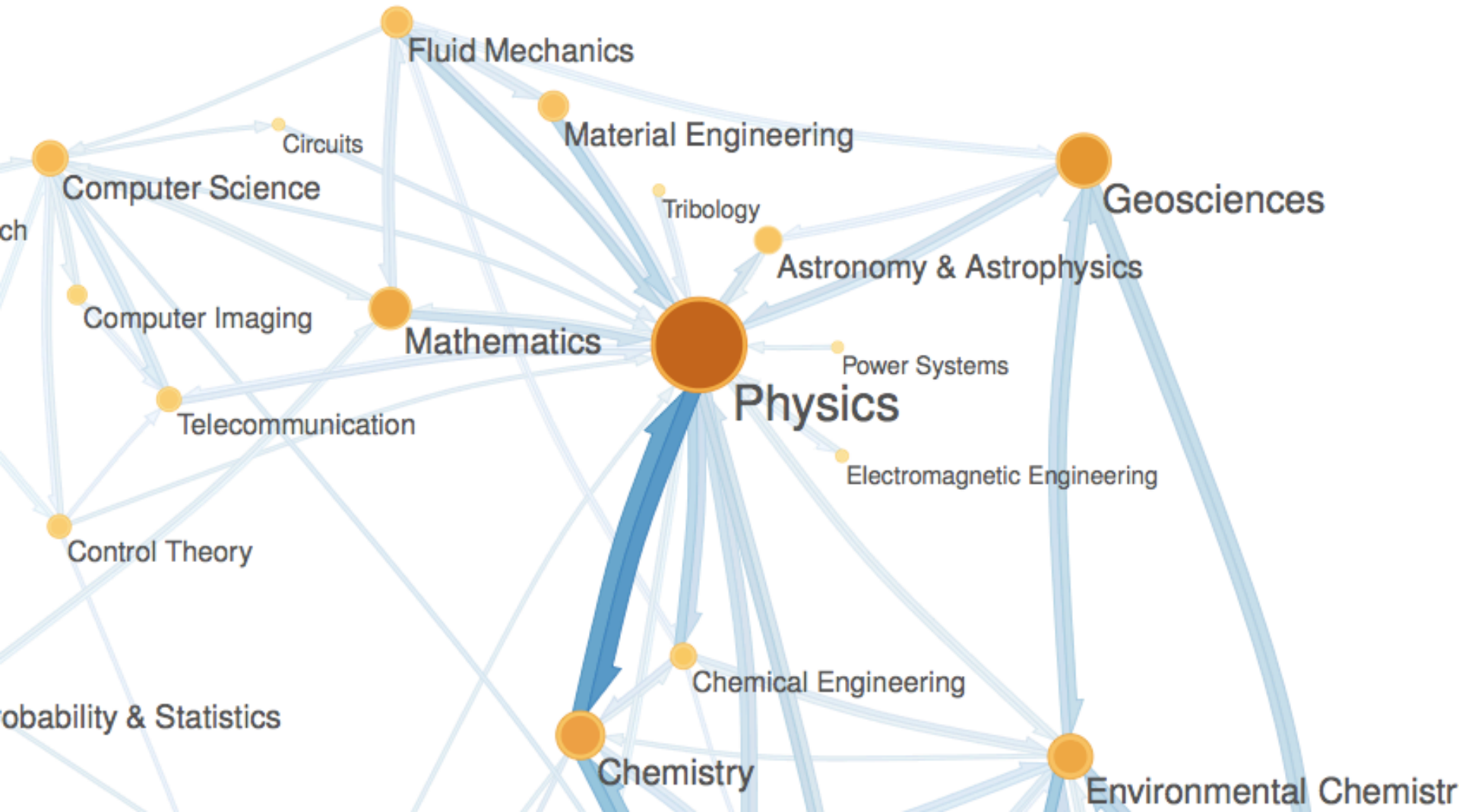
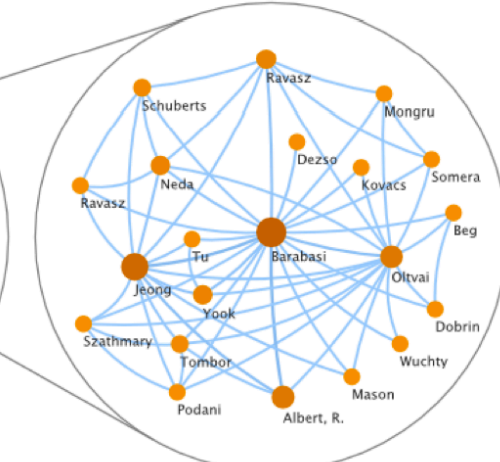
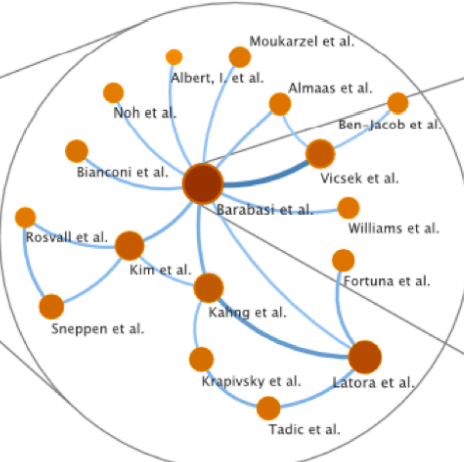
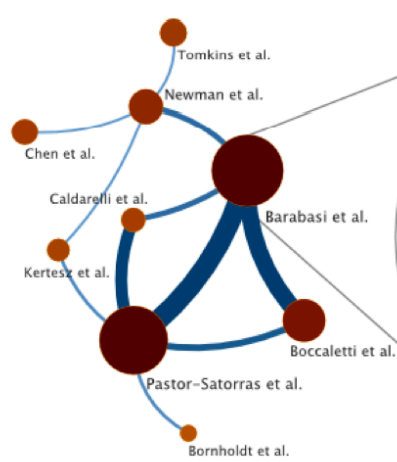
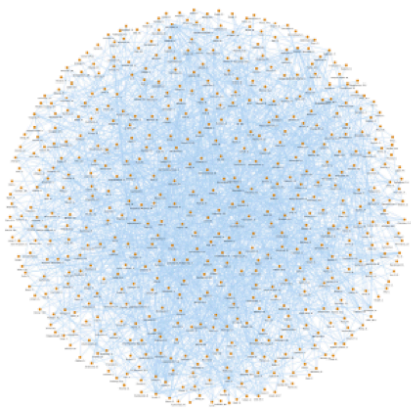
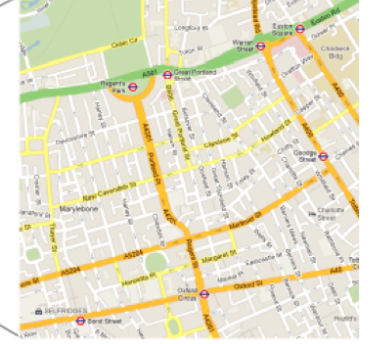
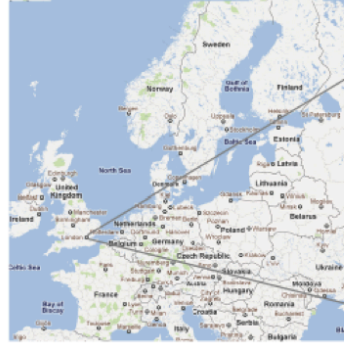


Mapping Knowledge Networks

Jevin West, Information School, University of Washington



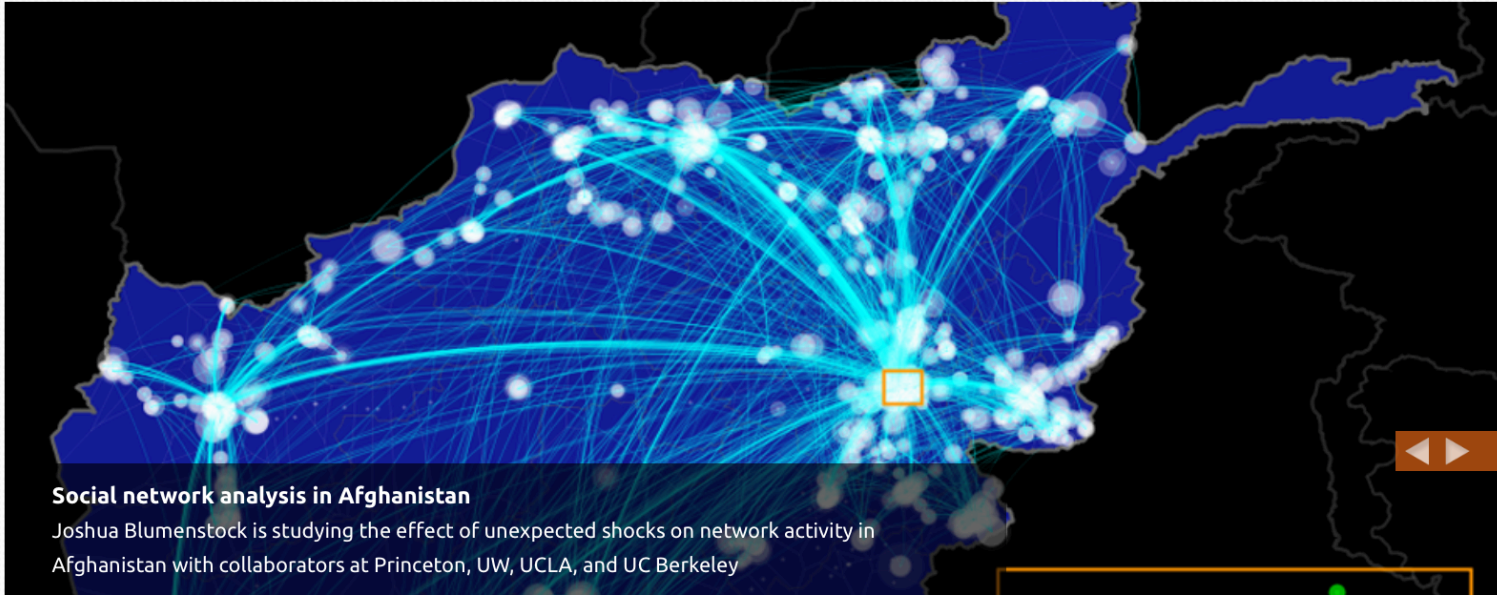


Data

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.

