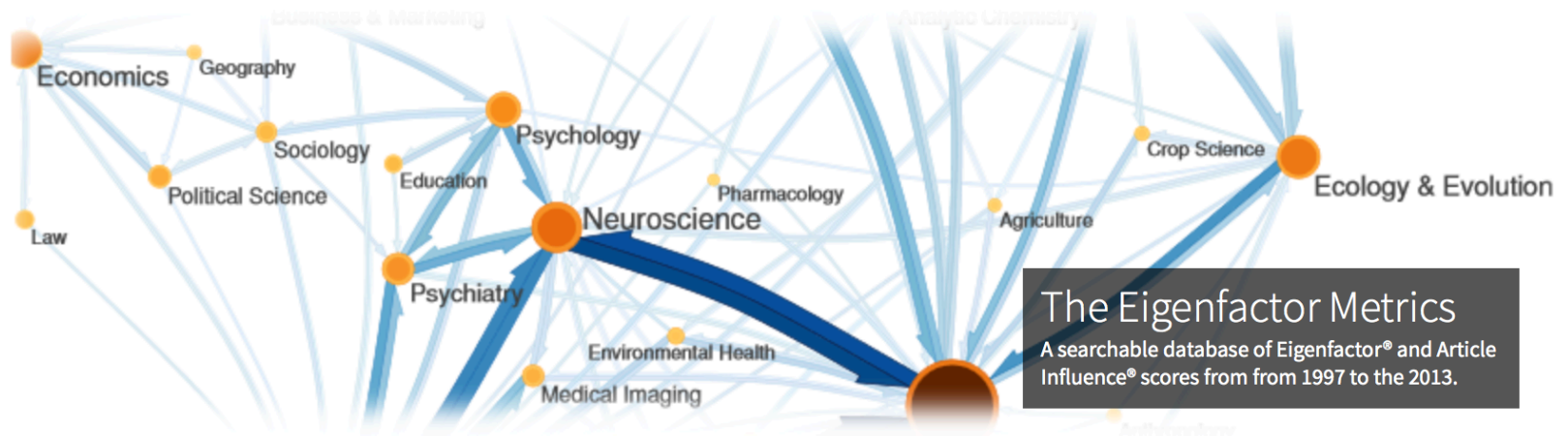
Social Networks
 Connectivity Profiles
 Stochastic equivalence
 SBMs, LFR
 p -values, hypothesis testing
 Bipartite treatment
 Predict missing links

System behavior, processes
 Non-adjacency focused
 Airline network
 Markovian diffusion process
 Undirected, Directed
 InfoMap



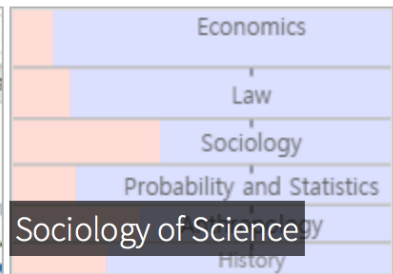
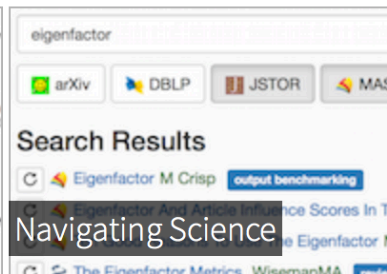
Summary

- Community detection – one size does not fit all
- Citation networks - dynamical perspective
- Memory - higher order dynamics
- Unrecorded teleportation to links (*undir*) improves ranking and hierarchical clustering
- Next steps – building benchmarks and methods for evaluating the different rankings and hierarchical clusterings (refer to Jennifer Webster's talk tomorrow)



The Eigenfactor Metrics
 A searchable database of Eigenfactor® and Article Influence® scores from from 1997 to the 2013.

RESEARCH AREAS



NEWS

23 Nov. **JEVIN WEST ON MEGAJOURNALS IN THE *CHRONICLE OF HIGHER EDUCATION***
 Jevin West discusses the rise of the megajournal and our [open access cost effectiveness tool](#) in the *Chronicle of Higher Education*.

23 Nov. **EIGENFACTOR TEAM PLACES SECOND IN MICROSOFT RESEARCH'S WSDM CUP**
 The [WSDM Cup Challenge](#) asked teams to use 30GB of data from the Microsoft Academic Graph to rank the

Acknowledgements

Carl Bergstrom, Department of Biology, University of Washington

Martin Rosvall, Department of Physics, Umea University

Seung-Hee Bae, Computer Science, Western Michigan

Jason Portenoy, Information School, University of Washington

Bill Howe, eScience, CSE, University of Washington

Jennifer Webster, Pacific Northwest National Laboratory

Jevin West

jevinw@uw.edu

jevinwest.org

@jevinwest