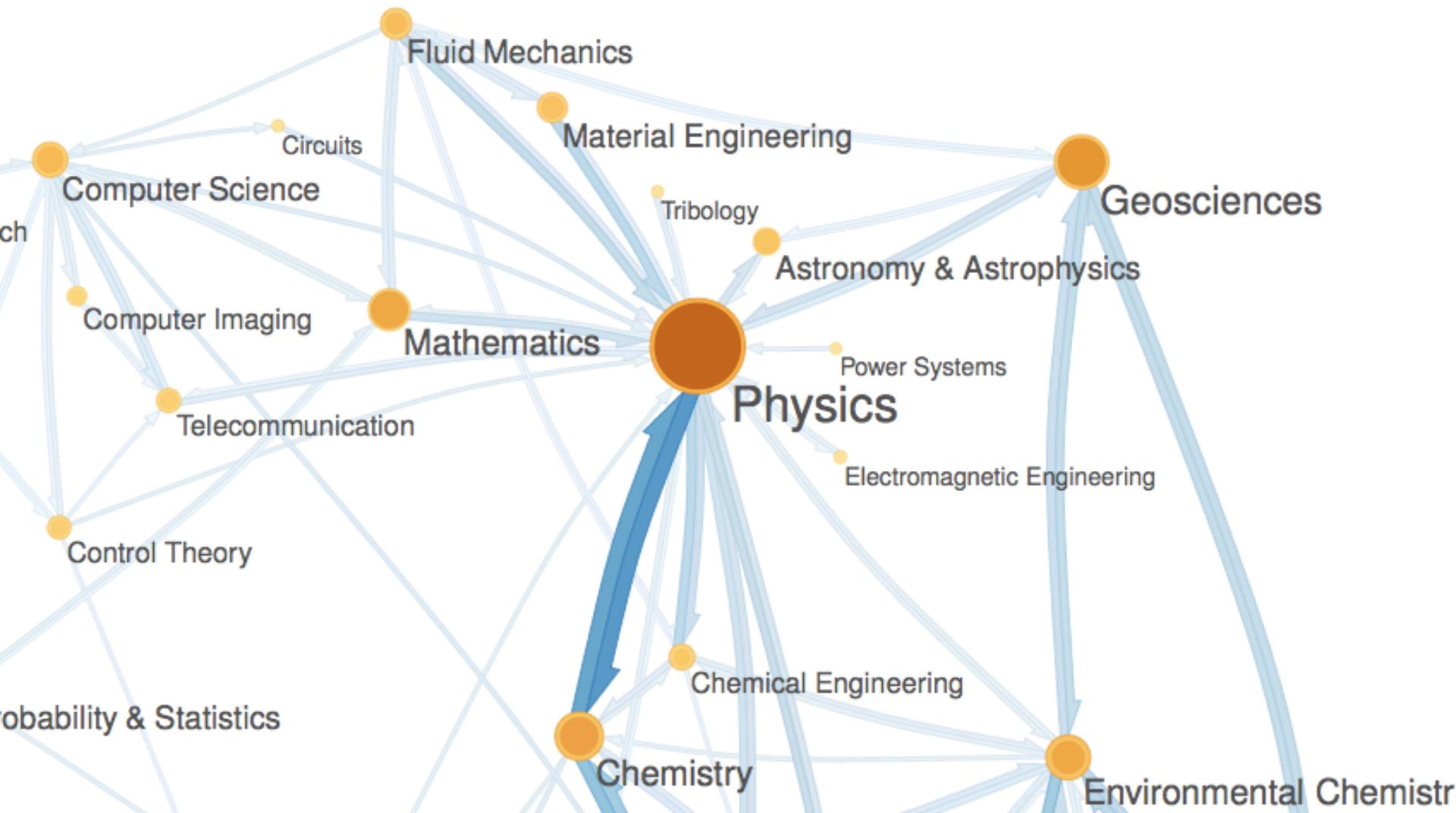


Facilitating discovery with zoomable maps

Jevin West, Information School, University of Washington



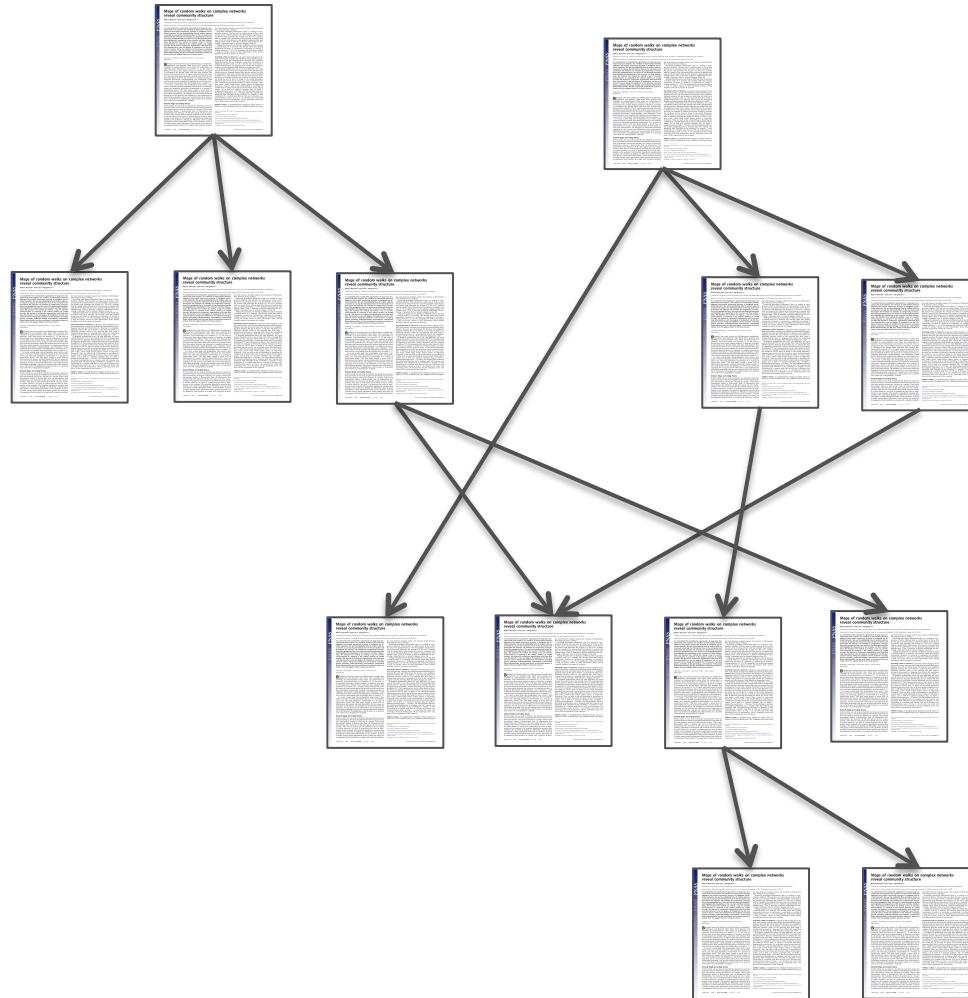
compressing \leftrightarrow finding patterns

5.8 MB (TIFF) \longrightarrow 0.9 MB (TIFF + LZW)



5.8 MB (TIFF) \longrightarrow 2.8 MB (TIFF + LZW)

Citations form a vast network



de Solla Price, Science (1965)