Minimal Description Length (MDL) Statistics



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



Data for Development



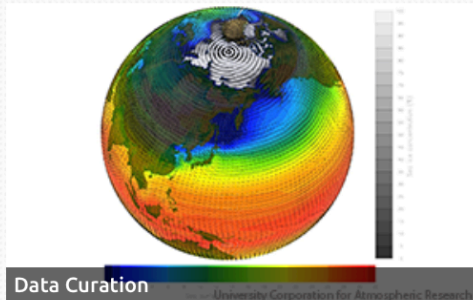
Social Networks



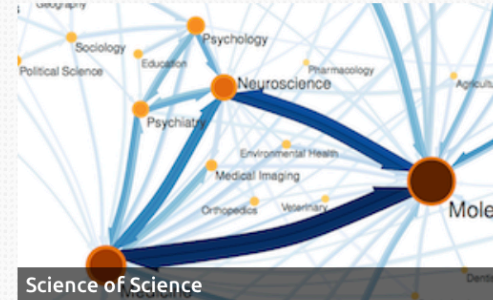
Data Visualization



Computational Social Science



Data Curation



Science of Science



What We Do



Overview

Over the course of the last decade many disciplines have evolved from recording observations in laboratory notebooks to the use of instruments capable of digitally recording many gigabytes of data in a day. This abundance of data provides unprecedented opportunities for discovery. Tapping its potential requires the application of sophisticated new computational techniques operating on large scale storage, computational and network resources. Since its creation in 2008, the eScience Institute has worked to create the intellectual and physical infrastructure needed to meet this challenge.

At the core of the eScience Institute are individuals who have proven track records in developing and applying advanced computational methods and tools to real world problems. Their task is to seek out and engage researchers across disciplines where eScience approaches are likely to have the greatest impact. To ensure that researchers have access to the necessary physical infrastructure, the Institute has undertaken coordinated planning and support for advanced local and remote computational platforms. This includes developing relationships with commercial and non-commercial service providers as well as the development of shared facilities on campus. This support extends to assistance in the preparation of select proposals where we are able to focus resources, improving their chances for success.

Also in... What We Do

[Appliance Gallery](#)

Find and use the eScience Institute's virtual machines equipped with software useful for specific applications.

[Campus Compute & Storage](#)

Learn about what UW is doing to support scalable scientific computing on campus

[Consulting & Services](#)

From algorithm development to database creation to cloud computing, we can help.

[Projects](#)

Explore some of our current collaborations with research scientists.

[Relevant Courses](#)

View a list of courses offered in eScience disciplines.

[SQLShare Success Stories](#)

[Tools](#)

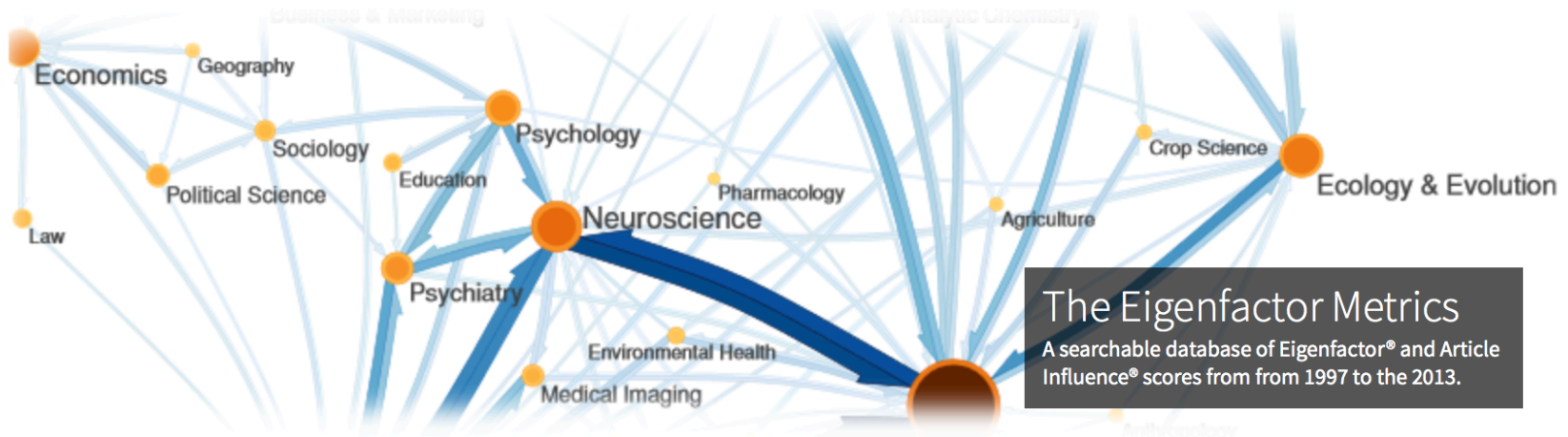
Whether it's database management, visualization, or developer tools, learn about tools we can help you use.

Latest eScience News

[Data Science Incubation Program - Winter 2016](#)

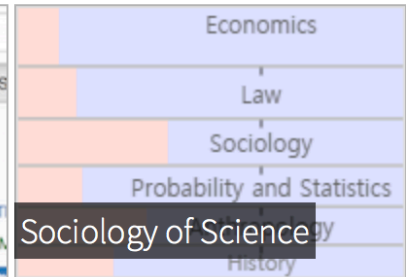
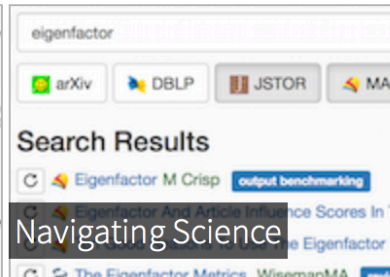
2 hours 4 min ago

[Ben Marwick On How Computers Broke Science](#)



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

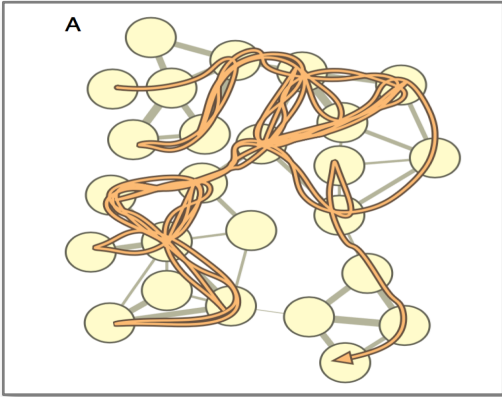
RESEARCH AREAS



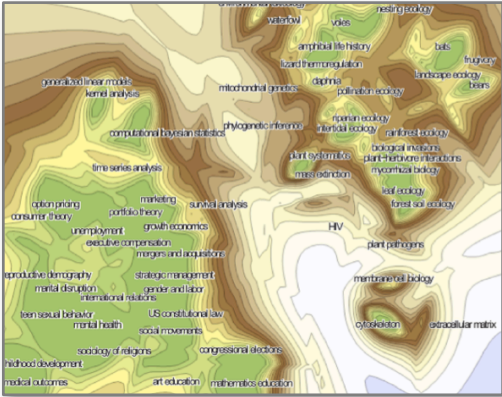
NEWS

5 **NATURE ON SELF-CITATION**
 July. Nature featured our paper on gender differences in self citation.

18 **THE ECONOMIST ON THE VIZIOMETRICS PROJECT**
 June. The Economist published an article describing our Viziometrics.org project.

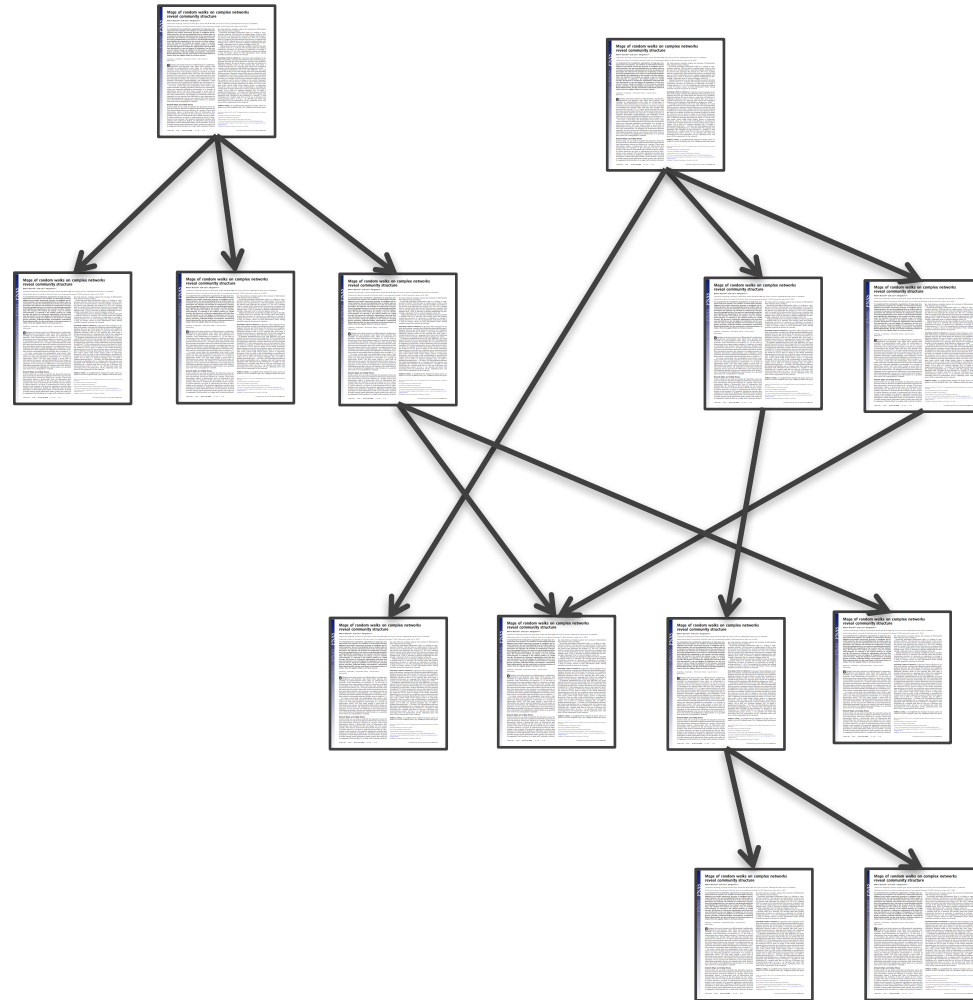


Science of Mapping



Mapping of Science

Citations form a vast network



de Solla Price, *Science* (1965)



The Scholarly Graph



THOMSON REUTERS

PatentVector™



PNAS



dblp
computer science bibliography





The Scholarly Graph



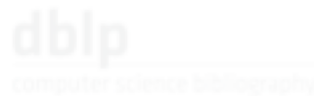
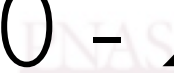
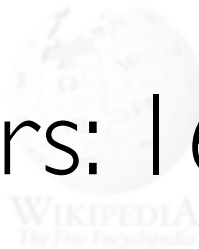
Tens of millions articles, patents, books

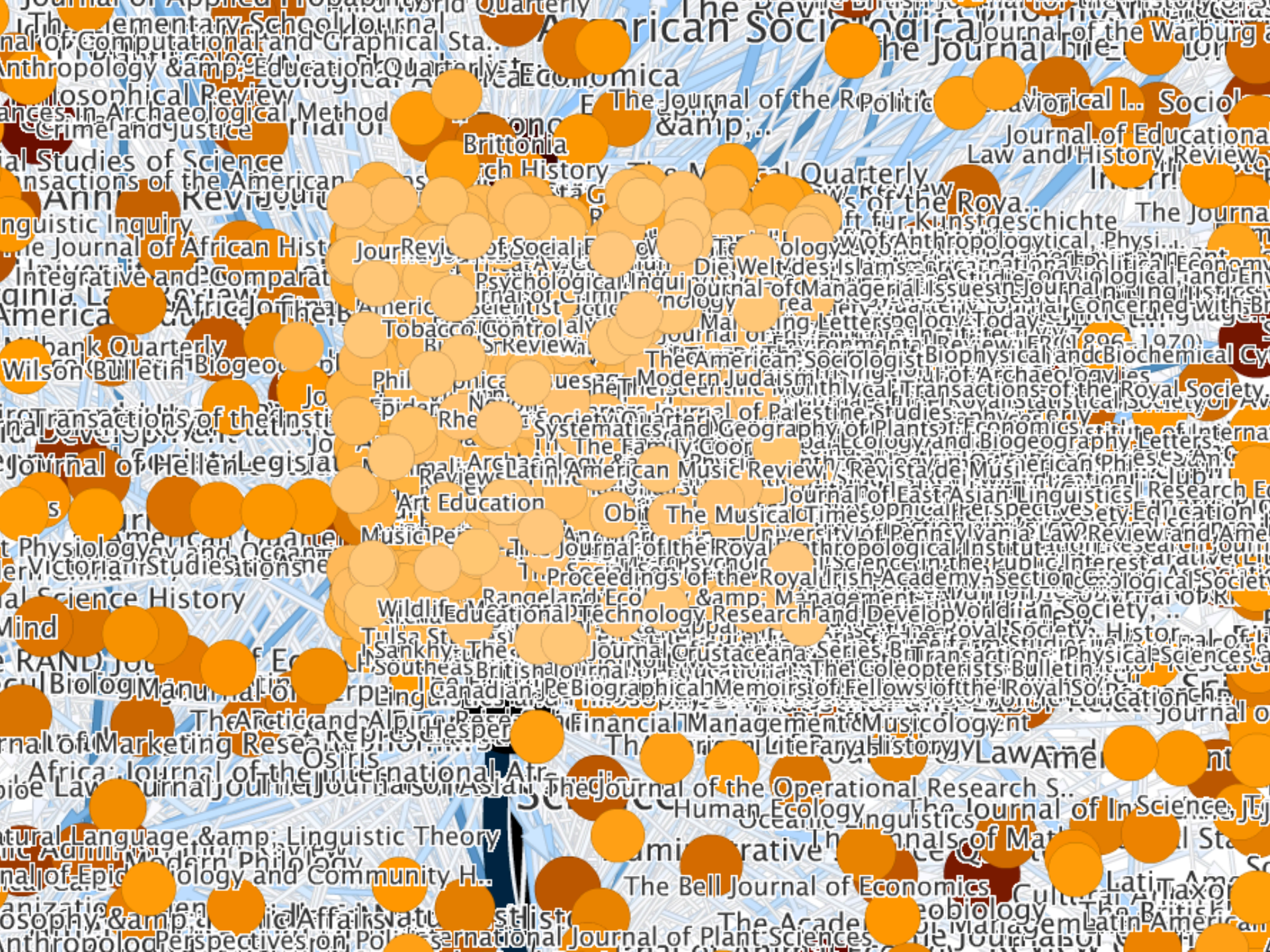


Billions of citation links



Years: 1600 - 2016





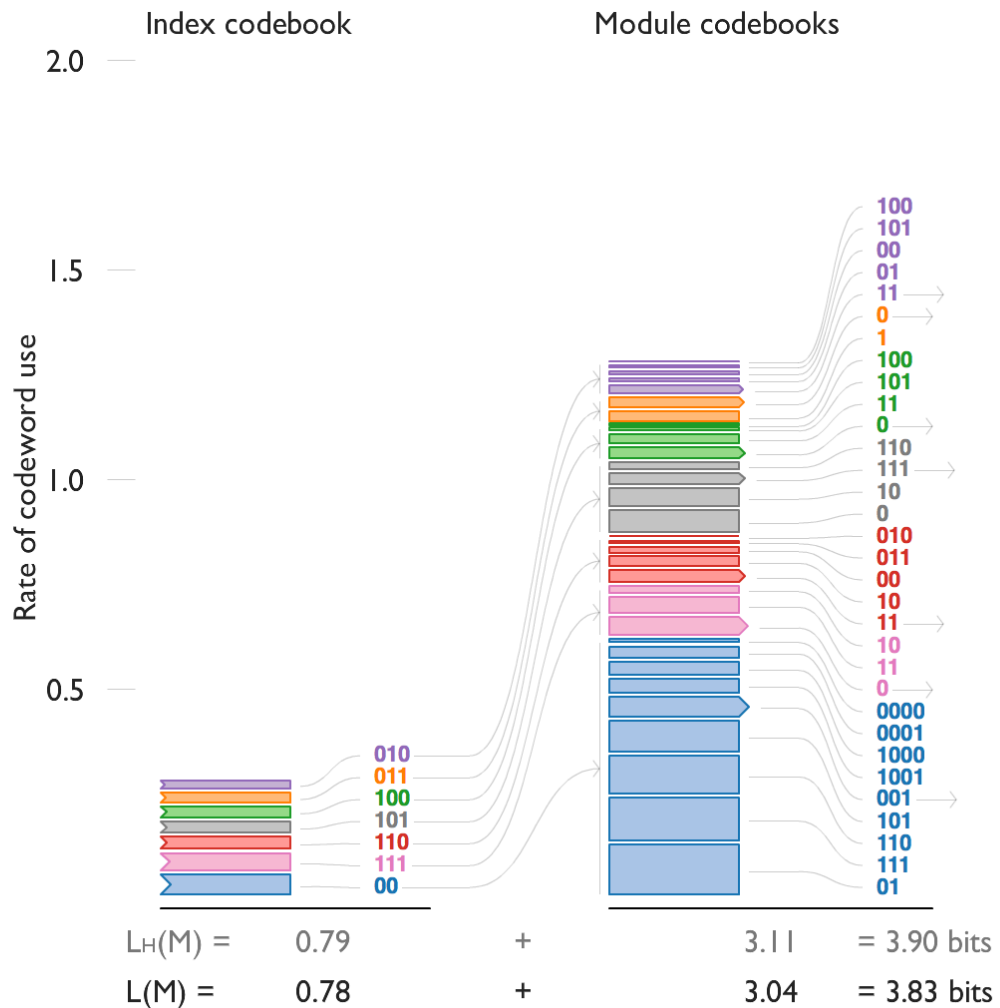
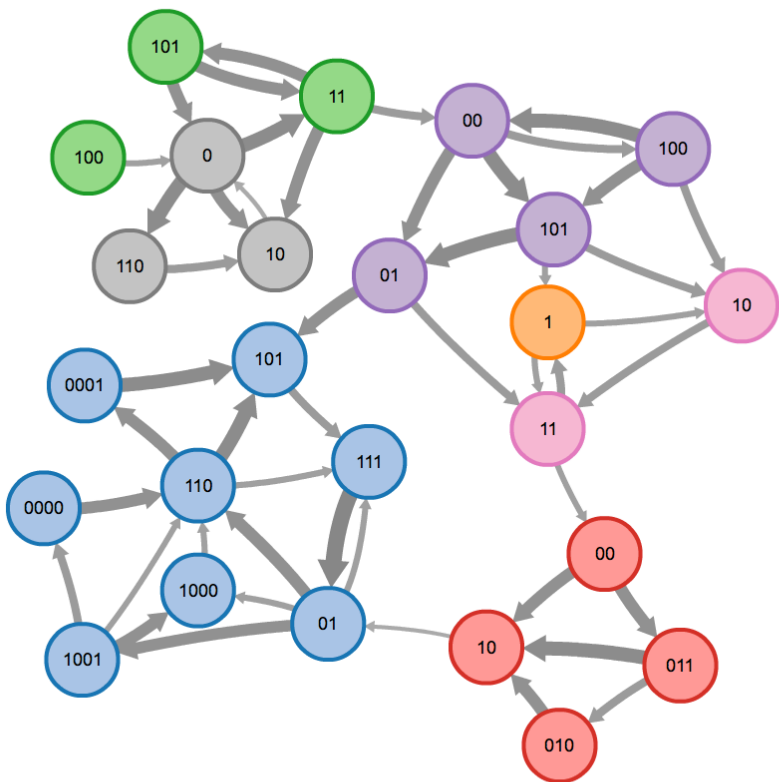
The map equation