The Scholarly Graph



PatentVector™

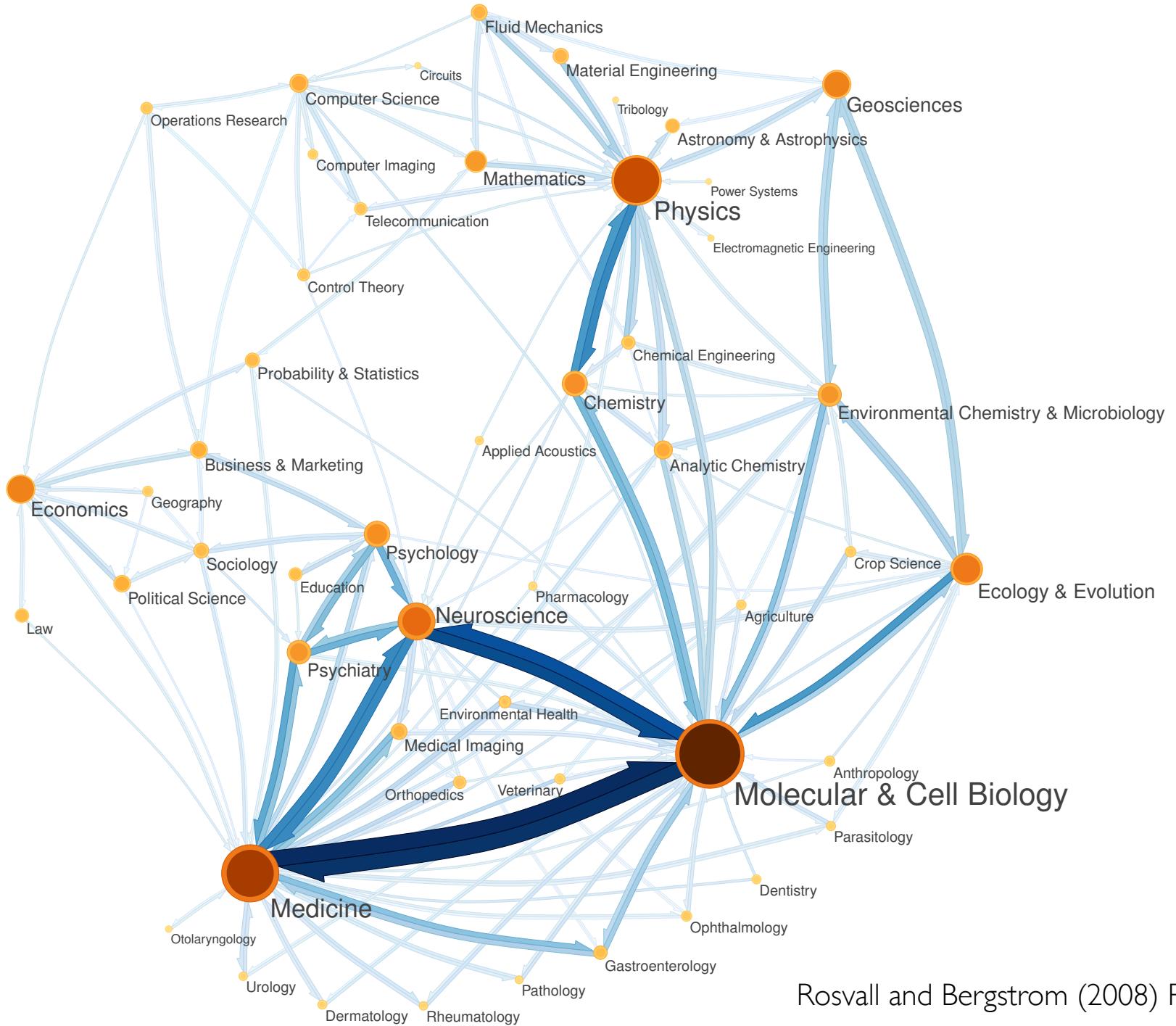


WIKIPEDIA
The Free Encyclopedia



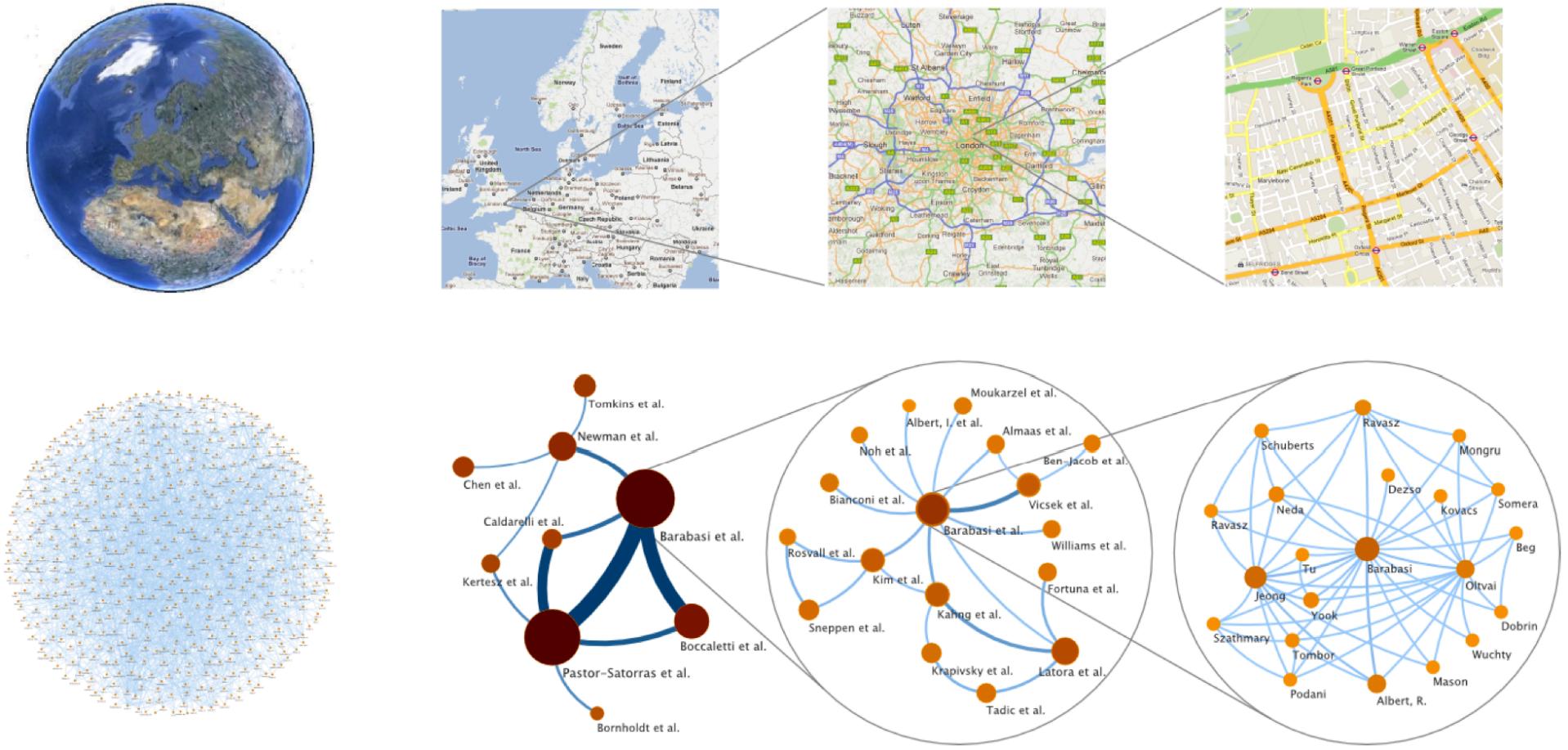
PNAS

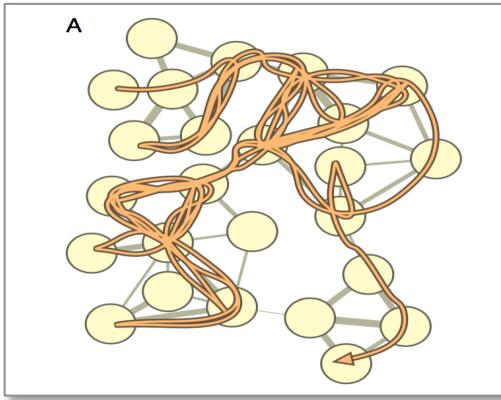




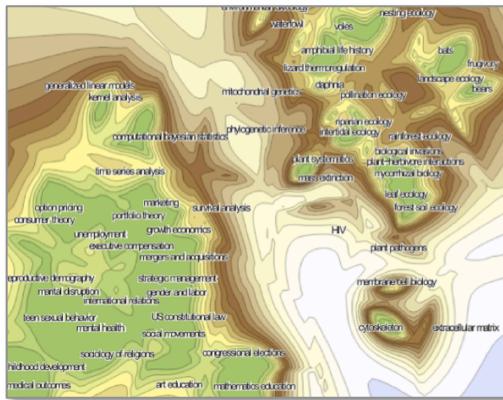
Rosvall and Bergstrom (2008) PNAS

Facilitating discovery with zoomable maps

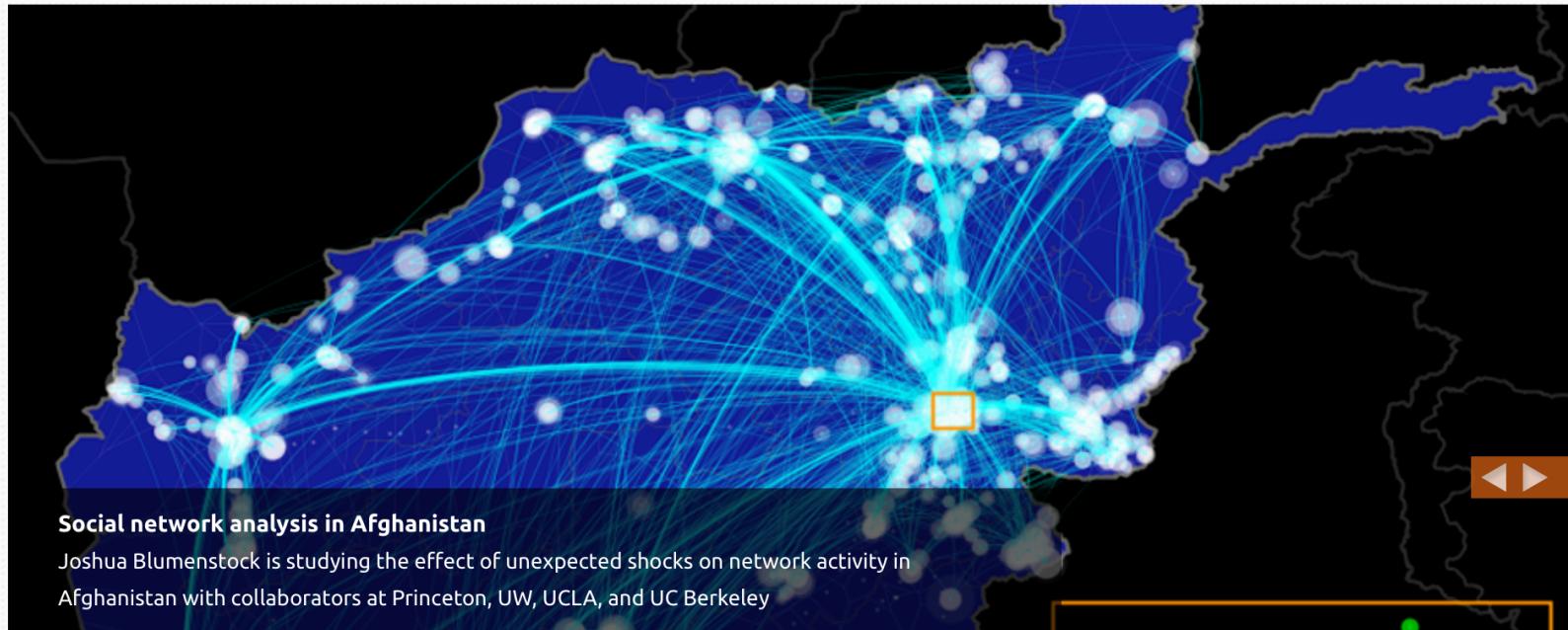




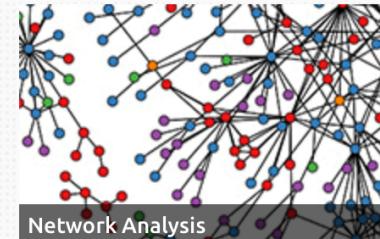
Science of Mapping



Mapping of Science



Research Focus Areas



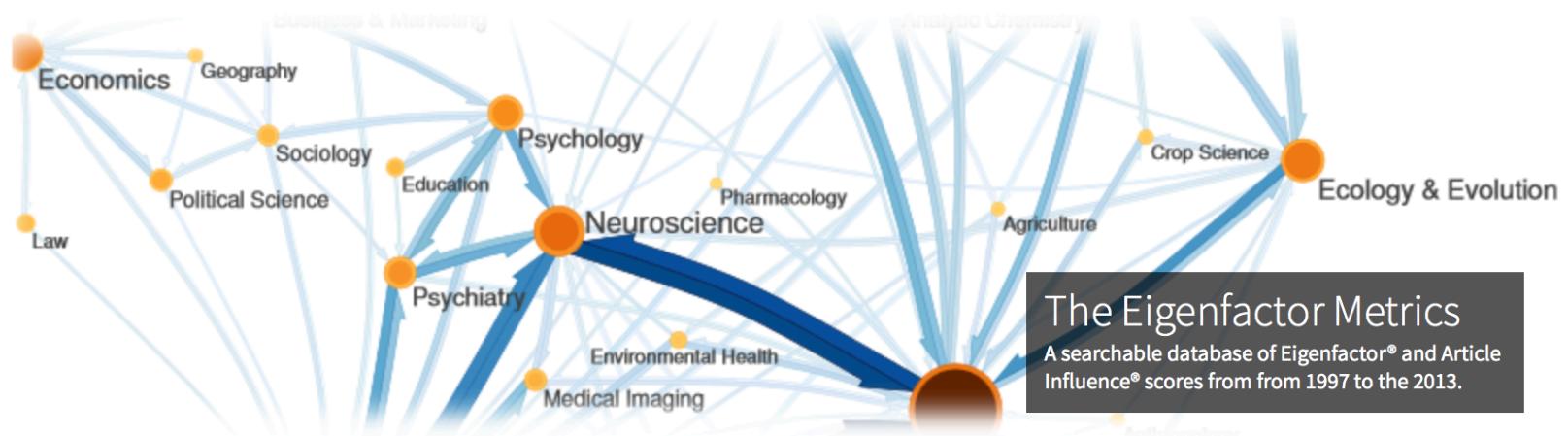
News and Updates

28

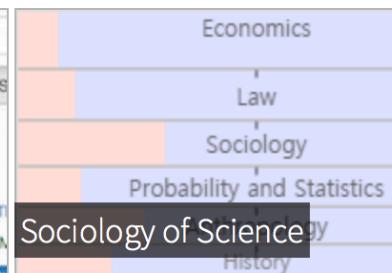
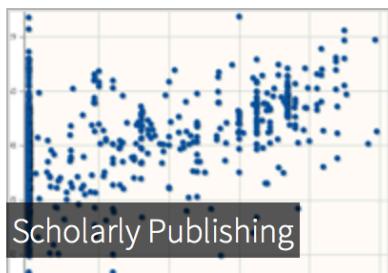
Blumenstock at Population Association of America

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



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 importance of individual articles. This was one of the article-level Eigenfactor algorithms used to rank the

oren etzioni



S

- Statistical Methods For Analyzing Speedup Learning Experiments O Etzioni satisfaction programs
- Face And Computer-Mediated Communities Amitai Etzioni, Oren Etzioni 1998 resources sustained
- Document Clustering O Zamir document clustering
- Communities: Virtual Vs. Real A Etzioni 1996 implications internet
- Statistical Methods For Analyzing Speedup Learning Experiments. O Etzioni 1993 scheduling problems
- Statistical Methods For Analyzing Speedup Learning Experiments O Etzioni 1993 generating abstractions

Get Related



Get Related



Get Related



Get Related



« Previous

1

2

3

4

5

6

7

8

9

10

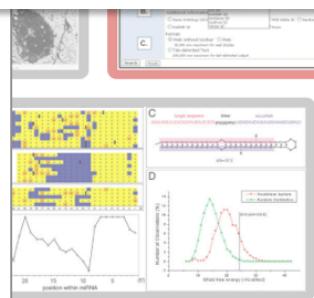
Next »

Papers related to

- Statistical Methods For Analyzing Speedup Learning Experiments O Etzioni satisfaction programs
- Automatically Configuring Constraint Satisfaction Programs: A Case Study S Minton 1995 satisfaction programs
- Abstraction Via Approximate Symmetry T Ellman 1992 satisfaction programs
- Integrating Heuristics For Constraint Satisfaction Problems: A Case Study S Minton 1992 satisfaction programs
- An Analytic Learning System For Specializing Heuristics S Minton 1992 satisfaction programs
- Automated Synthesis Of Constrained Generators W Braudaway 1988 satisfaction programs



Poshen Lee



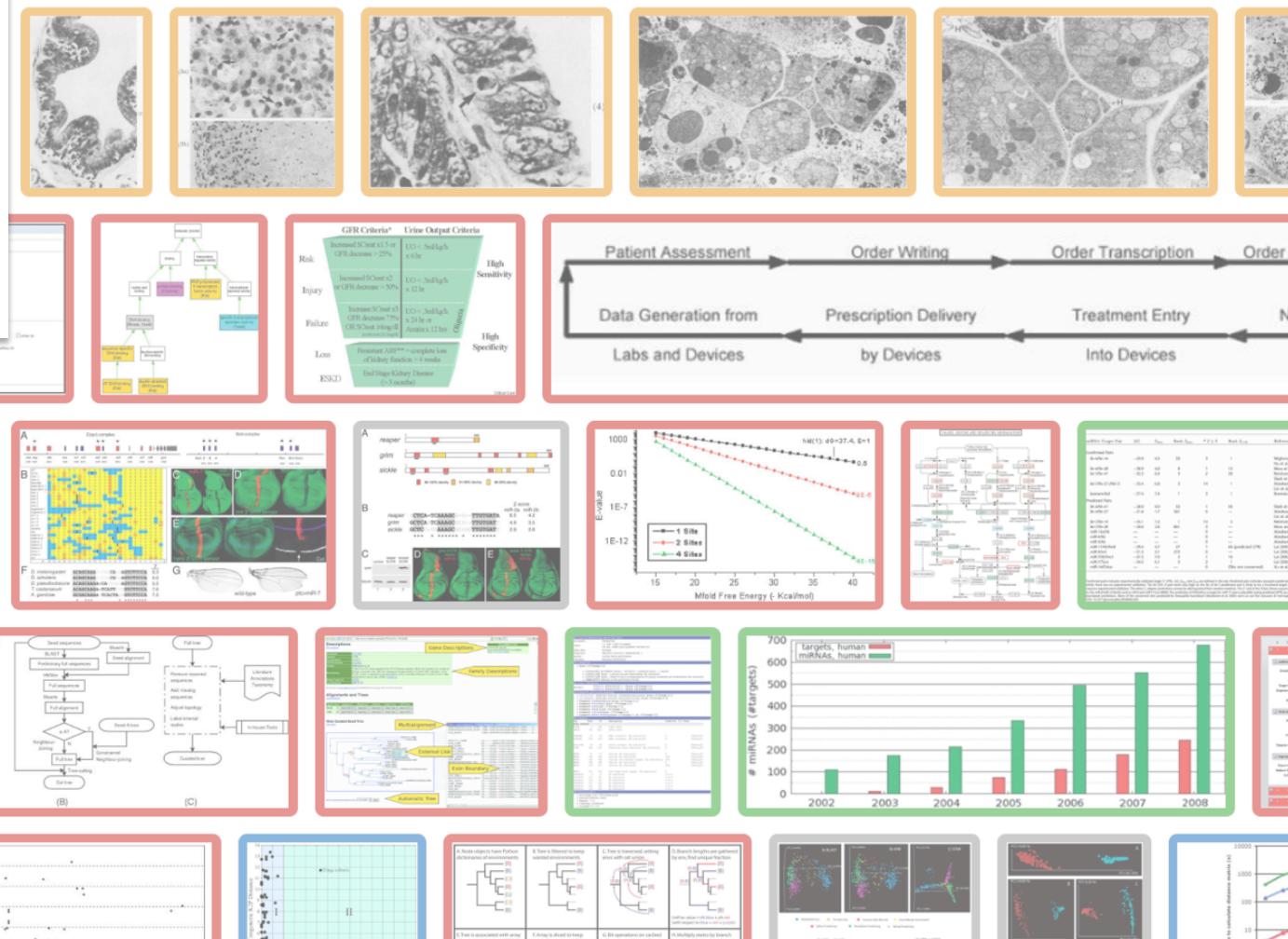
Composite
Equation
Diagram
Photo
Plot
Table

Impact

Keywords or Cluster, Result Ordered by Impact

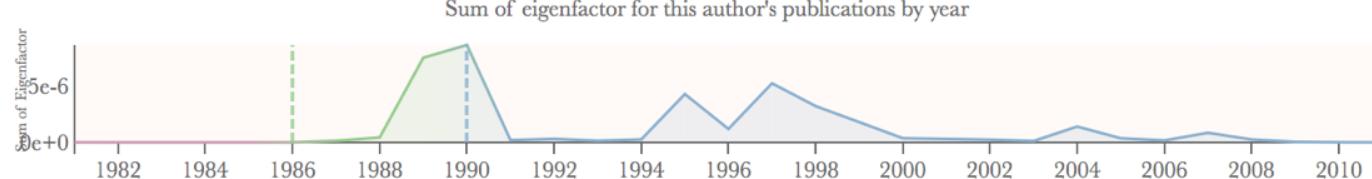
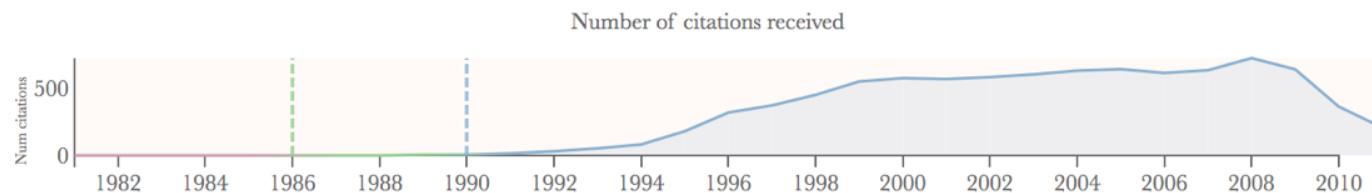
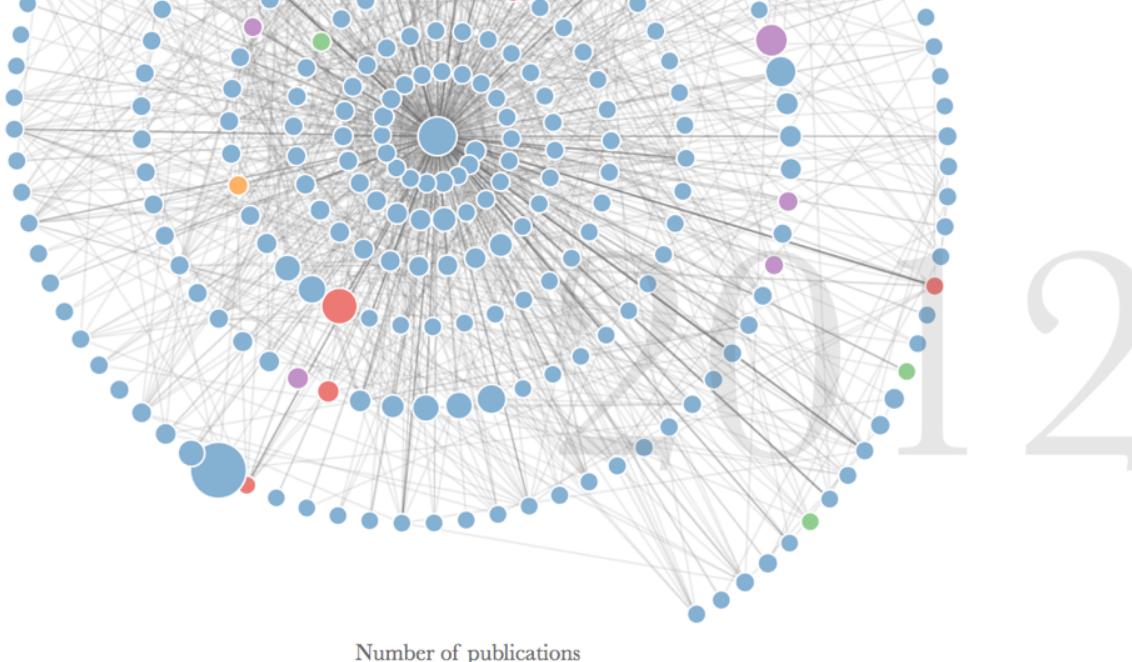
Search

Composite Equation Diagram Photo Plot Table



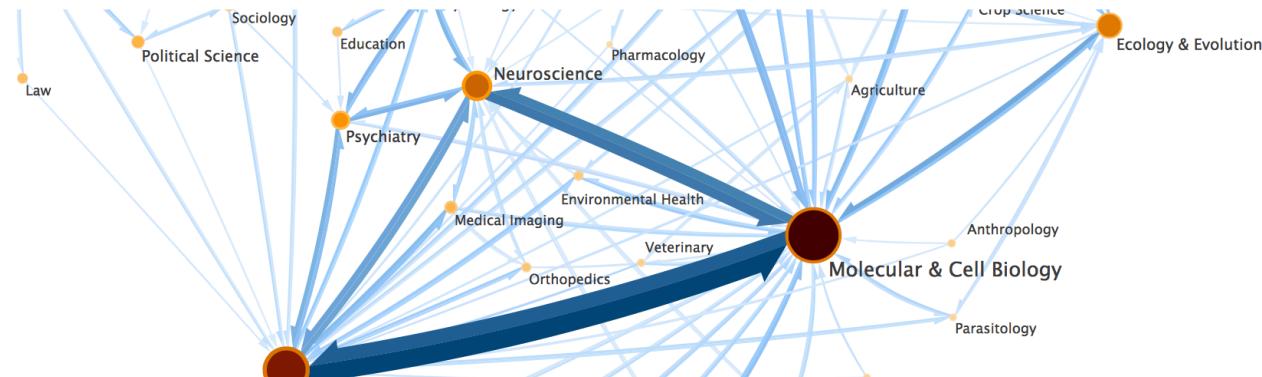


Jason Portenoy

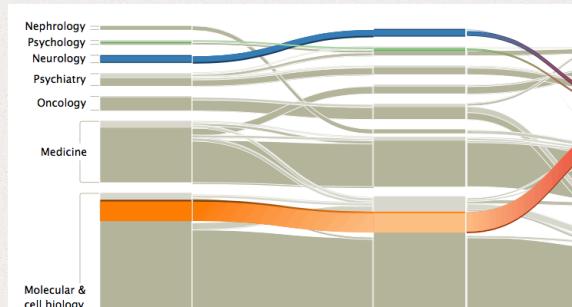


scholar.eigenfactor.org

Simplify and highlight important structures in complex networks



Apps »



Code »

```
using infomath::plog;
for (unsigned int i = 0; i < numNodes; ++i)
{
    enter_log_enter += plog(m_moduleFlowData[i].enterFlow);
    exit_log_exit += plog(m_moduleFlowData[i].exitFlow);
    flow_log_flow += plog(m_moduleFlowData[i].exitFlow);
    enterFlow += m_moduleFlowData[i].enterFlow;
}
enterFlow += exitNetworkFlow;
enterFlow -> enterFlow = nlogCenterFlow);
```

Publications »

Maps of information flow reveal community structure in complex networks

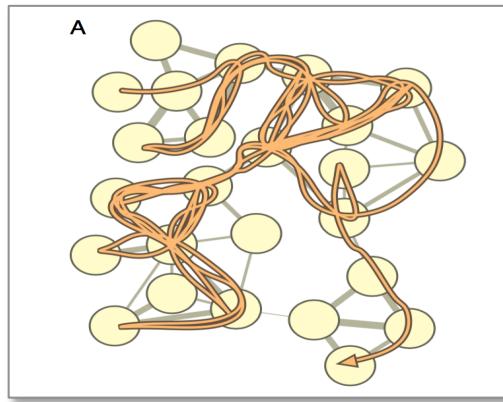
Martin Rosvall and Carl T. Bergstrom
PNAS **105**, 1118 (2008). [arXiv:0707.0609]



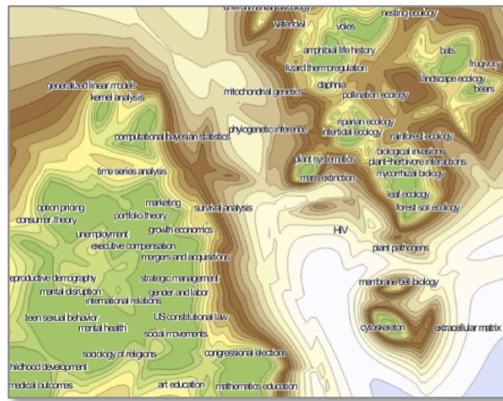
To comprehend the multipartite organization of large-scale biological and social systems, we introduce a new information-theoretic approach to reveal community structure in

News

- August 13, 2015 [Interactive storyboard – Multilevel network sampling](#) – infer network modes from multiple samples
- August 13, 2015 [Interactive storyboard – Higher-order Markov models](#) – identify flows on memory and multilayer networks
- July 23, 2015 [Source code – Infomap](#) – updates to memory and multilayer algorithms



Science of Mapping



Mapping of Science

Data

Compressing



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.