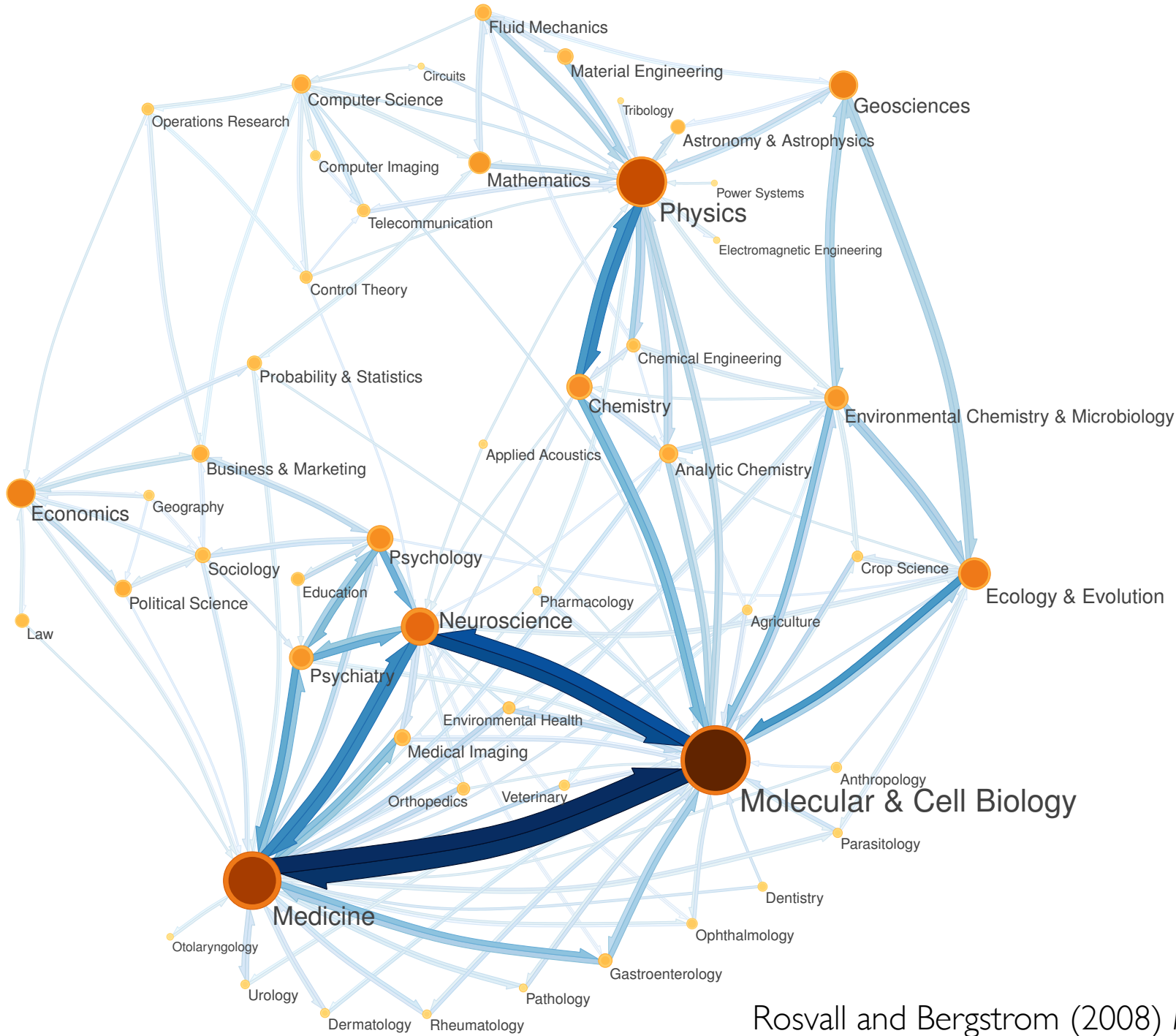
frequency of inter-module movements

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

code length of module names

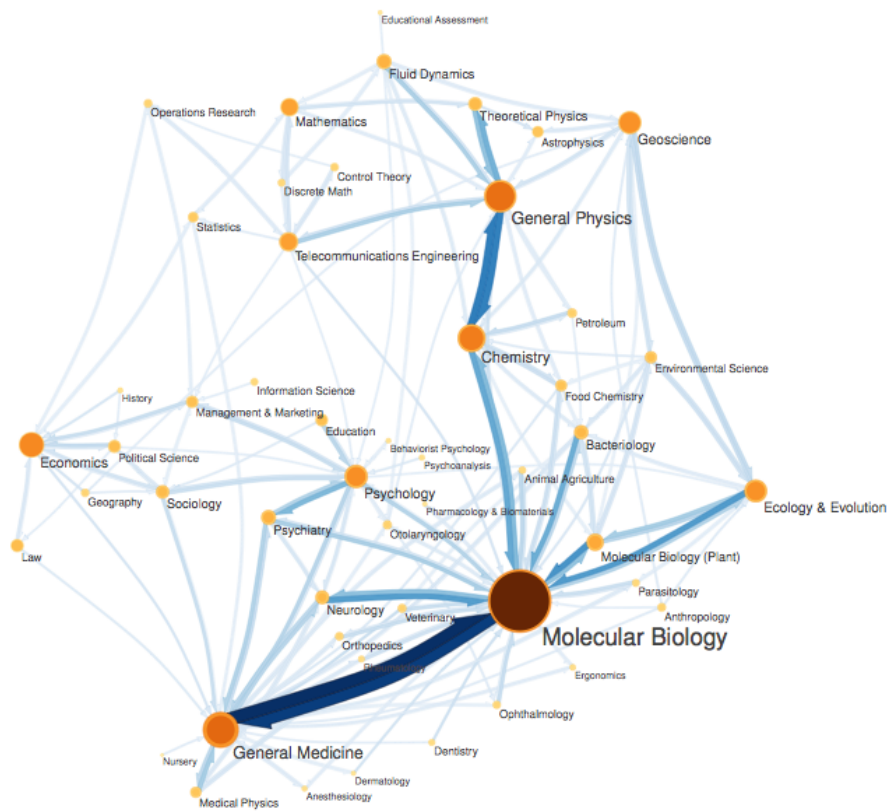
code length of node names in module i



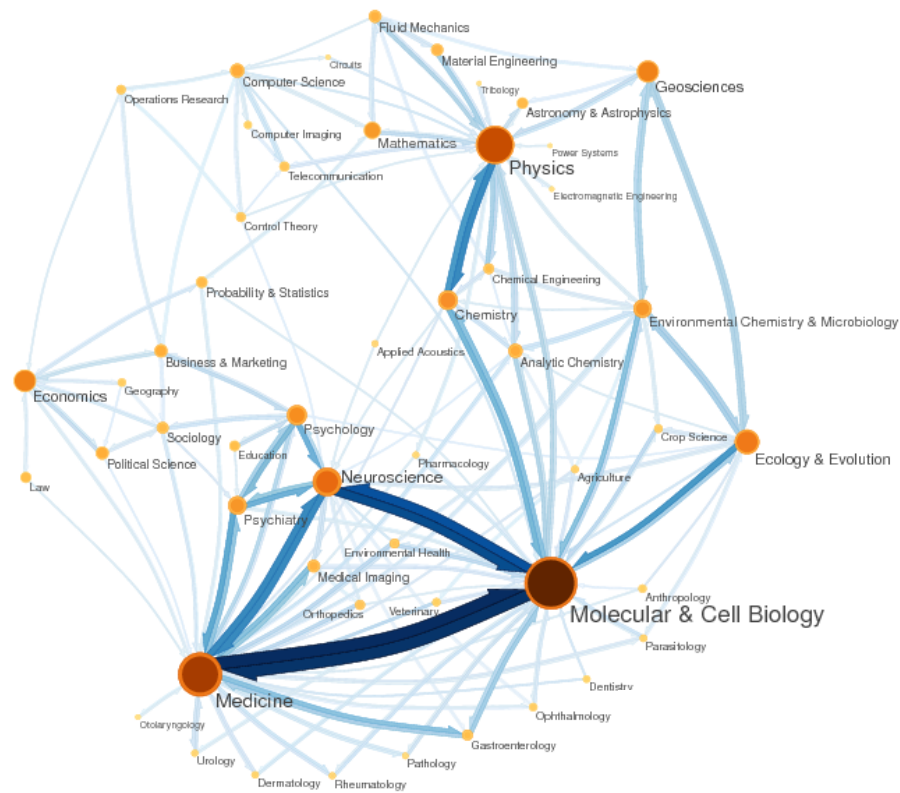