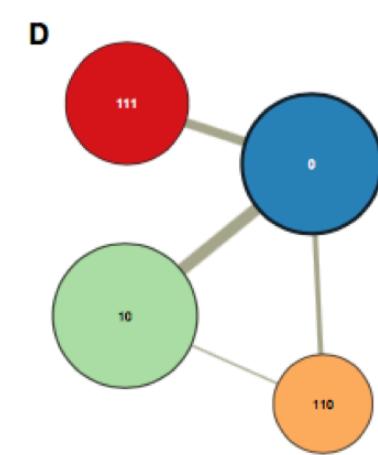
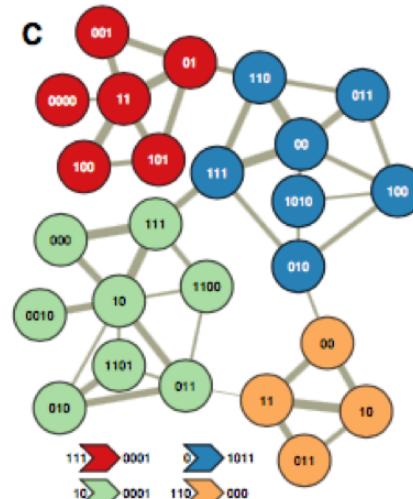
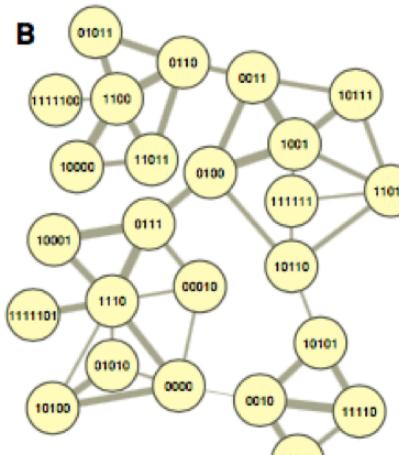
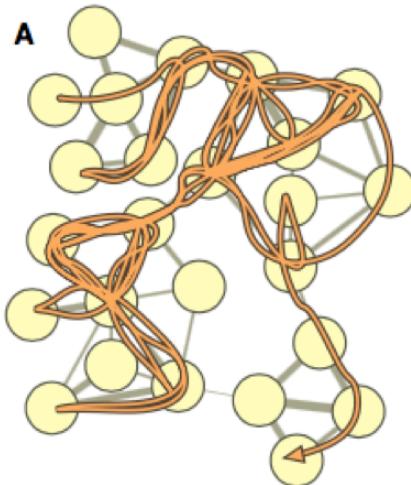
The map equation

frequency of inter-module movements

$$L(M) = q_{\sim} H(Q) + \sum_{i=1}^m p_i^{\circ} H(P^i)$$

The diagram illustrates the components of the map equation. It features a central equation: $L(M) = q_{\sim} H(Q) + \sum_{i=1}^m p_i^{\circ} H(P^i)$. To the left of the first term, a blue vertical bar is labeled "frequency of inter-module movements". To the right of the second term, a red vertical bar is labeled "frequency of movements within module i ". Below the first term, another blue vertical bar is labeled "code length of module names". Below the second term, a red vertical bar is labeled "code length of node names in module i ".

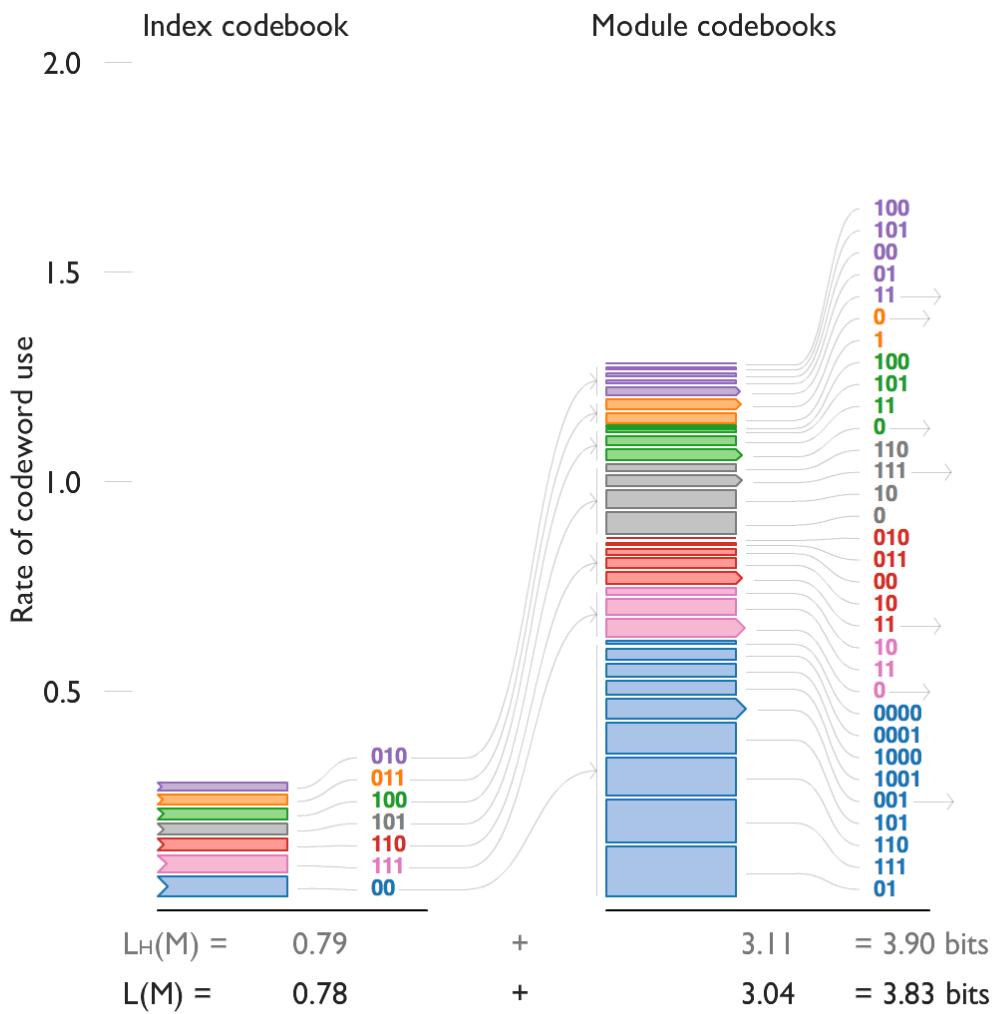
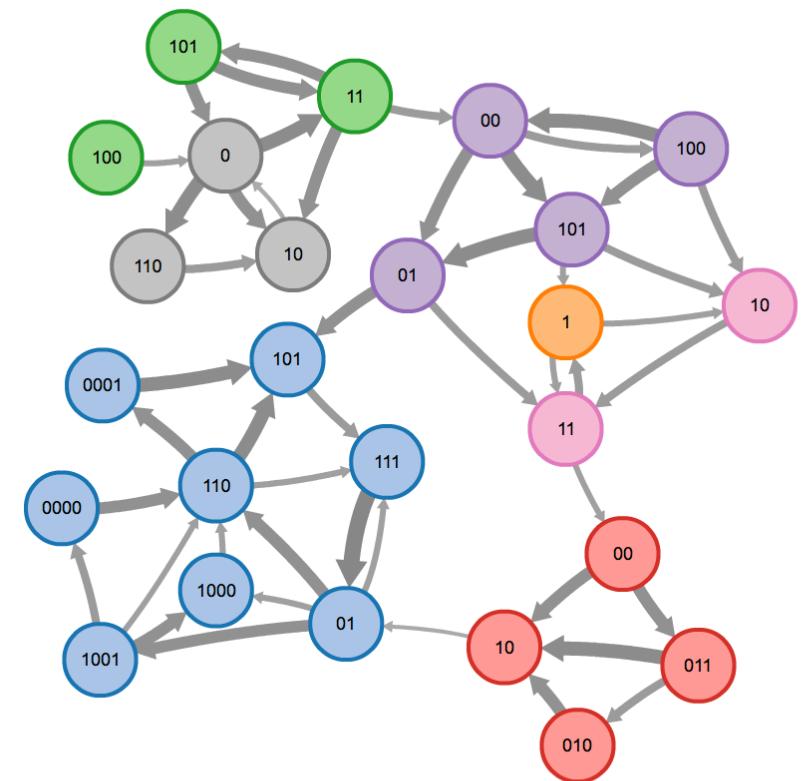
Finding regularities in the dynamics on networks



1111100 1100 0110 11011 10000 11011 0110 0011 10111 1001 0011
1001 0100 0111 10001 1110 0111 1100 0000 1110 10001 0111
0111 1110 0111 1110 1111001 1110 0000 10100 0000 1110 10001 0111
0100 10110 11010 10111 1001 0100 1001 10111 1001 0100 1001 0100
0011 0100 0011 0110 11011 0110 0011 0100 1001 10111 0011 0100
0111 10001 1110 10001 0111 0100 10110 111111 10101 11110
00011

111 0000 11 01 101 100 101 01 0001 0 110 011 00 110 00 111 10111 10
111 000 10 111 000 111 10 011 10 000 111 10 111 10 0010 10 011 010
011 10 000 111 000 0 111 010 100 011 00 111 00 011 00 111 00 111
110 111 110 1011 111 01 101 01 0001 0 110 111 00 011 110 111 1011
10 111 000 10 000 111 0001 0 111 010 1010 010 1011 110 00 10 011
00 111 000 10 000 111 0001 0 111 010 1010 010 1011 110 00 10 011

111 0000 11 01 101 100 101 01 0001 0 110 011 00 110 00 111 10111 10
111 000 10 111 000 111 10 011 10 000 111 10 111 10 0010 10 011 010
011 10 000 111 000 0 111 010 100 011 00 111 00 011 00 111 00 111
110 111 110 1011 111 01 101 01 0001 0 110 111 00 011 110 111 1011
10 111 000 10 000 111 0001 0 111 010 1010 010 1011 110 00 10 011



Performance

PHYSICAL REVIEW E **80**, 056117 (2009)

Community detection algorithms: A comparative analysis

Andrea Lancichinetti^{1,2} and Santo Fortunato¹

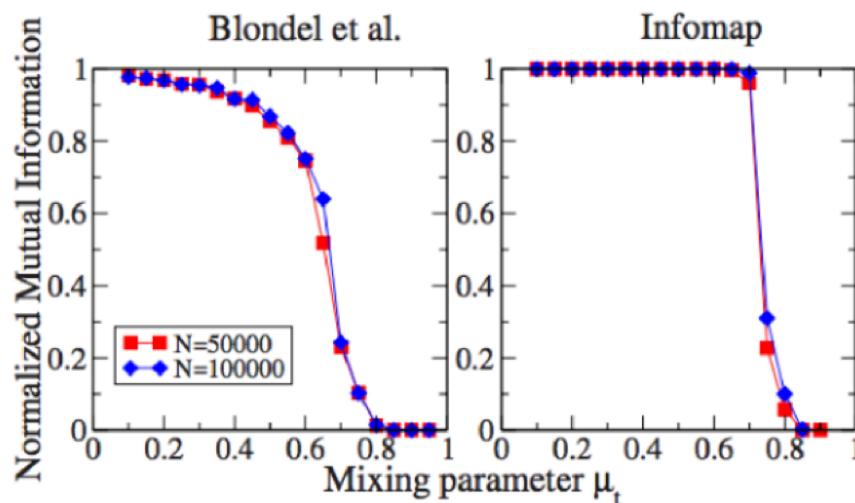
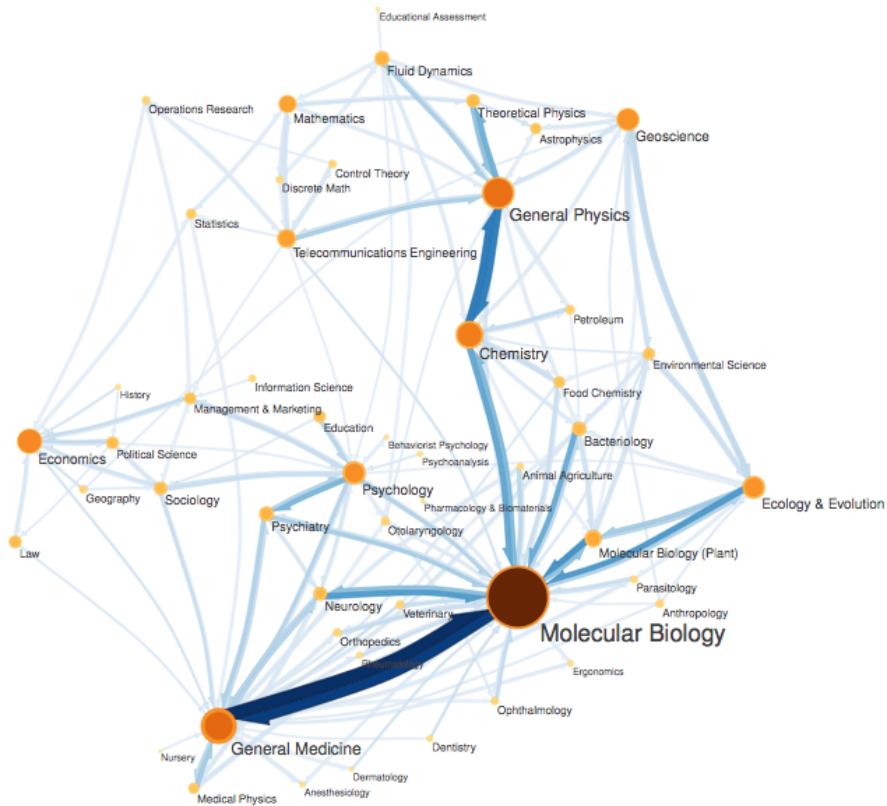
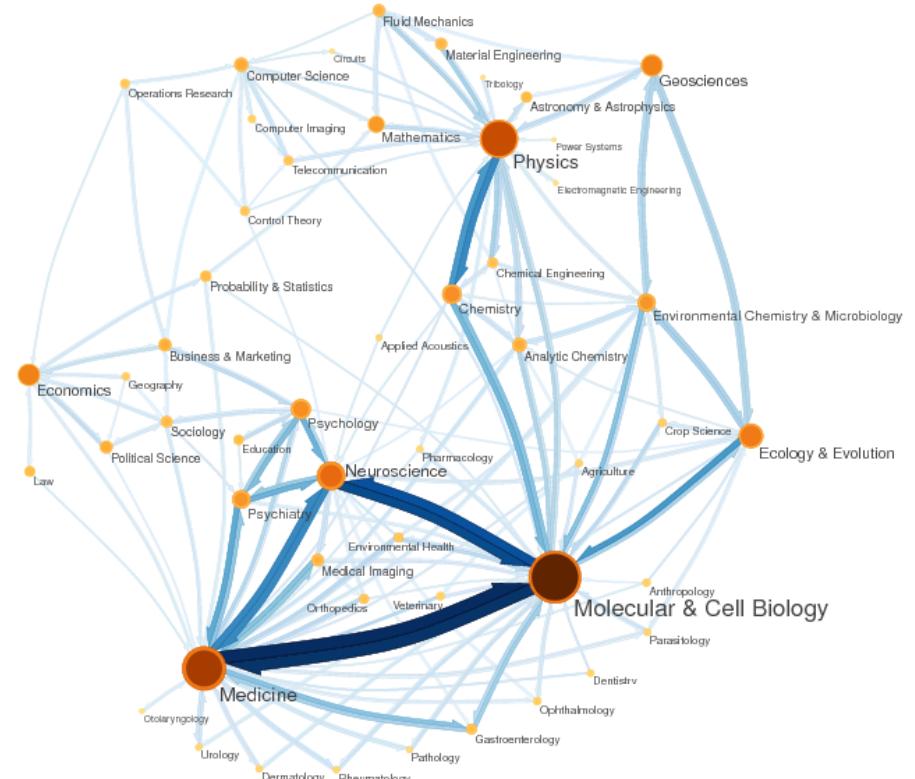


FIG. 3. (Color online) Tests of the algorithm by Blondel *et al.* and Infomap on large LFR benchmark graphs with undirected and unweighted links.