Rosvall and Bergstrom (2008) *PNAS*

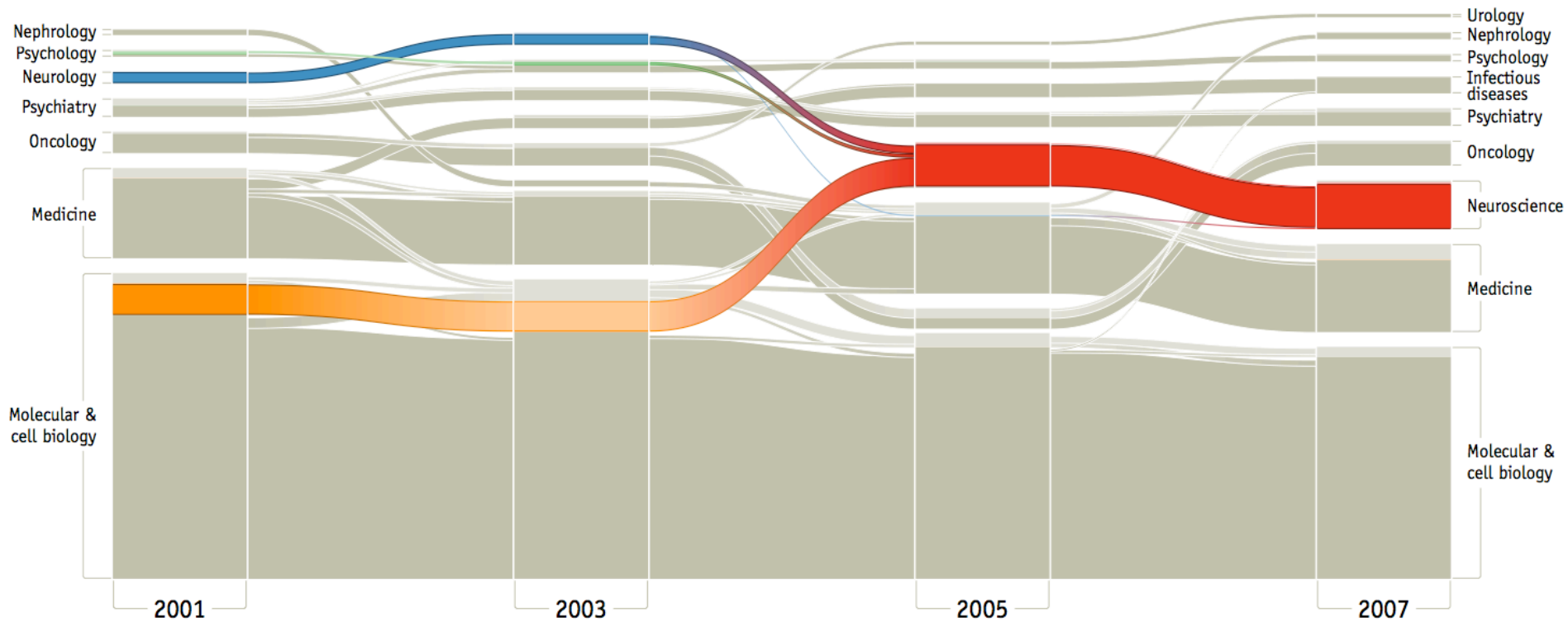
1995

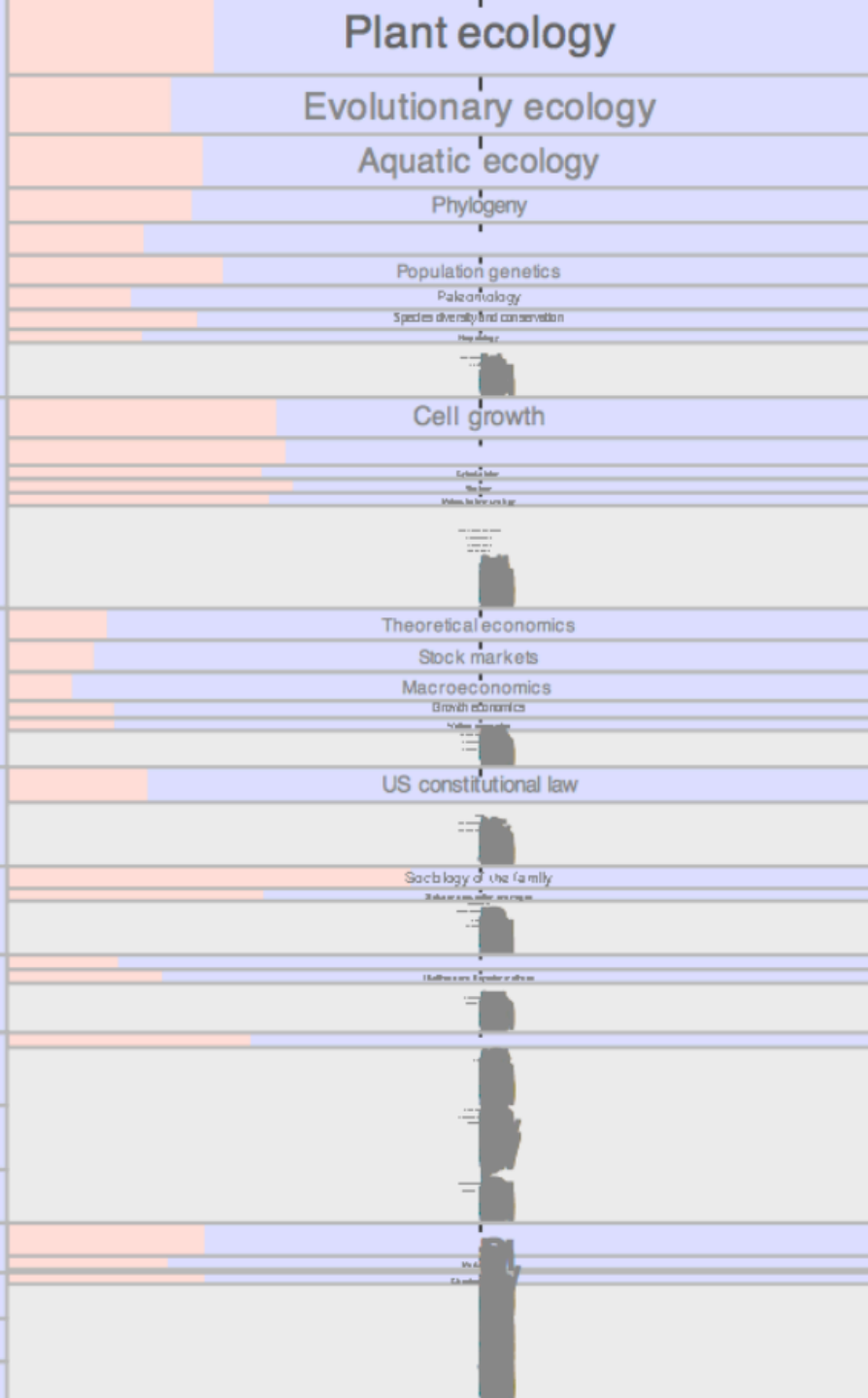
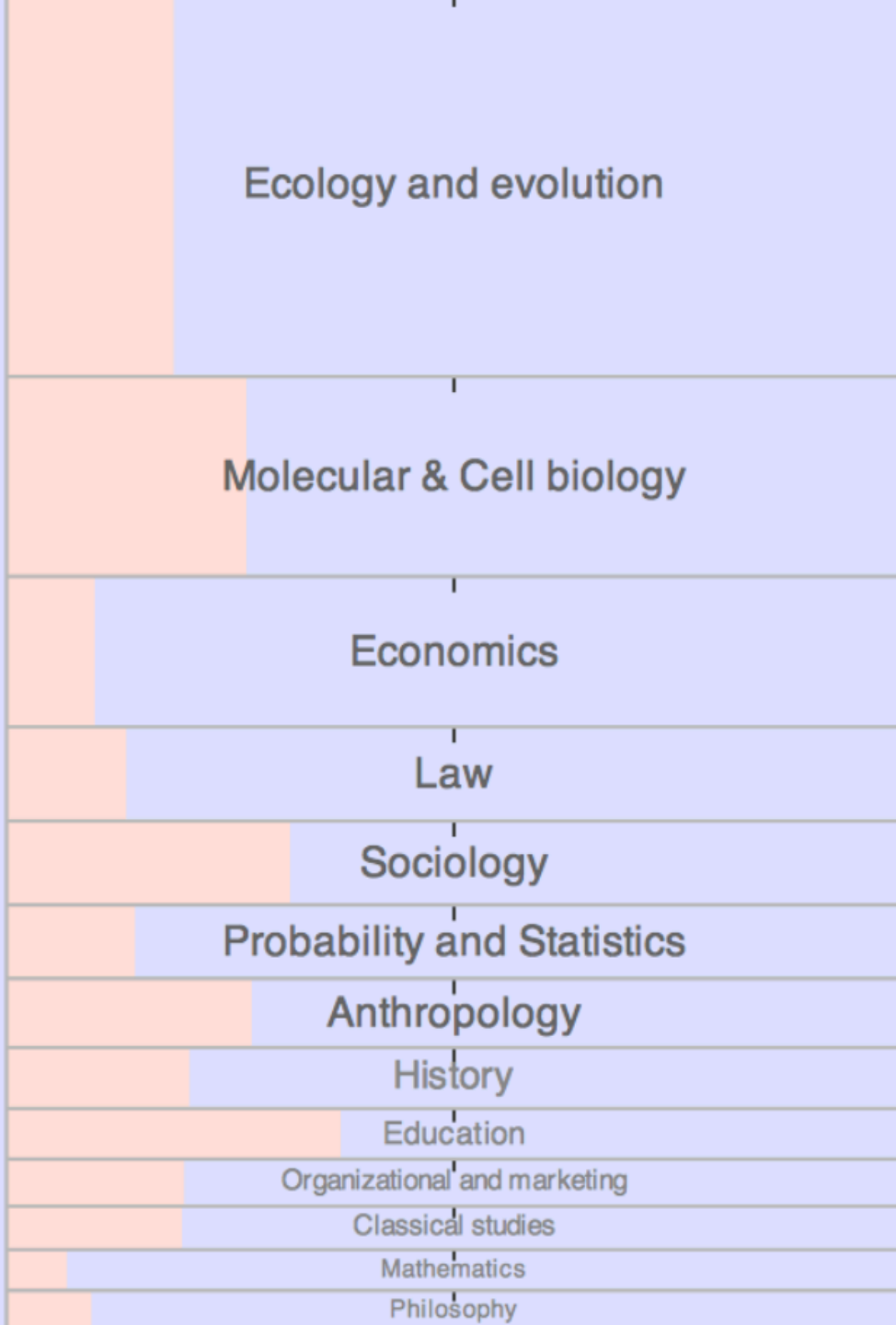


2004

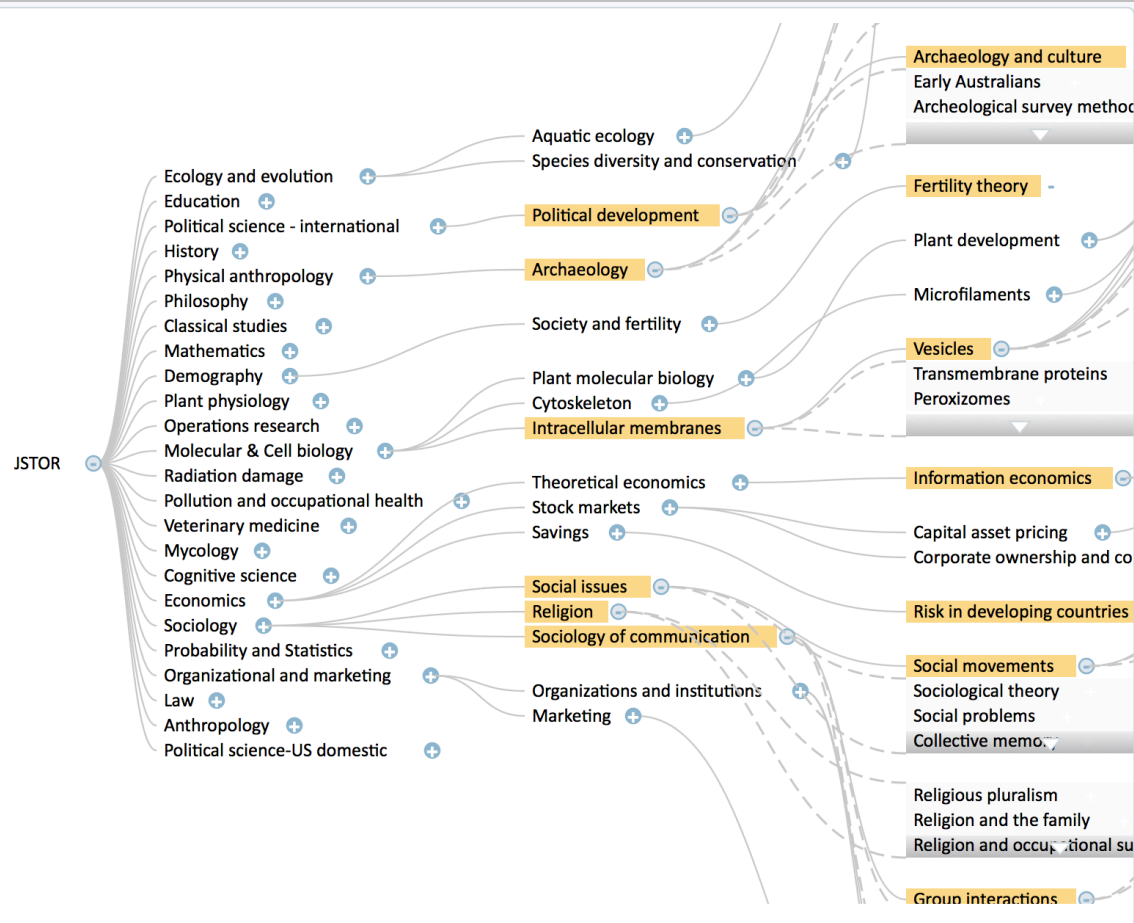


The Emergence of Neuroscience





“Network”