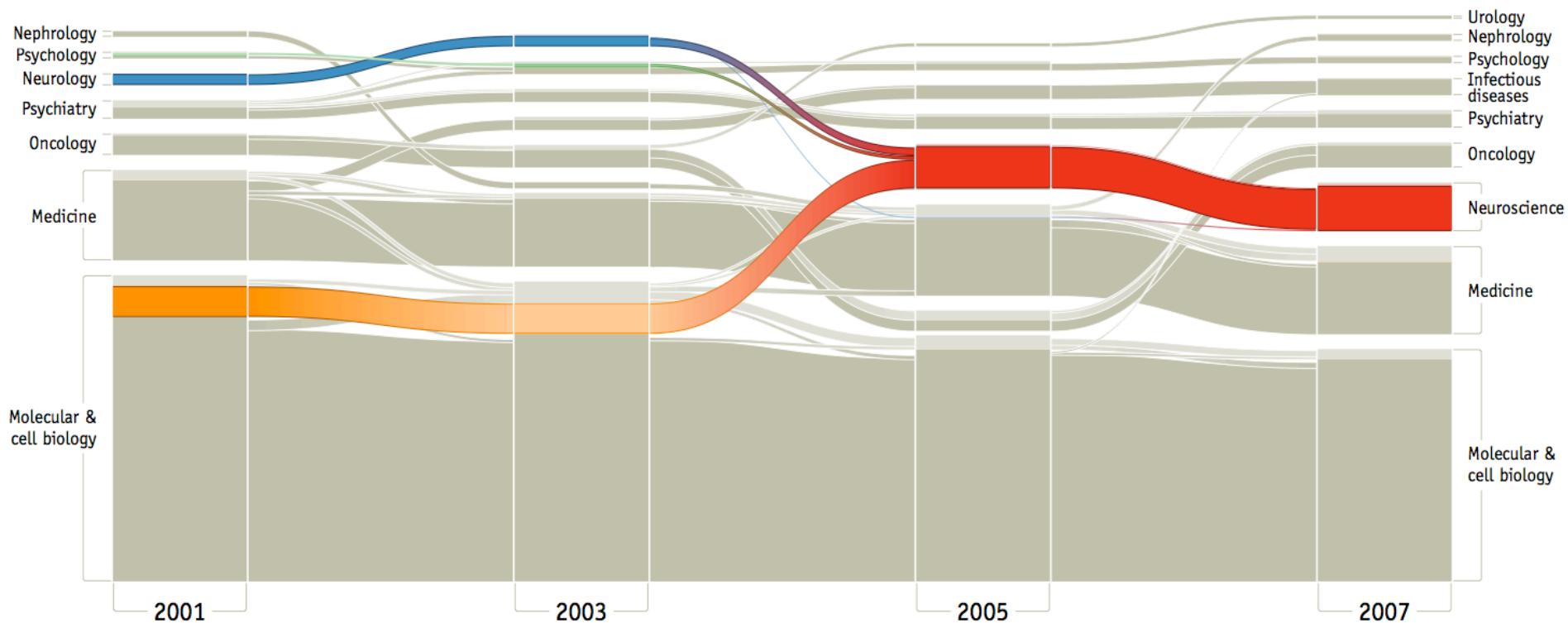
1995



2004

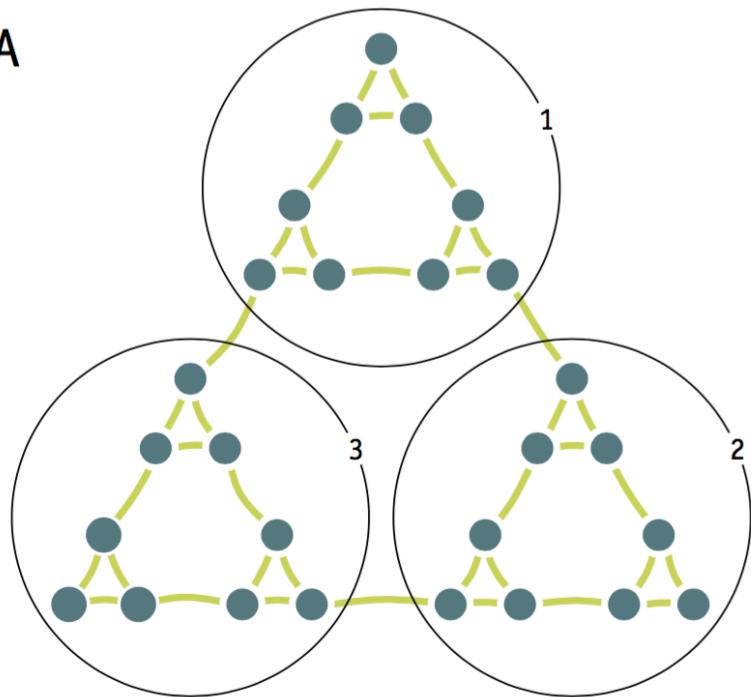


The emergence of Neuroscience

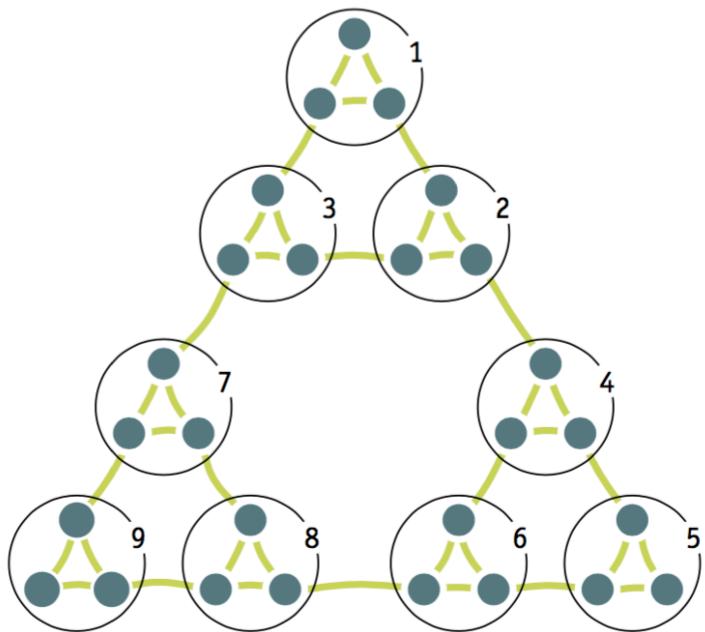


Two-level partitions with the map equation

A



B



$$L(\mathbf{M}) = q_{\cap} H(\mathcal{Q}) + \underbrace{\left\{ p_{\cup}^1 H(\mathcal{P}^1) \atop p_{\cup}^2 H(\mathcal{P}^2) \atop p_{\cup}^3 H(\mathcal{P}^3) \right\}}_{3.56 \text{ bits}} = 3.68 \text{ bits.}$$

0.12 bits 3.56 bits

$$L(\mathbf{M}) = q_{\cap} H(\mathcal{Q}) + \underbrace{\left\{ p_{\cup}^1 H(\mathcal{P}^1) \atop p_{\cup}^2 H(\mathcal{P}^2) \atop p_{\cup}^3 H(\mathcal{P}^3) \atop p_{\cup}^4 H(\mathcal{P}^4) \atop p_{\cup}^5 H(\mathcal{P}^5) \atop p_{\cup}^6 H(\mathcal{P}^6) \atop p_{\cup}^7 H(\mathcal{P}^7) \atop p_{\cup}^8 H(\mathcal{P}^8) \atop p_{\cup}^9 H(\mathcal{P}^9) \right\}}_{2.60 \text{ bits}} = 3.57 \text{ bits.}$$

0.97 bits 2.60 bits

Advantages of InfoMap

- Outperforms most other competing algorithms in accuracy benchmark tests
- Hierarchical version from first principles
- Fast, scalable codebase in C++
- Multiplex versions
- Continual development
- Visual tools associated with the code
- Simple information theoretic framework

The background of the image features a series of thin, wispy blue lines that curve and flow across the frame, resembling smoke or a stylized ribbon. These lines are more concentrated on the left side and spread out towards the right. The overall effect is dynamic and organic.

Dynamics

Journal Ranking

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

Matrix representing the random walk over citations

Probability of not teleporting

Cross-citation Matrix dictating the structure of the citation network

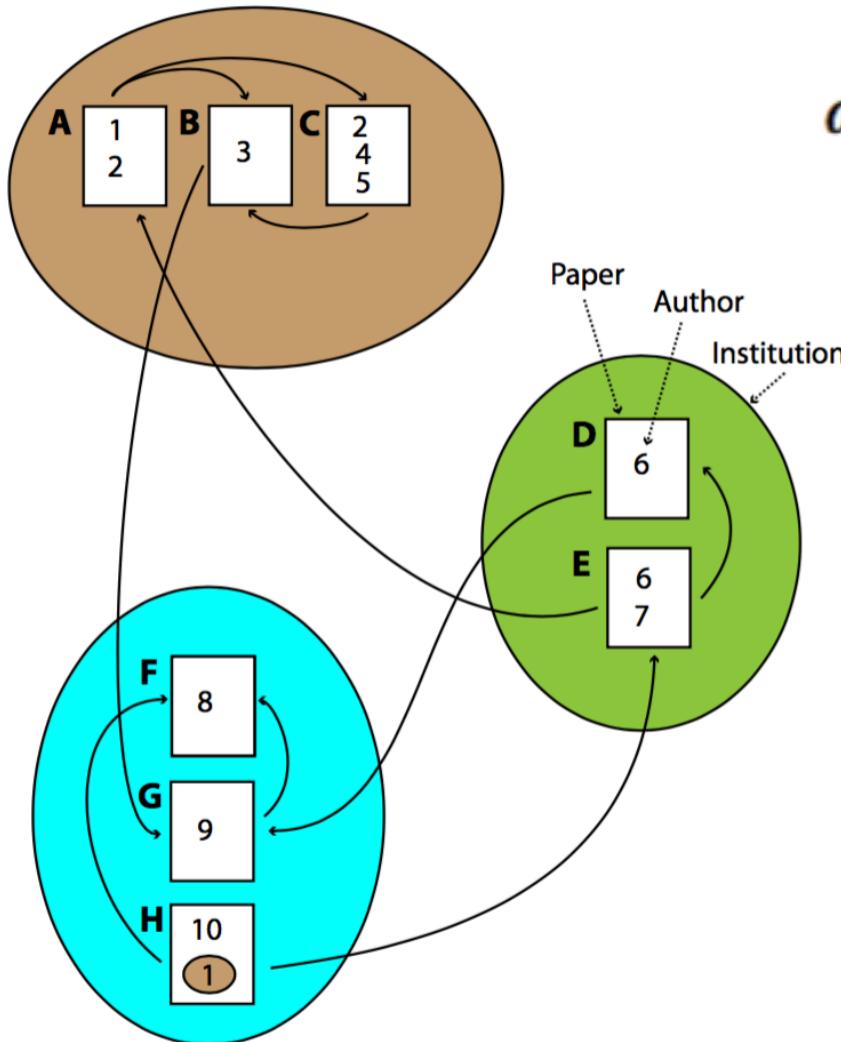
Probability of teleporting to completely new journal weighted by the number of articles in that journal

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

Leading eigenvector of the random walk matrix P .

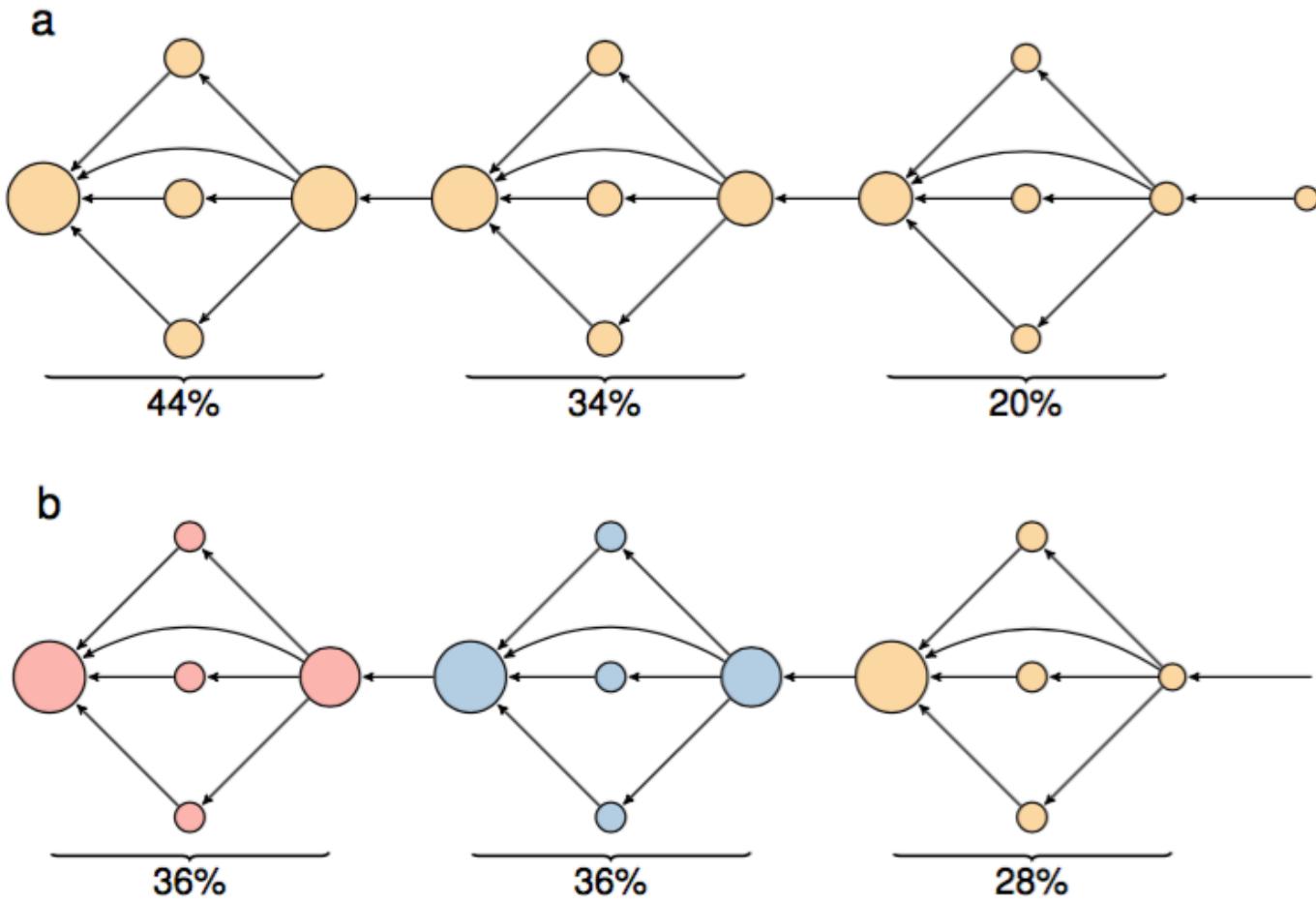
Normalization

Author-level Ranking

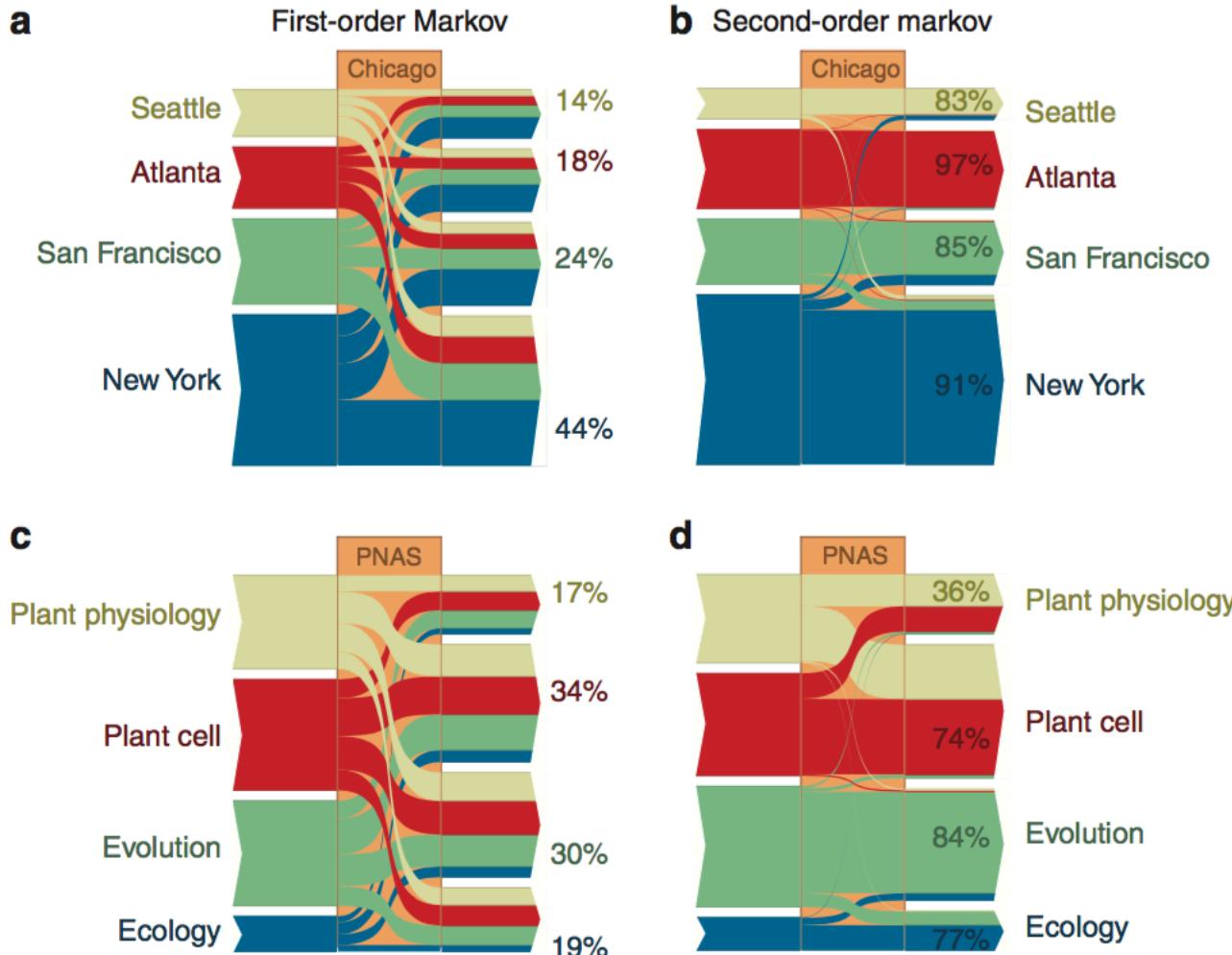


$$\omega = \frac{1}{c(X)} \frac{1}{m} \frac{1}{n}$$

Article-level Ranking



Higher Order Dynamics



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

WSDM CUP CHALLENGE

SIGN-UPS FOR THE WSDM CUP CHALLENGE ARE NOW CLOSED

The Graph

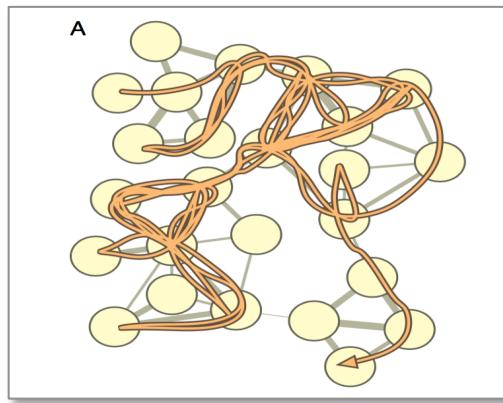
The Microsoft Academic Graph is a heterogeneous graph containing scientific publication records, citation relationships between publications, as well as authors, institutions, journal and conference "venues," and fields of study.

The Data

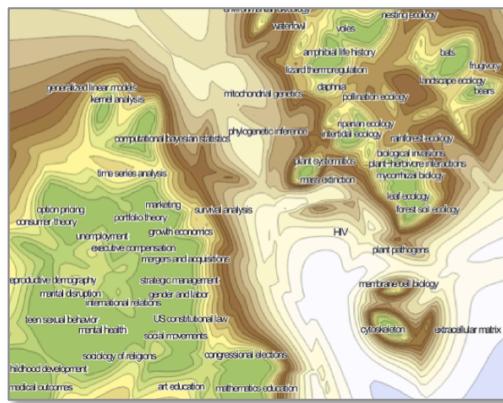
This data is available as a set of zipped text files stored in Microsoft Azure blob storage and available via HTTP. The file size (zipped) is ~30GB and may be downloaded [here](#).

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.



Science of Mapping



Mapping of Science

What can you do with citation maps?

Recommend Articles

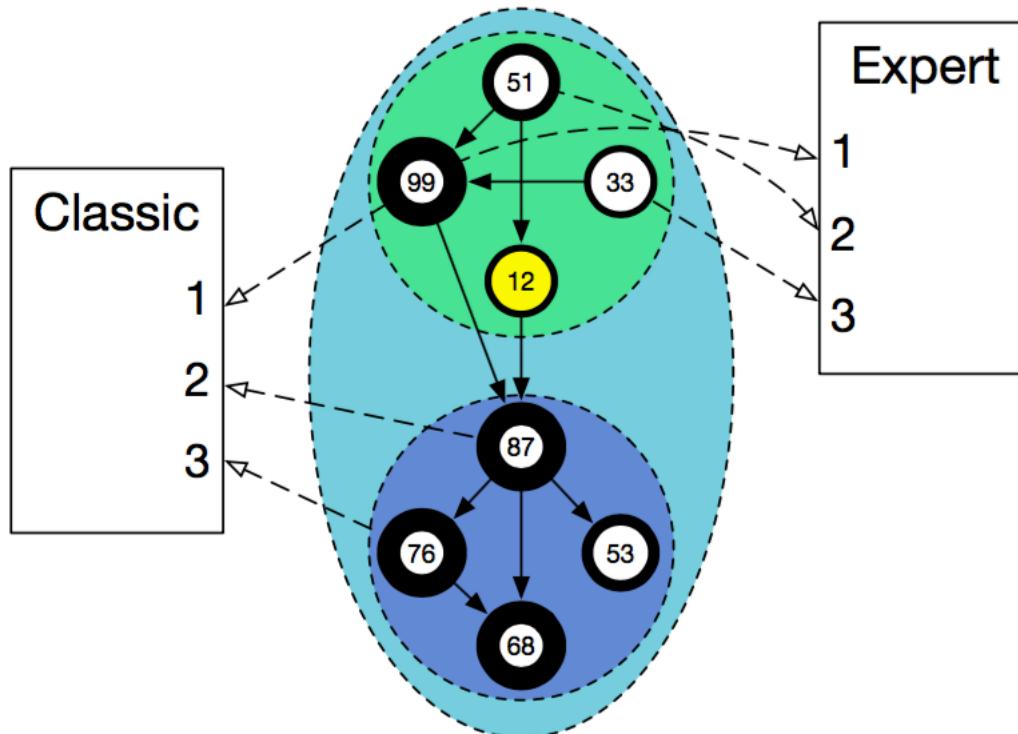
Study the Science of Science

Author Evaluation

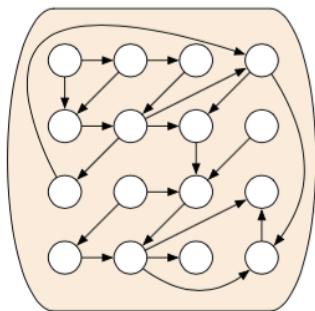
Search and Discovery

Visual Interfaces

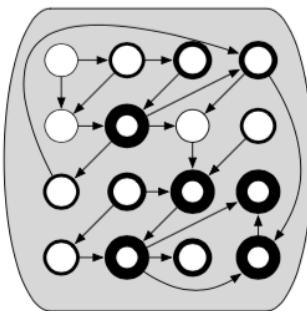
Recommend



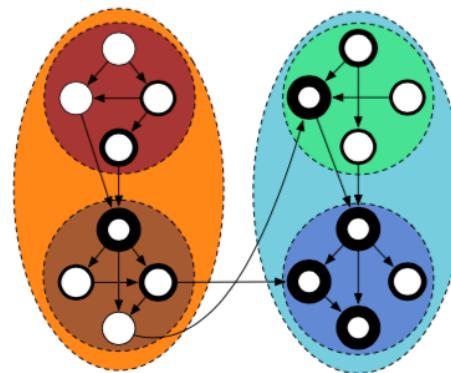
Assemble



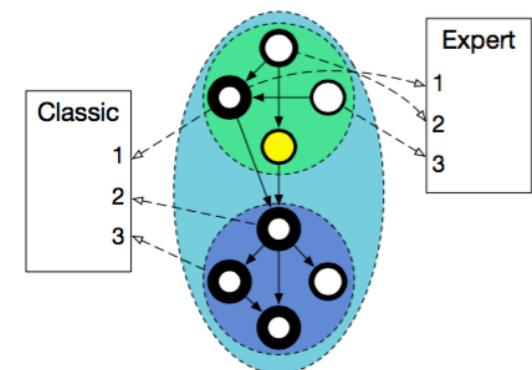
Rank



Cluster



Recommend



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

oren etzioni



S

- Statistical Methods For Analyzing Speedup Learning Experiments O Etzioni satisfaction programs
- Face And Computer-Mediated Communities Amitai Etzioni, Oren Etzioni 1998 resources sustained
- Document Clustering O Zamir document clustering
- Communities: Virtual Vs. Real A Etzioni 1996 implications internet
- Statistical Methods For Analyzing Speedup Learning Experiments. O Etzioni 1993 scheduling problems
- Statistical Methods For Analyzing Speedup Learning Experiments O Etzioni 1993 generating abstractions

Get Related



Get Related



Get Related



Get Related



« Previous

1

2

3

4

5

6

7

8

9

10

Next »

Papers related to

- Statistical Methods For Analyzing Speedup Learning Experiments O Etzioni satisfaction programs
- Automatically Configuring Constraint Satisfaction Programs: A Case Study S Minton 1995 satisfaction programs
- Abstraction Via Approximate Symmetry T Ellman 1992 satisfaction programs
- Integrating Heuristics For Constraint Satisfaction Problems: A Case Study S Minton 1992 satisfaction programs
- An Analytic Learning System For Specializing Heuristics S Minton 1992 satisfaction programs
- Automated Synthesis Of Constrained Generators W Braudaway 1988 satisfaction programs



Scholar Mirrors

Bibliometrics, Scientometrics, Informetrics, Webometrics, and Altmetrics
in Google Scholar Citations, ResearcherID, Researchgate, Mendeley, and Twitter

[HOME](#)[ABOUT](#)[METHODOLOGY](#)[OUR TEAM](#)[OTHER PROJECTS](#)[AUTHORS](#)[DOCUMENTS](#)[JOURNALS](#)[PUBLISHERS](#)

General overview

Displaying core authors 1-20 of 398. Sorted by GS citations (last 5 years), decreasingly.

 Check to display related authors as well Search an author

Name	Online presence	Google Scholar +		ResearcherID +		ResearchGate +		Mendeley +		Twitter +	
		Citations	H Index	Citations	H Index	RG Score	Downloads	Readers	Followers	Tweets	Followers
Loet Leydesdorff	 	26484	73	6444	44	45.14	32165	0	11	84	375
Eugene Garfield*	 	22622	55	8790	153	-	-	-	-	-	-
Mike Thelwall	 	13840	61	3593	32	42.64	24989	7423	36	85	522
Derek J. de Solla Price	 	13263	33	-	-	-	-	-	-	-	-
Francis Narin	 	11297	45	-	-	32.38	795	-	-	-	-
Wolfgang Glänzel	 	10796	54	4924	38	41.16	10572	-	-	-	-
Ronald Rousseau	 	9570	42	NA	NA	42.75	8066	-	-	-	-
Chaomei Chen	 	9512	43	1740	20	34.65	31579	965	3	67	65
Anthony (Ton) F.J. van Raan	 	9200	53	-	-	38.47	6014	-	-	58	166
Ben R Martin	 	8975	39	-	-	-	-	-	-	-	-
András Schubert	 	8655	45	4121	31	39.24	1962	-	-	-	-
Peter Ingwersen	 	8356	35	NA	NA	30.64	8600	-	-	-	-
Henk F. Moed	 	8256	46	-	-	-	-	-	-	-	-
Blaise Cronin	 	7347	43	-	-	33.9	1891	-	-	-	-
Henry Small	 	7307	32	3360	23	-	-	-	-	-	-

Visualizing Scholarly Influence Over Time

Influence of Pew Scholars

Roberta A. Gottlieb

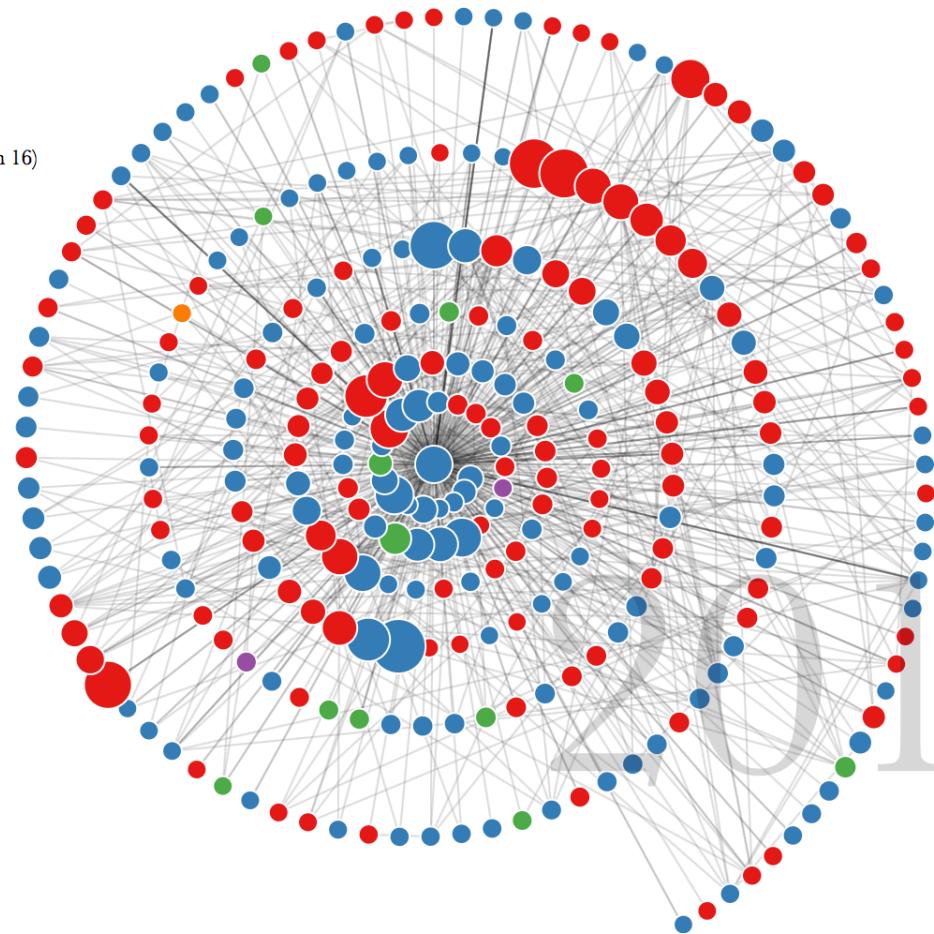
[Learn More](#)

- █ Papers in category "Medicine" (domain 6)
- █ Papers in category "Biology" (domain 4)
- █ Papers in category "Chemistry" (domain 5)
- █ Papers in category "Unknown" (domain 0)
- █ Papers in category "Agriculture Science" (domain 16)

Roberta A.
Gottlieb

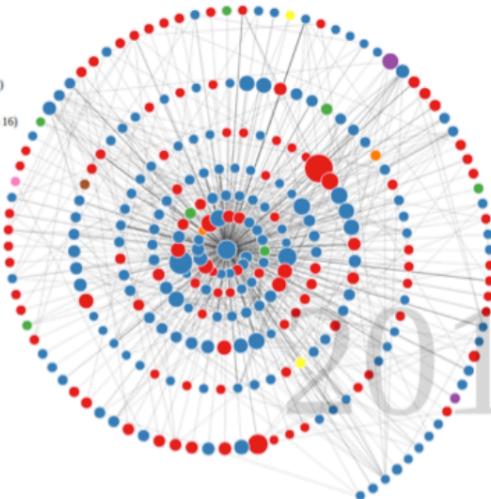


Pew Scholar
1997



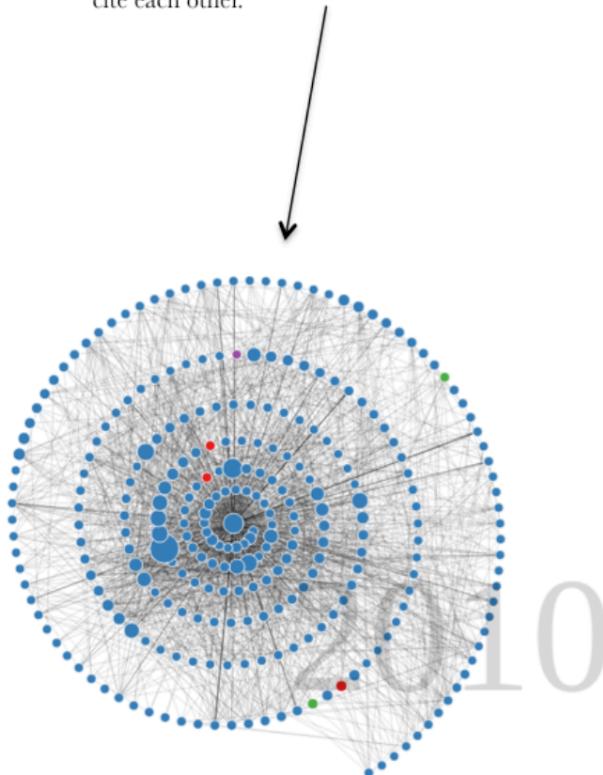
Comparing Authors

- Papers in category "Medicine" (domain 6)
- Papers in category "Biology" (domain 4)
- Papers in category "Chemistry" (domain 5)
- Papers in category "Engineering" (domain 8)
- Papers in category "Material Science" (domain 12)
- Papers in category "Physics" (domain 19)
- Papers in category "Agriculture Science" (domain 16)
- Papers in category "Social Science" (domain 22)



A more sparse network indicates fewer citations between papers shown in the network. This could be a result of the central scholar having impact across a wider set of academic communities.

A denser network means that the papers that cite the central author also tend to cite each other.



Visualizing Scholarly Influence Over Time

Influence of Pew Scholars

Mark W. Grinstaff

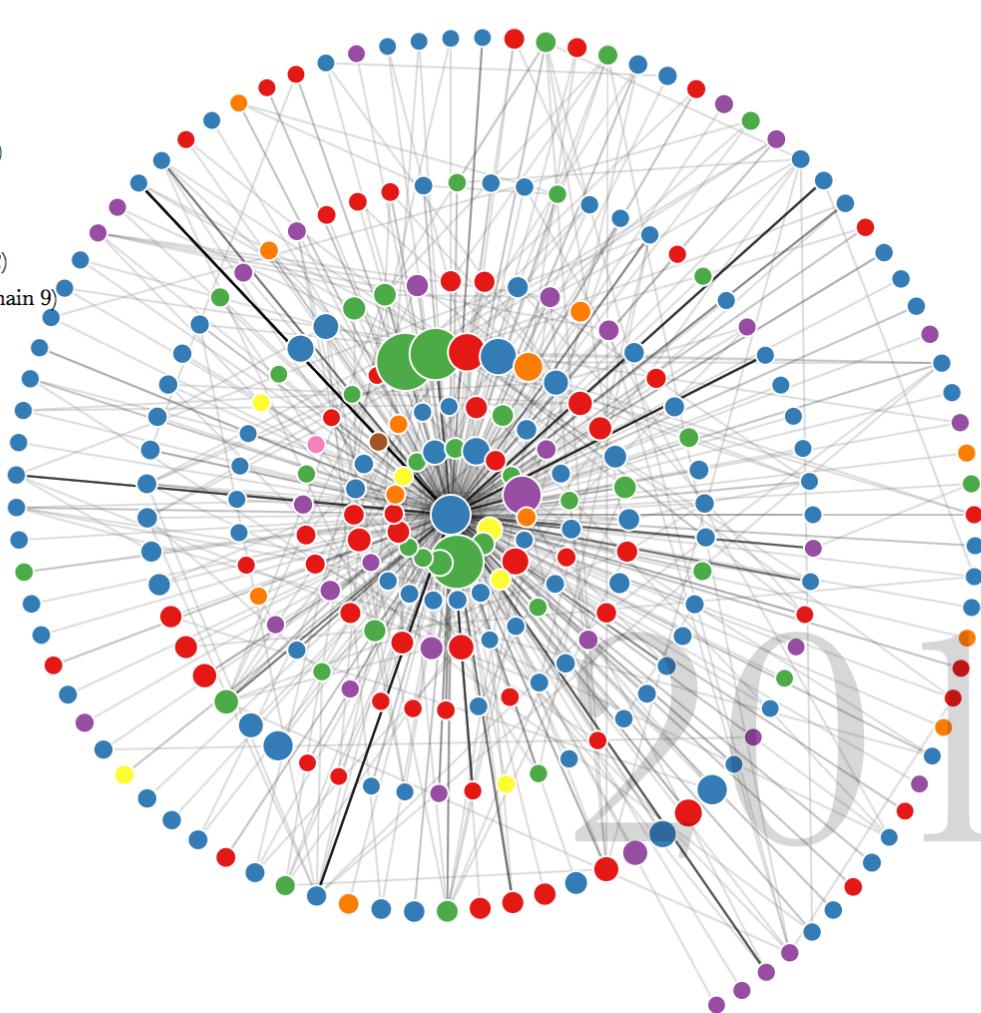
[Learn More](#)

- Papers in category "Chemistry" (domain 5)
- Papers in category "Medicine" (domain 6)
- Papers in category "Biology" (domain 4)
- Papers in category "Material Science" (domain 12)
- Papers in category "Engineering" (domain 8)
- Papers in category "Physics" (domain 19)
- Papers in category "Computer Science" (domain 2)
- Papers in category "Environmental Sciences" (domain 9)