Find Papers

by title

by field

by author

by journal

Active Queries: [clear all](#)

✕

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)

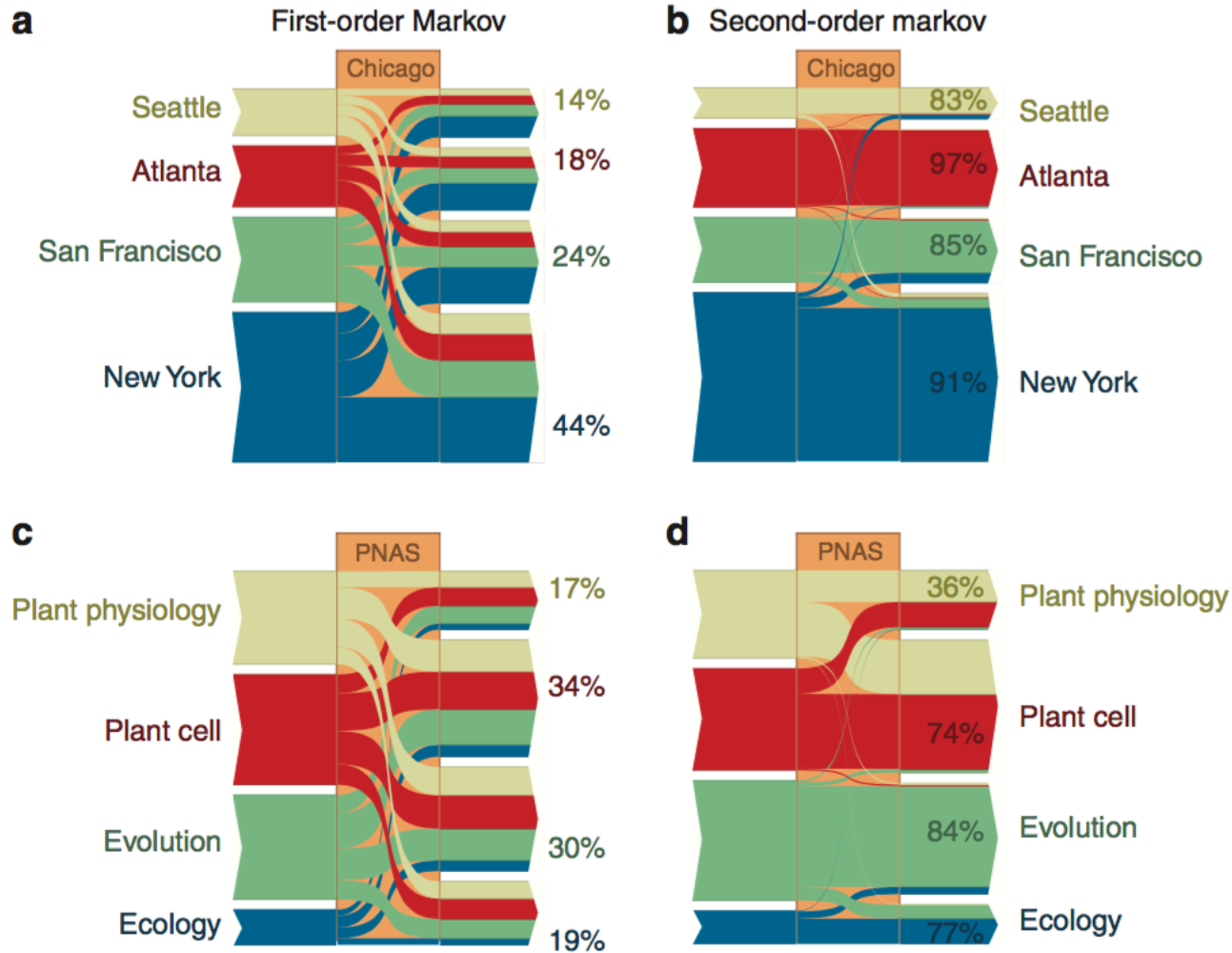
[The Density of Social Networks and Fertility Decisions: Evidence from South Nyanza District, Kenya](#)
Hans-Peter Kohler - *Demography* (2001)

[A New Dynamin-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)

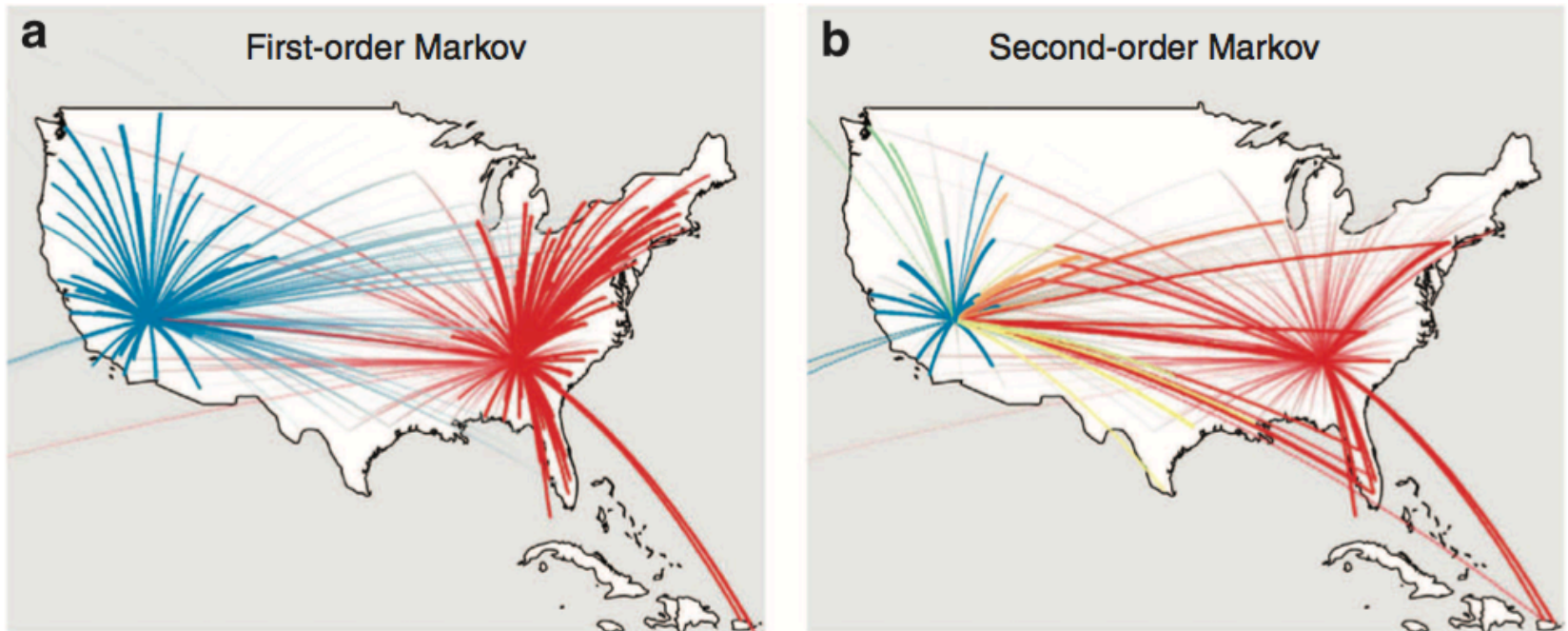
[A Noncooperative Model of Network Formation](#)
Venkatesh Bala - *Econometrica* (2000)

Higher Order Dynamics



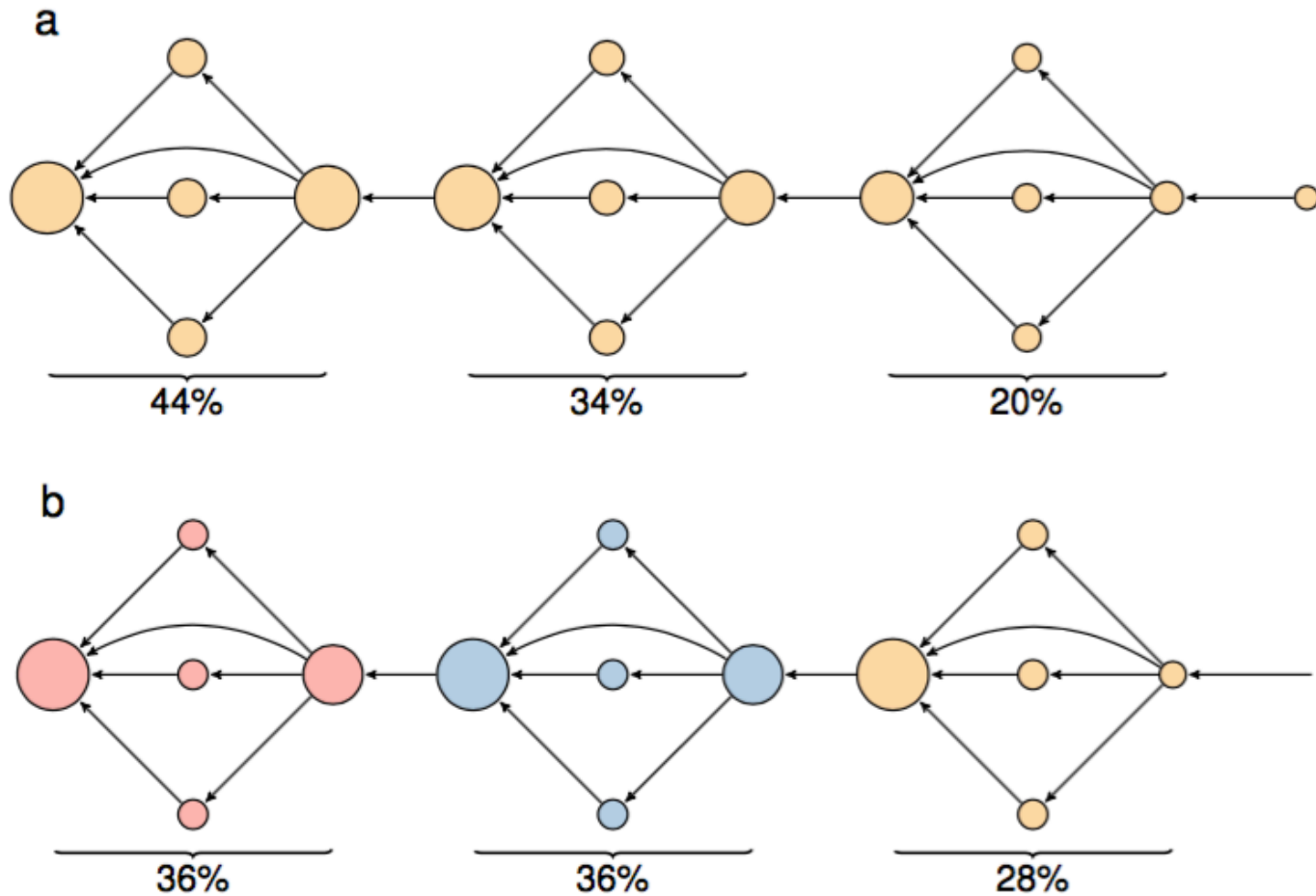
Rosvall et al. (2014) Memory in network flows and its effects on spreading dynamics and community detection. *Nature Communications*