Mark W.
Grinstaff



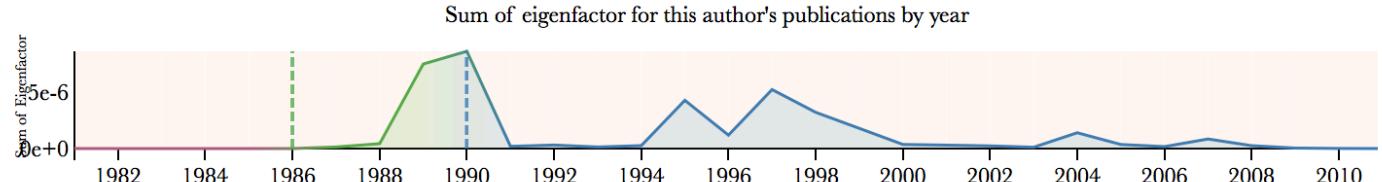
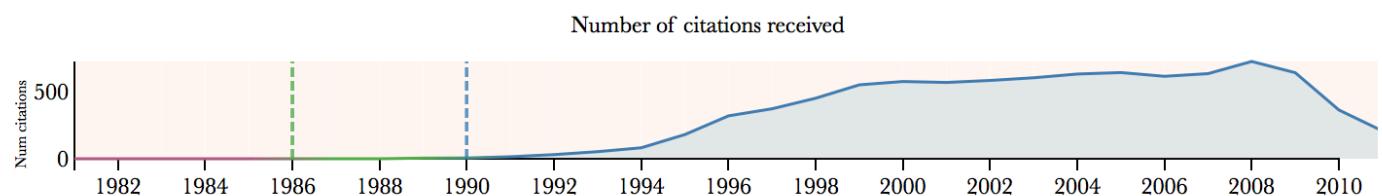
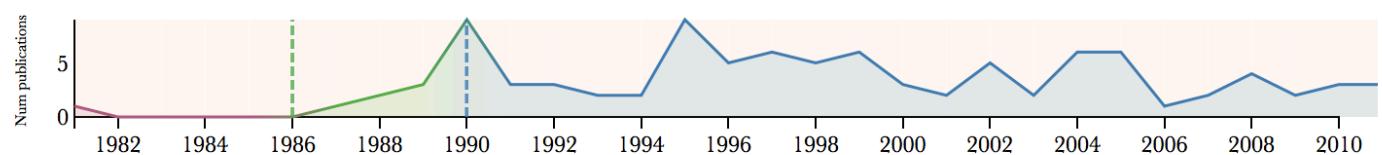
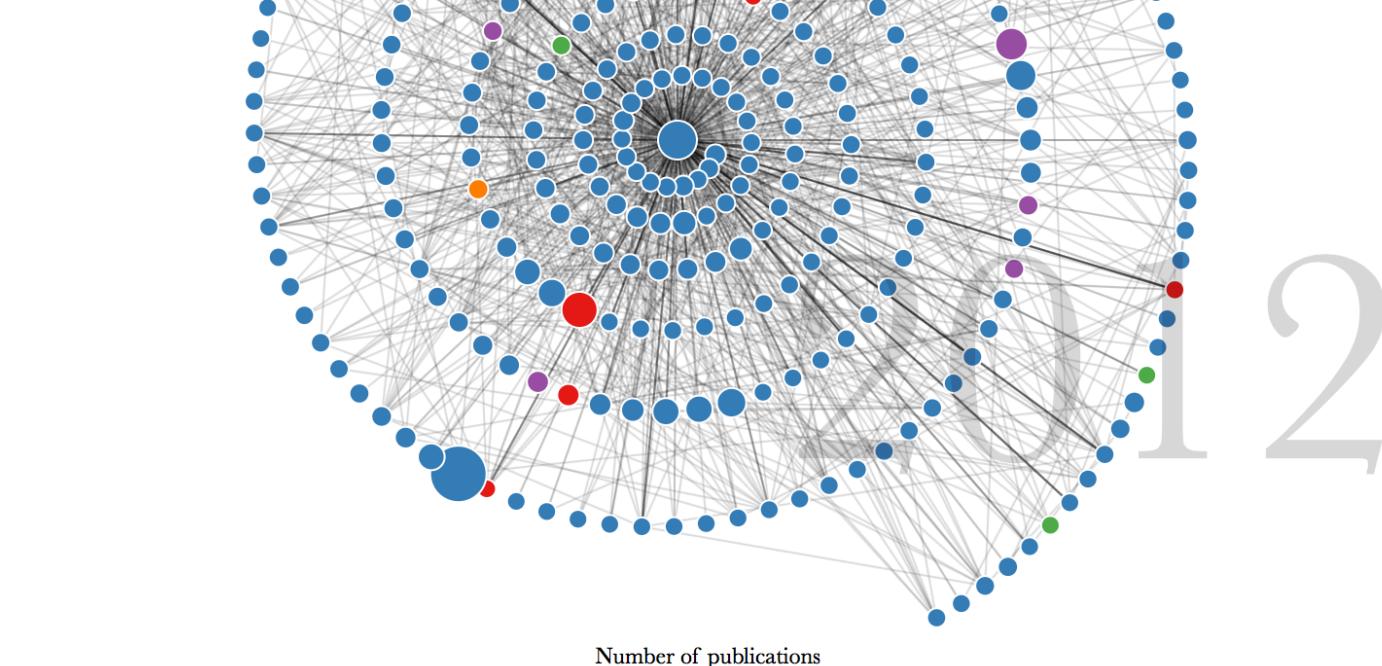
Pew Scholar
1999



Philip A.
Hieter



Pew Scholar
1986



Academic Migration

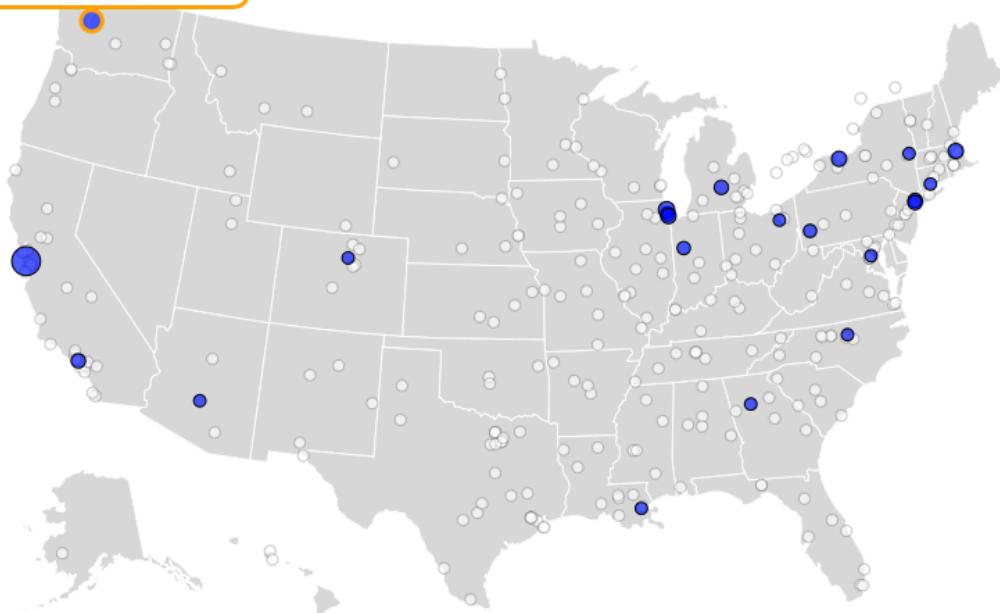
U.S. ACADEMIC MIGRATION MAP

Where do people who pursue academic careers in the U.S. go on to land faculty positions after earning their advanced degrees?

Where do faculty come from? Click on a school to explore.

[ABOUT THIS PROJECT](#)

University of Washington



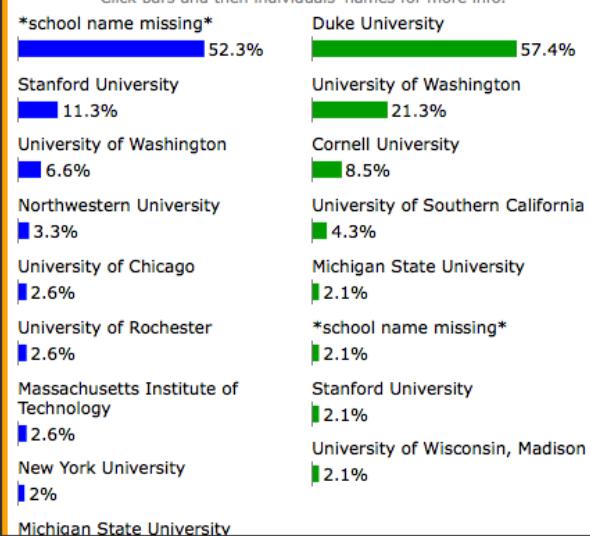
[Clear selection](#)

[Resize & re-center map](#)

University of Washington

Current sample: 198 individuals (151 faculty and 47 grads)

Faculty have come from
these schools:



Click to view grad
destinations on map

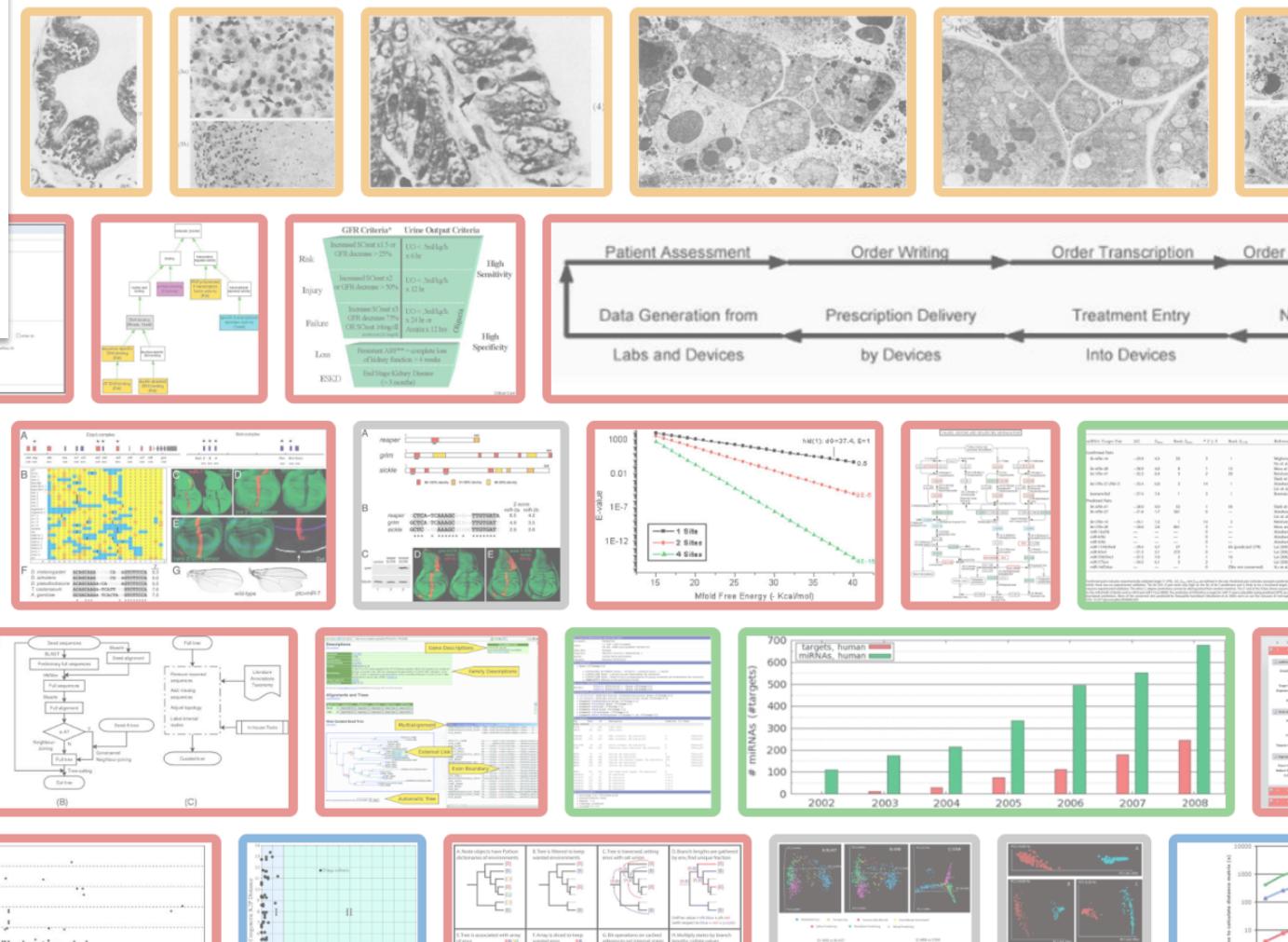


Poshen Lee

Impact

Keywords or Cluster, Result Ordered by Impact

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?

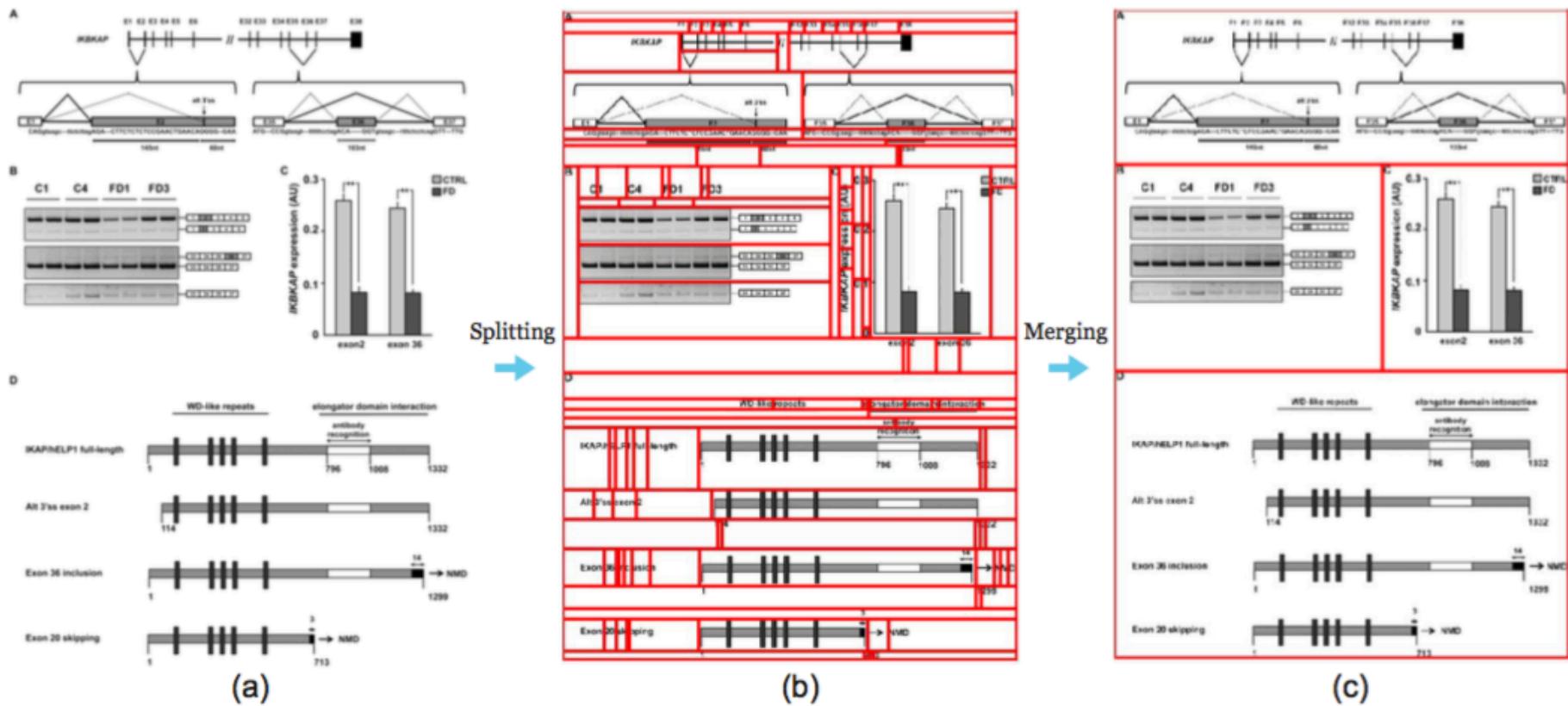
*How can we better utilize visual information
in the search and navigation process?*



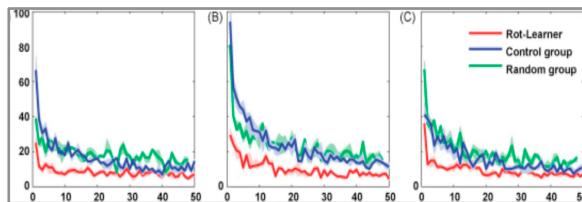
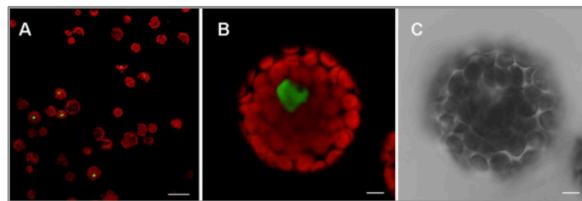
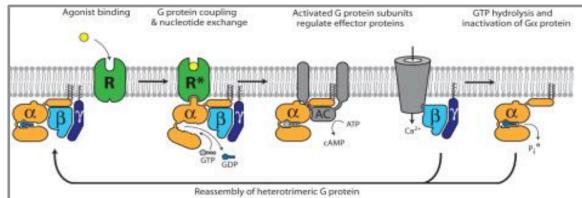
650,000 papers

5 million images

Composite Figure Dismantling



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



	PW reading	PW reading	PW reading	W reading	W reading
	W RT	PW RT	CTL	W RT	PW RT
MOG → LOT	0.28	0.18	0.58	-0.70	-0.50
MOG → LP	-0.22	-0.52	-0.04	0.27	-0.03
LOT → LP	0	0.10	0.24	-0.56	-0.60
LOT → IFG	0.38	0.17	0.40	0.43	0.13
LP → IFG	0.26	0.05	0.31	0.03	-0.03

Equations (394)

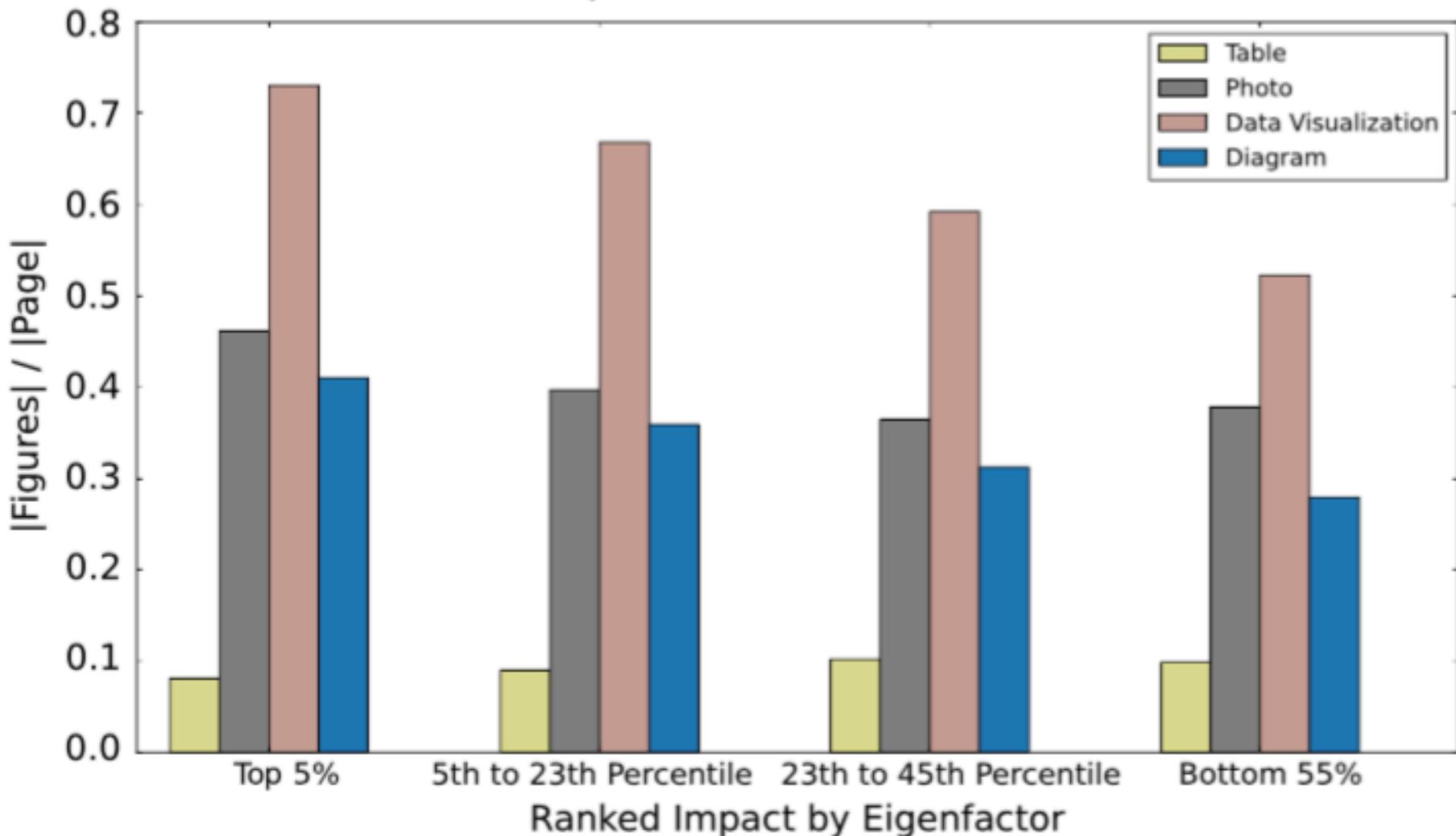
Schematics (769)

Photos (782)

Plots (890)

Tables (436)

Impact versus Figure Density



viziometrics.org

Viziometrics About Search Crowdsourcing

Random Keywords or Cluster, Result Ordered by Random Search

Composite Equation Diagram Photo Table Visualization

A project of the eScience Institute at the University of Washington

(b)

Viziometrics About Search Crowdsourcing

Random Keywords or Cluster, Result Ordered by Random Search

Composite Equation Diagram Photo Table Visualization

On the Number of New World Founders: A Population Genetic Portrait of the Peopling of the Americas
Clark Andy G. and Hey Jody

Malaria in American troops in the South and Southwest Pacific in World War II.

A project of the eScience Institute at the University of Washington

Vascular ossification – calcification in metabolic syndrome, type 2 diabetes mellitus, chronic kidney disease, and calciphylaxis – calcific uremic arteriolopathy: the emerging role of sodium thiosulfate
Khanna Rameesh, Sowers James R, Koltz Leea, Tyagi Suresh C and Hayden Melvin R
Cardiovascular Diabetology 2009

Abstract

Background/Background Vascular calcification is associated with metabolic syndrome, diabetes, hypertension, atherosclerosis, chronic kidney disease, and end stage renal disease. Each of the above contributes to an accelerated and premature demise primarily due to cardiovascular disease. The above conditions are associated with multiple metabolic toxicities resulting in an increase in reactive oxygen species to the arterial vessel wall, which results in a response to injury/wound healing (remodeling). The endothelium seems to be at the very center of these disease processes, acting as the first line of defense against these multiple metabolic toxicities and the first to encounter their damaging effects to the arterial vessel wall. Vascular calcification is associated with metabolic syndrome, diabetes, hypertension, atherosclerosis, chronic kidney disease, and end stage renal disease. Each of the above contributes to an accelerated and premature demise primarily due to cardiovascular disease. The above conditions are associated with multiple metabolic toxicities resulting in an increase in reactive oxygen species to the arterial vessel wall, which results in a response to injury/wound healing (remodeling). The endothelium seems to be at the very center of these disease processes, acting as the first line of defense against these multiple metabolic toxicities and the first to encounter their damaging effects to the arterial vessel wall. Results/Hypothesis The pathobiological mechanisms of vascular calcification are presented in order to provide the clinician – researcher a database of knowledge to assist in the clinical management of these high-risk patients and examine newer therapies. Calciphylaxis is associated with medial arterial vascular calcification and results in ischemic subcutaneous necrosis with vulnerable skin ulcerations and high mortality. Recently, this clinical syndrome once thought to be rare is presenting with increasing frequency. Consequently, newer therapeutic modalities need to be explored. Intravenous sodium thiosulfate is currently used as an antidote for the treatment of cyanide poisoning and prevention of severe rhabdomyolysis. It is also used as a food and medicinal preservative and frequently used as an antifungal medication. Conclusion/Conclusion A discussion of sodium thiosulfate's dual role as a potent antioxidant and chelator of calcium is presented in order to better understand its role as an emerging novel therapy for the clinical syndrome of calciphylaxis and its complications.

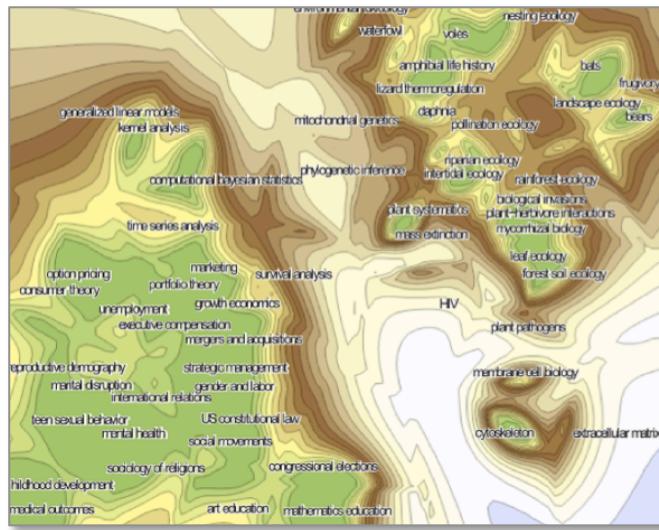
[Hide Abstract] [View Paper] [View Cluster]

Diagram Validate
Help us improve the accuracy

The central role of the endothelium in VOD and atherosclerosis. This image portrays the endothelium as the first line of defense against multiple injurious stimuli. Most of the injurious stimuli are represented by the A-FLIGHT-U toxicities found in table 3. When discussing the role of VOD and how it ties into atherosclerosis and the accelerated ASO associated with metS, predabetes, and overt T2DM it is important to include the various interactions of A-FLIGHT-U toxicities with the associated ROS and the sensitizers for calcium deposition and ossification (Ca++, Pi, FTH) within the AW (both the media and intima) in table 3. The role of ROS, inflammation monocyte-macrophage foam cells undergoing apoptosis necrosis with creation of a nodule and the stimuli for the VSMC and pericyte following the neovascularization (Vv) of the media and intima to result in ossoid formation and later the mineralization within these atherosclerotic plaques. This image also portrays the possible role for nondiabetes in addition to C. pneumoniae as well as other bacteria and viral infections as possible contributing factors for the calcification process.

[Hide Other Figures] [Show Related Figures]

Science of Science

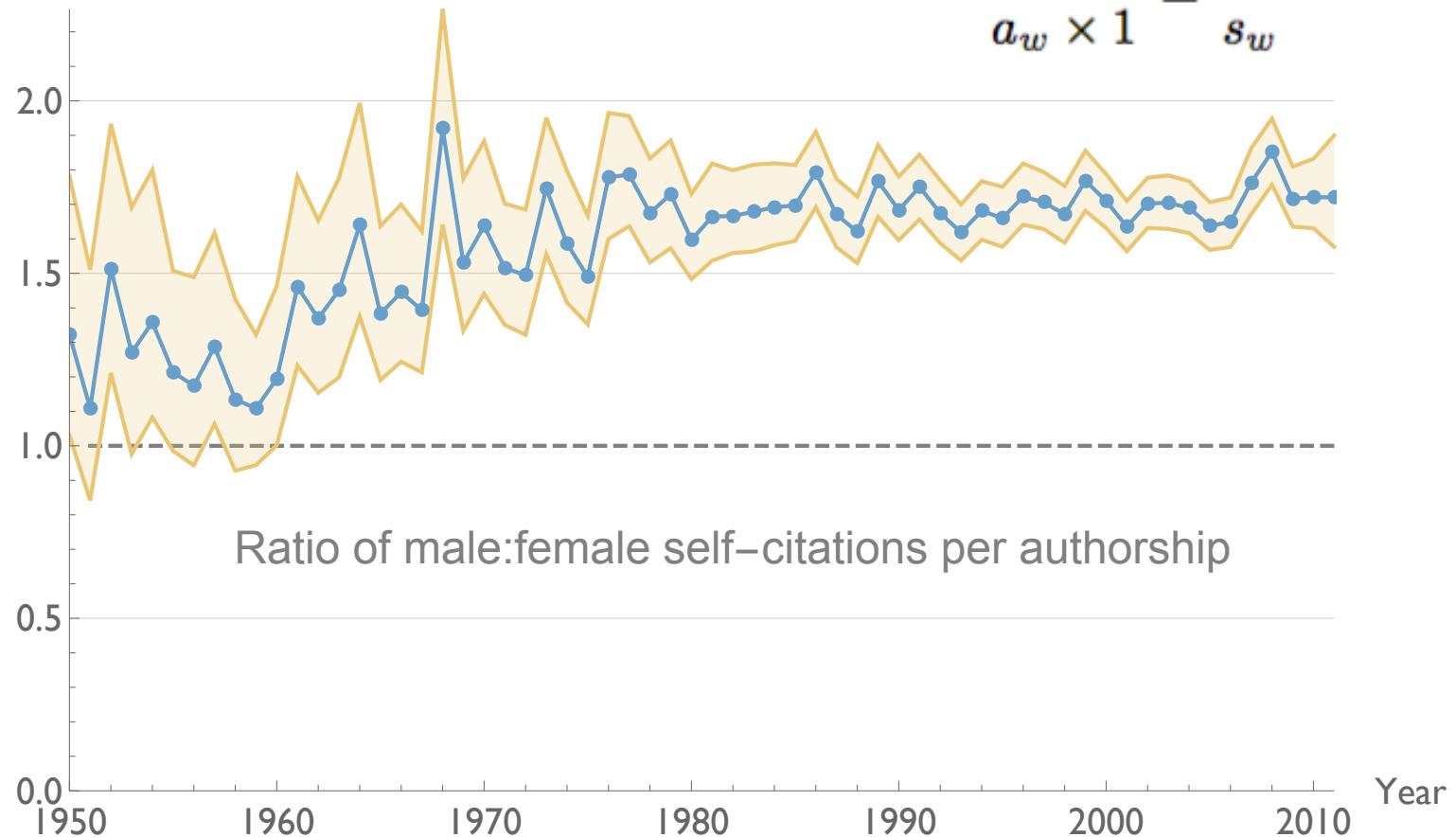


Incentives, reproducibility, biases, funding models, publishing economics...

Self-citation over time

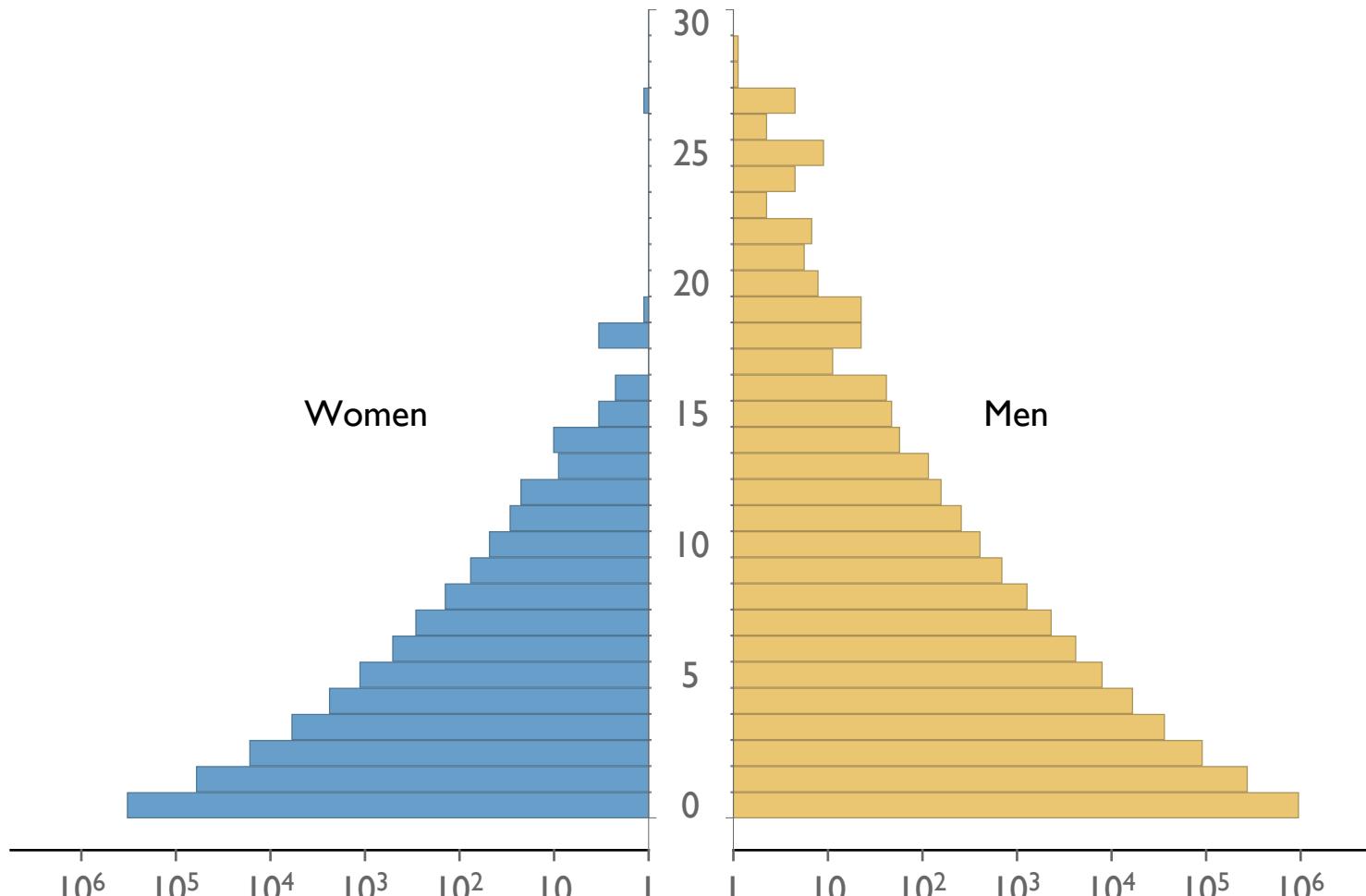
Ratio M:F

$$\frac{a_m \times k}{a_w \times 1} = \frac{s_m}{s_w}$$

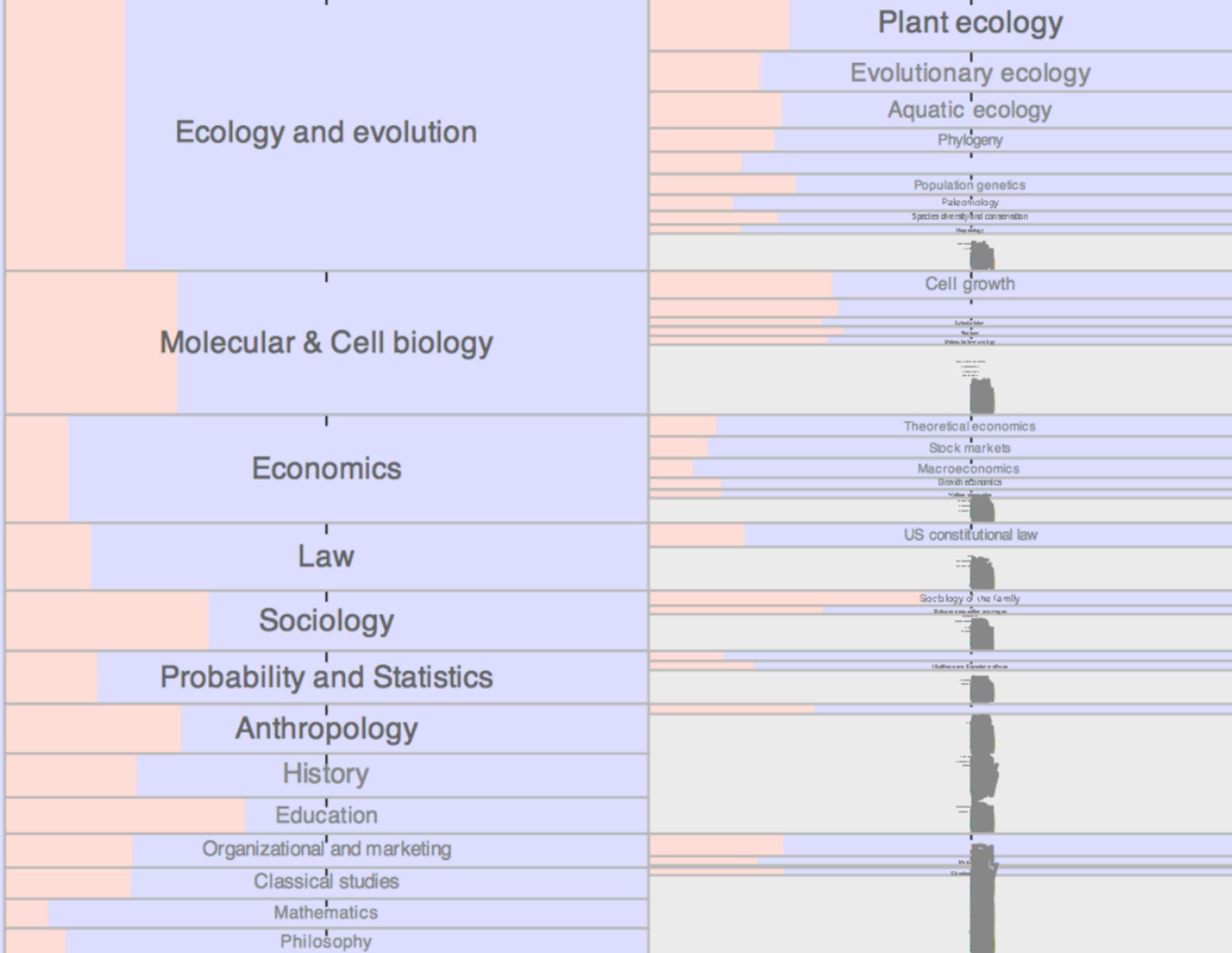


Excessive Self Citation

Number of authorships with n self-citations



King et al (2016) in prep.