Higher Resolution Maps



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

Article-level Ranking and Mapping



West et al. (2016) Ranking and mapping article-level citation networks. *in prep.*

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

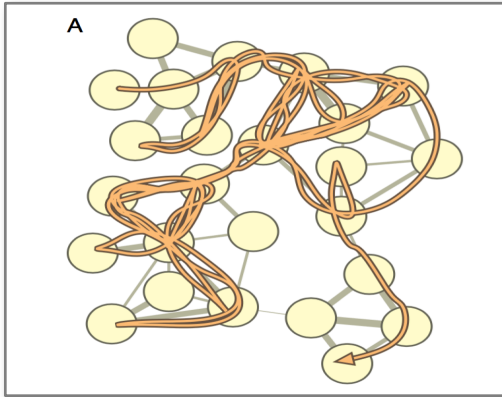
This data is available as a set of zipped text files stored in Microsoft Azure blob storage and available via HTTP. The file size (zipped) is ~30GB and may be downloaded [here](#).

The Challenge

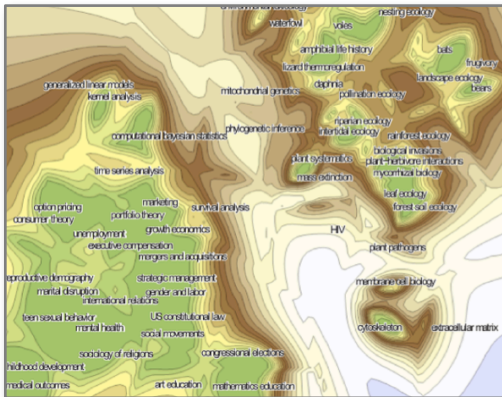
The goal of the Ranker Challenge is to assess the query-independent importance of scholarly articles, using data from the Microsoft Academic Graph.

Wesley-Smith et al. (2016) Static Ranking of Scholarly Papers using Article-Level Eigenfactor (ALEF).

arxiv:1606.08534

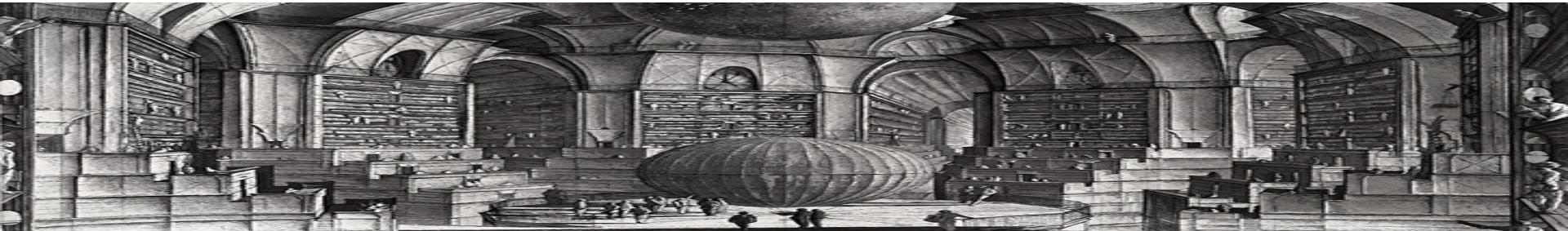


Science of Mapping



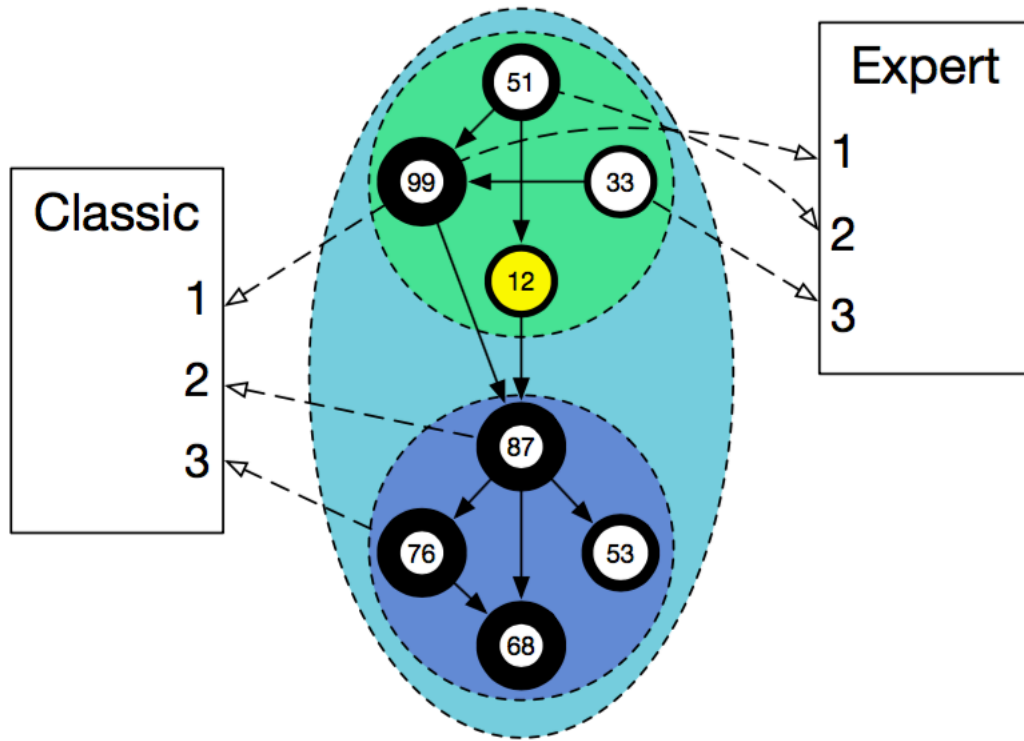
Mapping of Science

Explore the recommendations
babel.eigenfactor.org

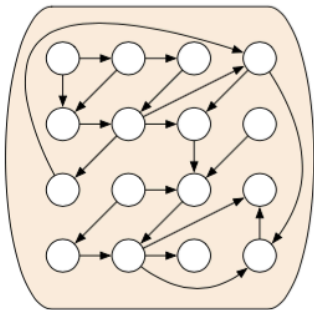


jevinw@uw.edu

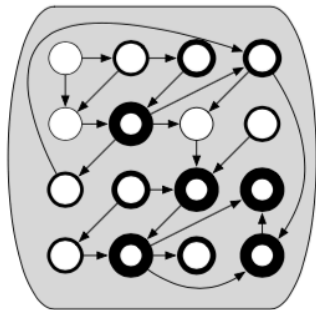
Recommend



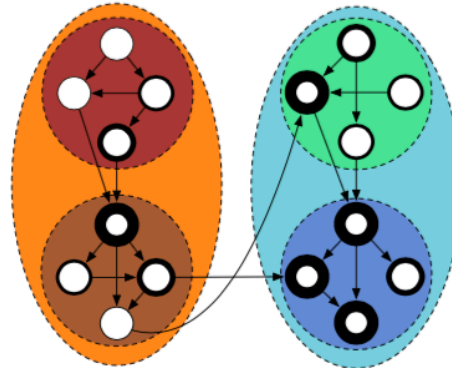
Assemble



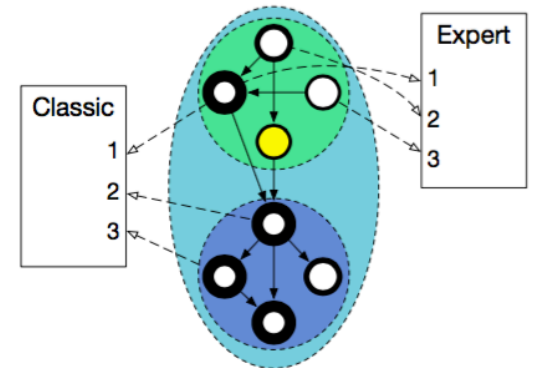
Rank



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*



Babel

- Facilitate research and *implementation* of recommendations
- Bibliographic data at scale
- Freely available and open source
- Evaluation of recommendations
- Audience: publishers, researchers, developers
- API Standardization & Endpoint Discovery

babel.eigenfactor.org



CONSERVATION BIOLOGY OVERVIEW



Conservation is the scientific study of the nature and of Earth's **biodiversity** with the aim of protecting species, their habitats, and **ecosystems** from excessive rates of extinction and the **erosion** of biotic interactions. It is an interdisciplinary subject drawing on natural and social sciences, and the practice of natural resource management. The conservation ethic is based on the findings of conservation biology.

Source: Wikipedia



Source: Wikimedia Commons

Influential Articles

1960s	1970s	1980s	1990s	2000s	2010s	2020
		<ul style="list-style-type: none"> The Canonical Distribution of Commonness and ... An Equilibrium Theory of Insular Zoogeography 	<ul style="list-style-type: none"> Turnover Rates in Insular Biogeography: ... 	<ul style="list-style-type: none"> The Statistics and Biology of the ... 		