Information Visualization

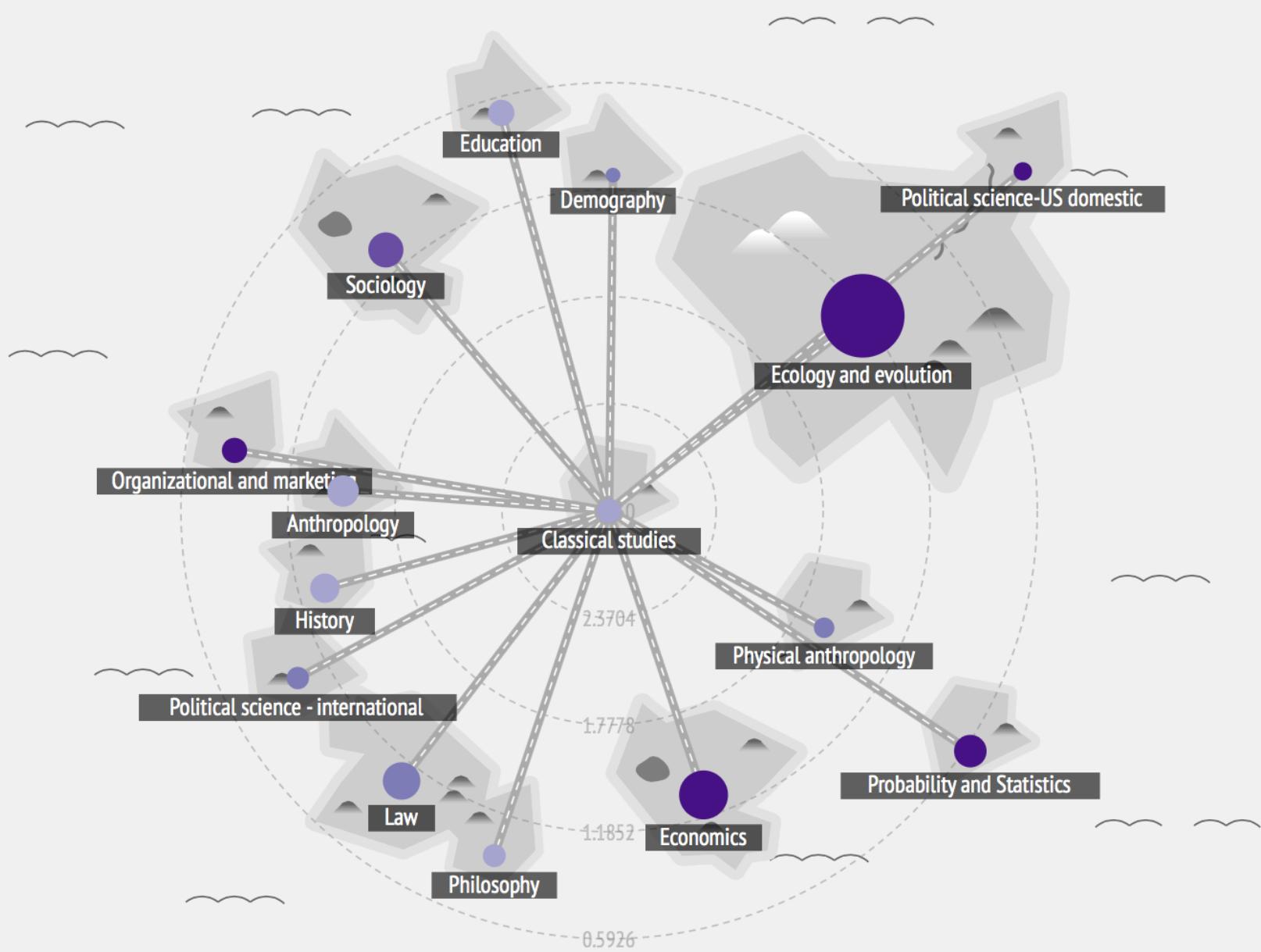


Cecilia Aragon Michael Brooks
UW, HCDE UW, HCDE

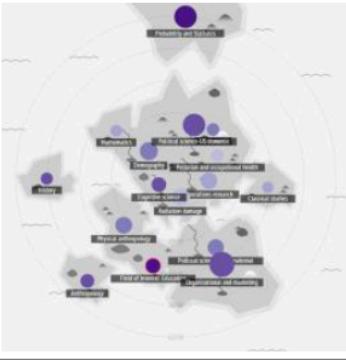
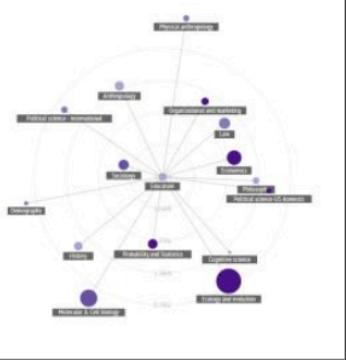
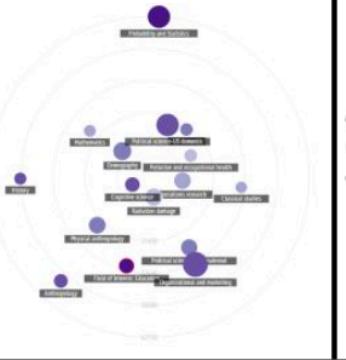
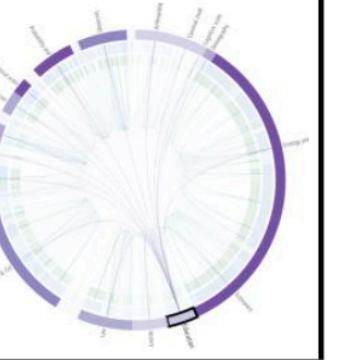


How do we augment human
memory when zooming in/out of
hierarchical trees?

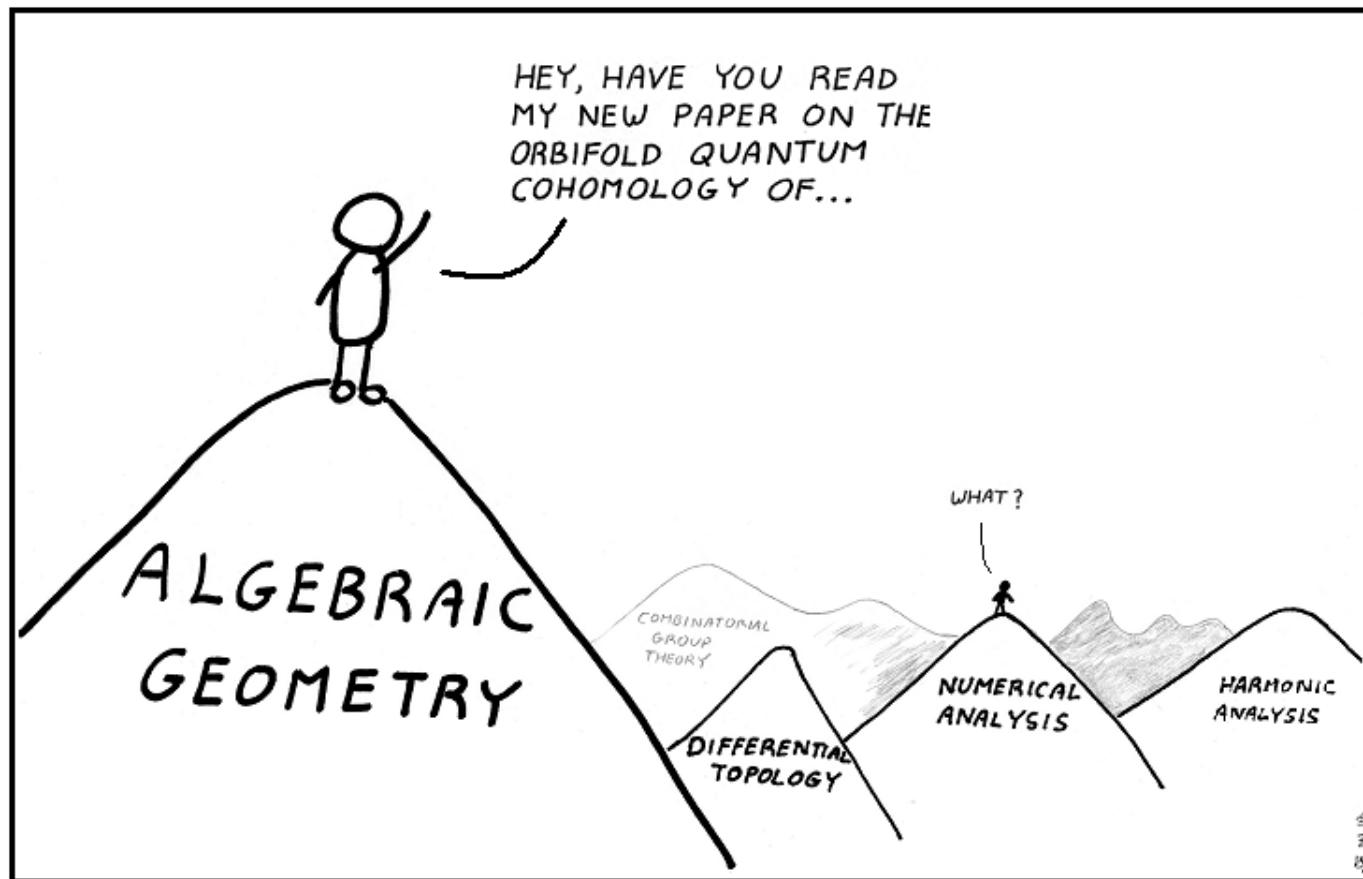
Brooks, M et al. (2013) *INTERACT*



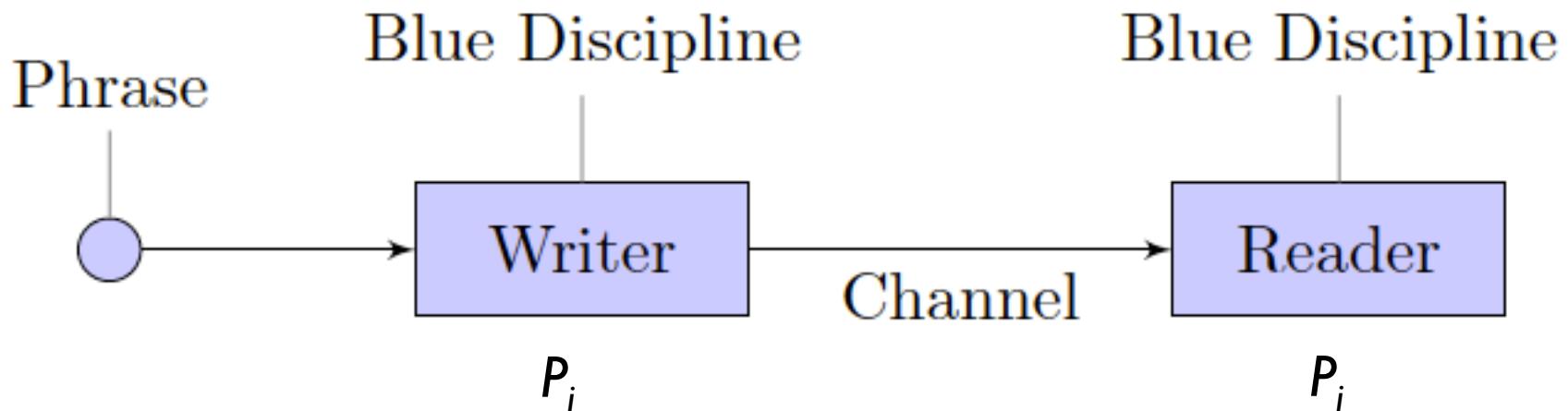
Navigating Hierarchical Knowledge Networks

1. Congruent Landscape	2. Incongruent Landscape	3. Congruent Abstract	4. Incongruent Abstract	5. Designer Baseline
 <p>Landscape visualization with data properties mapped to visual elements according to applicable image schemata</p>	 <p>Landscape visualization with data properties mapped to visual elements deliberately breaking with image schemata</p>	 <p>Identical to the Congruent Landscape tool but with all realistic details and overt “landscape” visuals removed</p>	 <p>Identical to the Incongruent Landscape tool but with all realistic details and overt “landscape” visuals removed</p>	 <p>Visualization designed by a hypothesis-blind designer attempting to make an effective visualization but without special emphasis on metaphor</p>

The jargon barriers of science



The Landscape of Modern Mathematics



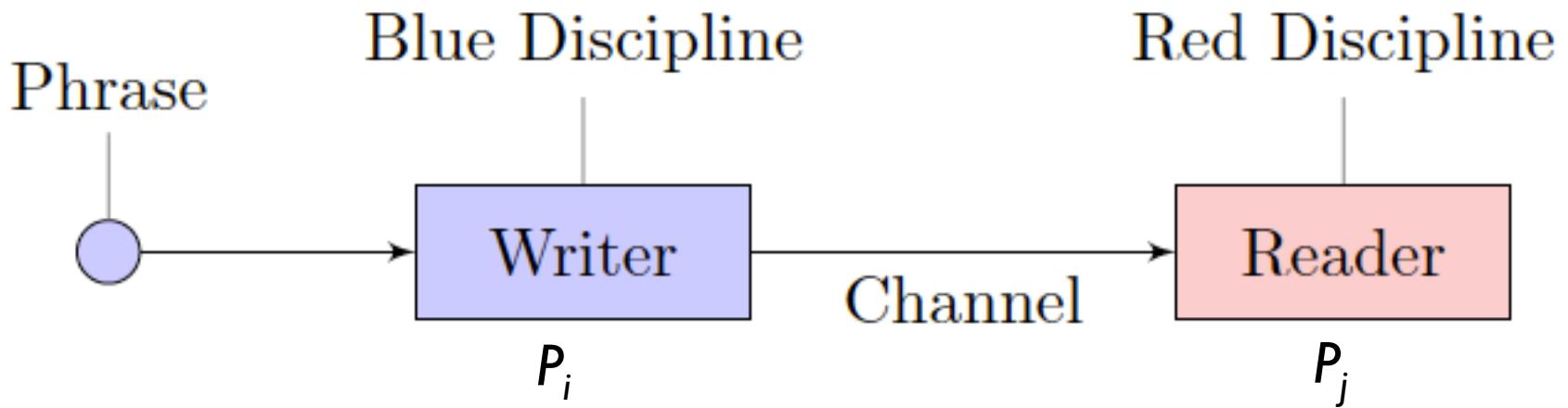
$X \sim$ space of all phrases

$P_i \sim$ probability distribution over x_i with values $x \in X$

- writer chooses phrases with probability $p_i(x)$
- optimal codeword has length $-\log_2 p_i(x)$

expected message length $H(X_i) = - \sum_{x \in \mathcal{X}} p_i(x) \log_2 p_i(x)$

assumption: language of each scientific field is *optimized* based on frequency of phrases



cross entropy
 ↓
 expected message length: $Q(p_i||p_j) = - \sum_{x \in \mathcal{X}} p_i(x) \log_2 p_j(x)$

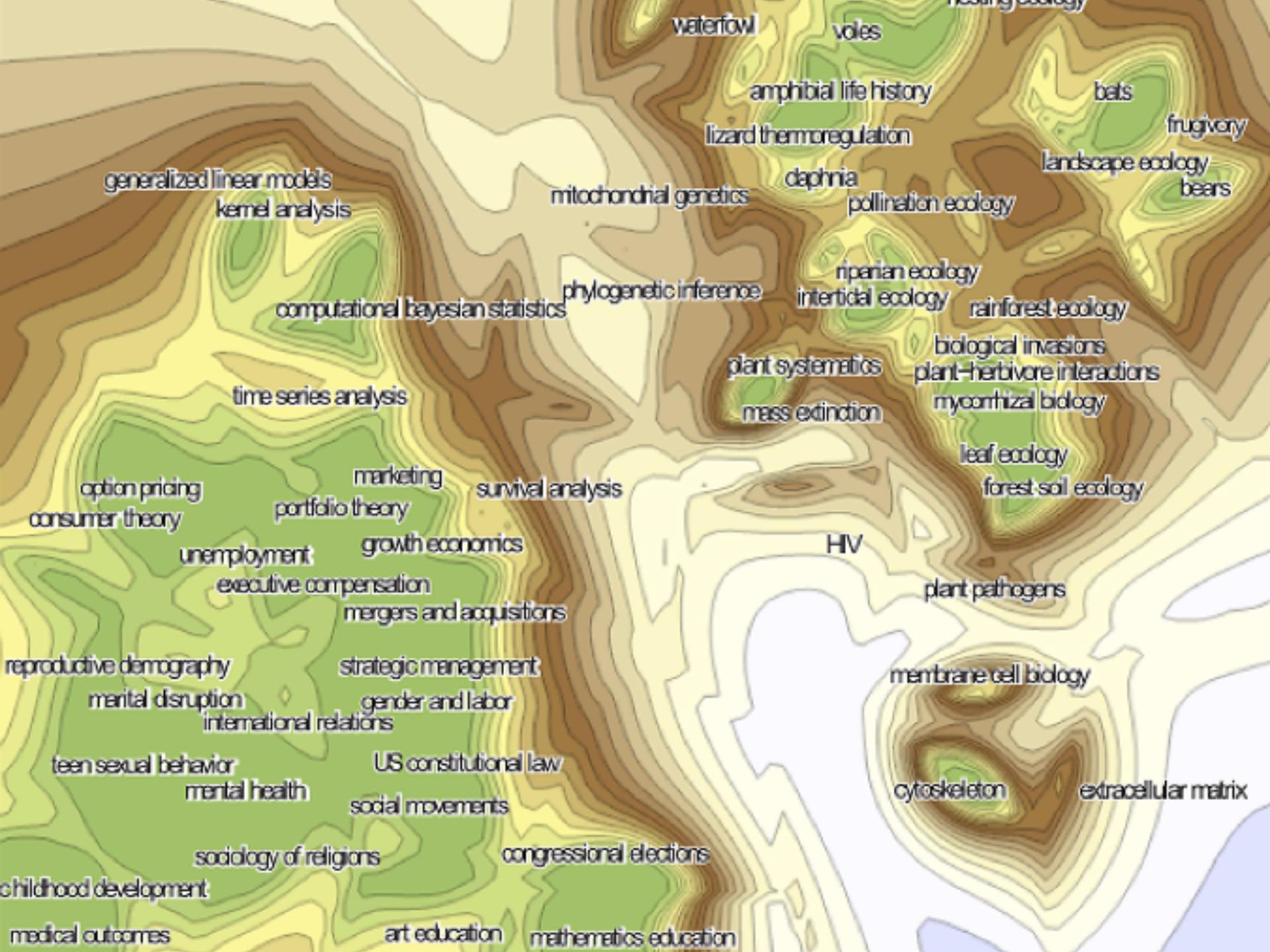
efficiency of communication

$$\downarrow$$

$$E_{ij} = \frac{H(X_i)}{Q(p_i||p_j)} = \frac{-\sum_{x \in \mathcal{X}} p_i(x) \log_2 p_i(x)}{-\sum_{x \in \mathcal{X}} p_i(x) \log_2 p_j(x)}$$

$$C_{ij} = 1 - E_{ij}$$

↑
 cultural hole



Challenges: Labeling

$\text{label}(\mathcal{C}_i) = \mathcal{N}_j$ where $j = \underset{j}{\text{argmax}} \hat{I}(f(\varphi, i), g(\varphi, j))$

$$\hat{I}(f_i(\varphi), g_j(\varphi)) = \frac{I(f_i(\varphi), g_j(\varphi))}{H(f_i(\varphi))}$$

Normalized Mutual Information

Semantic Scholar

- Babel (data and recommendation)
- InfoMap (clustering citation networks)
- Scholar profiles (share data)
- Auto-detecting neuroscience
- Viziometrics (labeling tools)
- Visual interface design
- Auto-Labeling Fields
- Data cleaning (author disambiguation)

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