Related Topics

- Species
- Habitat conservation

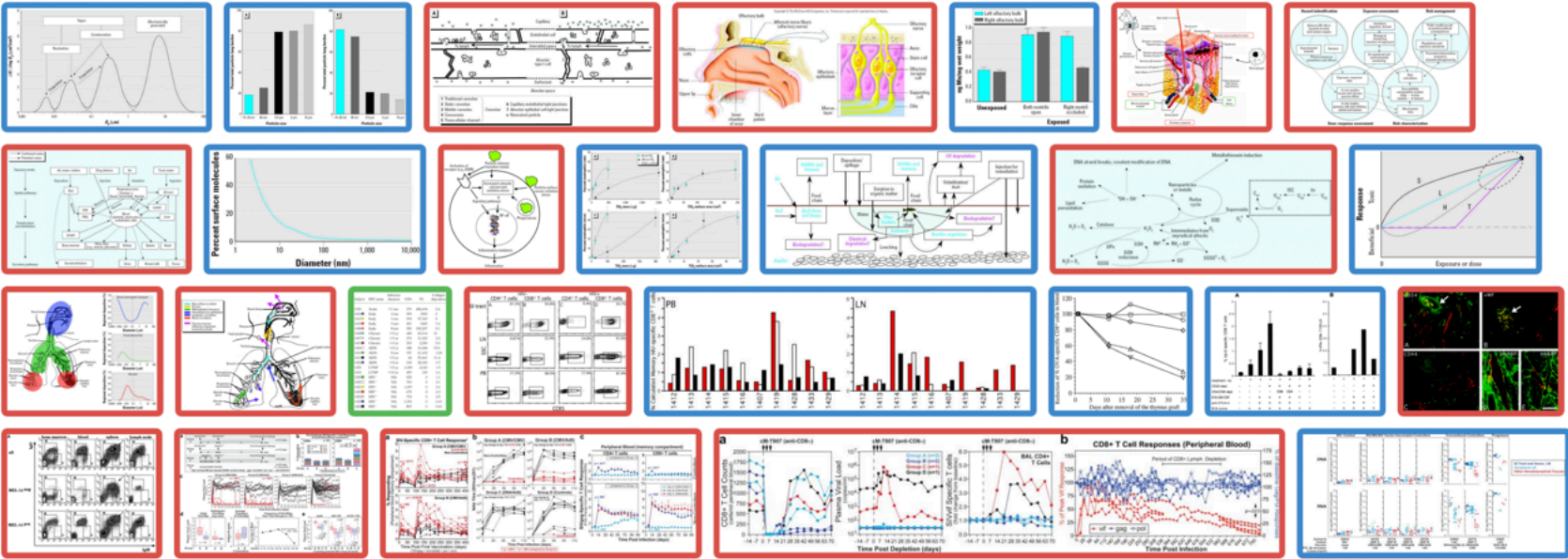
<http://labs.jstor.org/sustainability/topic/Biodiversity>

Figure-Centric Search Engine

Impact blood lymph Search



Composite Equation Diagram Photo Plot Table



Questions

- How do patterns of encoding visual information in the literature vary across disciplines?
- How have patterns of encoding visual information in the literature evolved over time?
- Is there any link between patterns of encoding visual information and scientific impact?

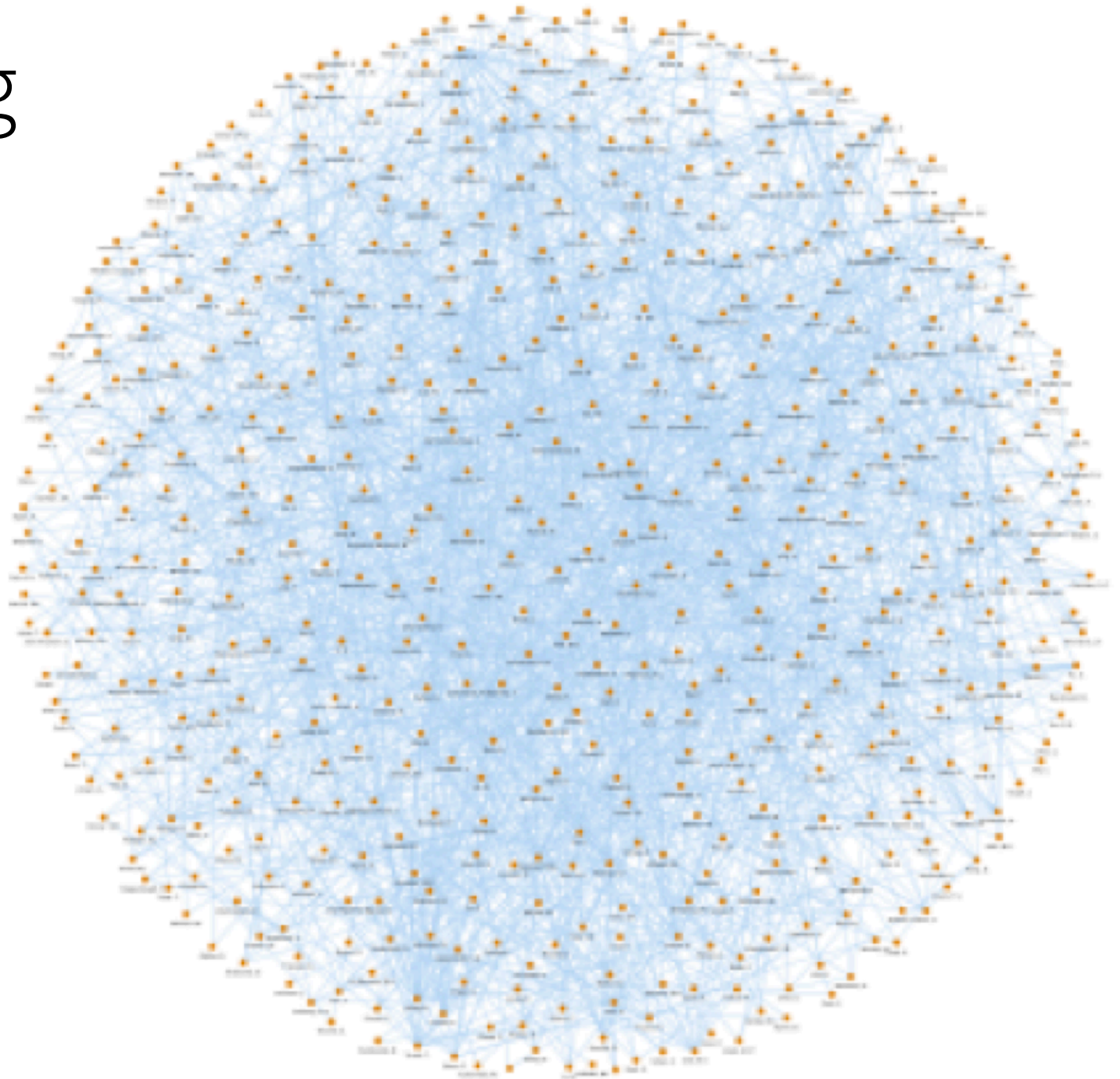
How can we better utilize visual information in the search and navigation process?

Challenges

1. scaling

2. mechanism

Scaling



RelaxMap

A

Scalable and Efficient Flow-Based Community Detection for Large-Scale Graph Analysis

SEUNG-HEE BAE, DANIEL HALPERIN, and JEVIN WEST, University of Washington
MARTIN ROSVALL, Umeå University
BILL HOWE, University of Washington

Community detection is a powerful approach to uncover important structures in large networks. For real networks that often describe the flow of some entity, flow-based community detection methods are particularly important. Infomap is a flow-based community detection algorithm that optimizes the objective function known as the map equation. Third-party benchmarks have found that Infomap is the most effective algorithm for identifying clusters in large graphs. Unfortunately, though Infomap works well, it is an inherently serial algorithm and thus cannot take advantage of multi-core processing in modern computers, limiting its use for analyzing large graphs quickly.

In this paper, we propose a novel algorithm to optimize the map equation called RelaxMap. RelaxMap provides two important improvements over Infomap: parallelization, so that the map equation can be optimized over much larger graphs, and prioritization, so that the most important work occurs first, iterations take less time, and the algorithm converges faster. We implement these techniques using OpenMP on shared-memory multicore systems, and evaluate our approach on a variety of graphs from standard graph clustering benchmarks as well as real graph datasets. Our evaluation shows that both techniques are effective: RelaxMap achieves 70% parallel efficiency on 8 cores, and prioritization improves algorithm performance by an additional 20%–50% in average, depending on the graph properties. Additionally, RelaxMap converges in the similar number of iterations and provides solutions of equivalent quality as the serial Infomap implementation.

Categories and Subject Descriptors: I.5.3 [Clustering]: Algorithms; F.1.2 [Modes of Computation]: Parallelism and Concurrency

General Terms: Design, Algorithms, Performance

Additional Key Words and Phrases: Community Detection, Graph Analysis, Parallelization, Prioritization

ACM Reference Format:

Seung-Hee Bae, Daniel Halperin, Jevin West, Martin Rosvall, and Bill Howe, 2014. Scalable and efficient flow-based community detection for large-scale graph analysis. *ACM Trans. Knowl. Discov. Data.* V, N, Article A (January 2014), 29 pages.

DOI: <http://dx.doi.org/10.1145/0000000.0000000>

1. INTRODUCTION

Community detection in large graphs is emerging as a first-class technique in a number of applications: functional similarity in biological networks [Gavin et al. 2006; Guimera and Amaral 2005], collaboration communities in research networks [Girvan and Newman 2002], and the macro-structure

Summary

- Study the *Science of Science*
- Assemble scholarly knowledge graph into machine readable formats
- Ask questions about the origin and evolution of ideas and fields, interdisciplinarity, impact assessment and sociology of science
- Develop clustering algorithms for automatically detecting scholarly communities in the literature
- Building statistical and visualization tools that improve navigation, make relevant connections and facilitate knowledge discovery
- Challenges: scaling, mechanism
- Eigenfactor.org, Viziometrics.org, Babel.eigenfactor.org

Acknowledgements

Carl Bergstrom, Department of Biology, University of Washington

Martin Rosvall, Department of Physics, Umea University

Ian Wesley-Smith, Information School, University of Washington

Jason Portenoy, Information School, University of Washington

Bill Howe, eScience, CSE, University of Washington

Poshen Lee, CSE, University of Washington

jevinw@uw.edu

jevinwest.org

Eigenfactor.org

@jevinwest