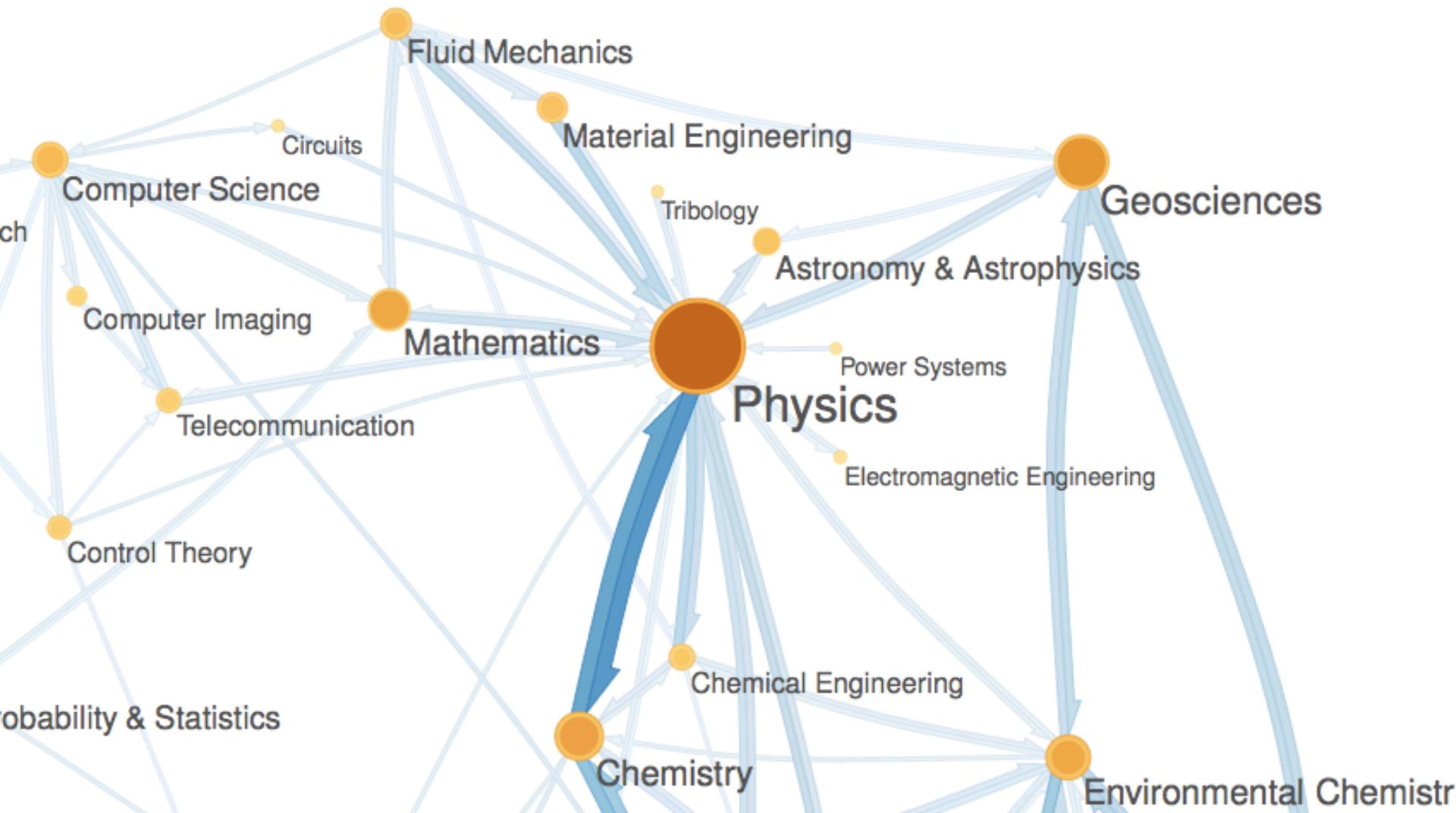
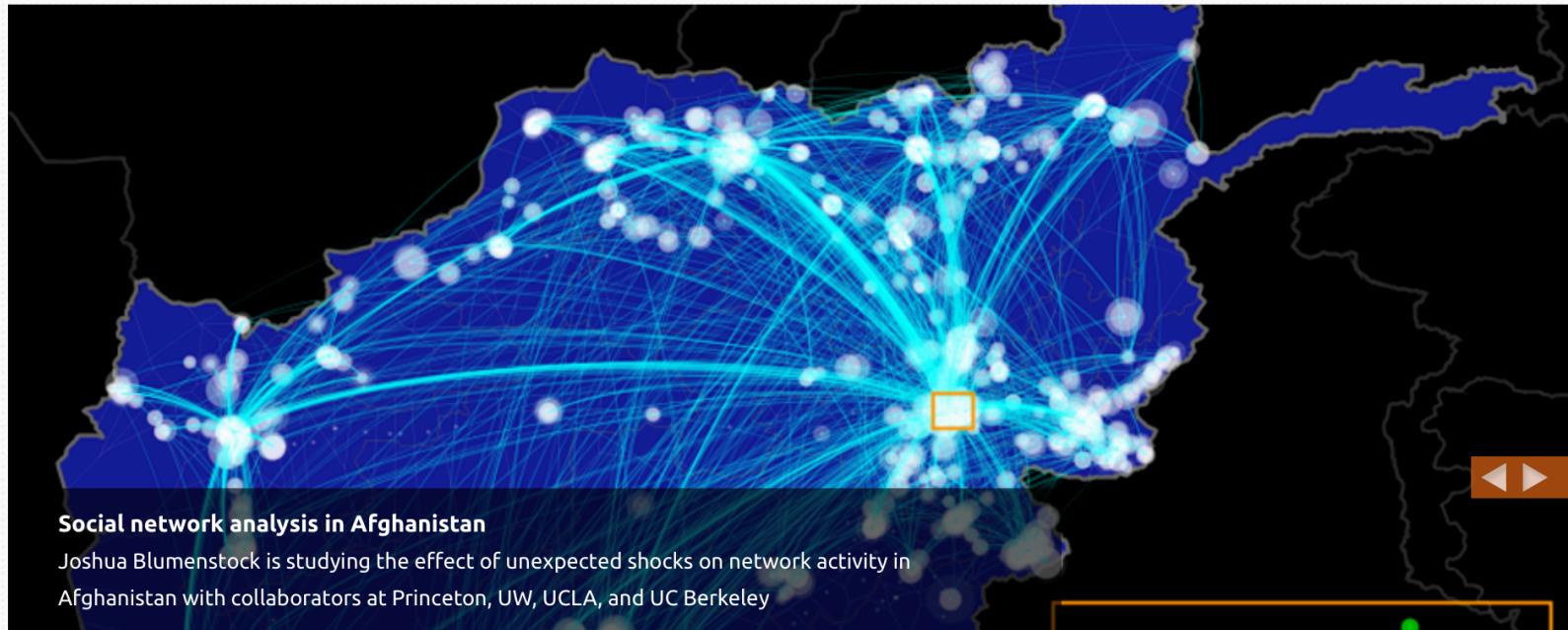


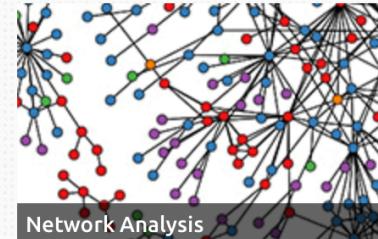
Mapping the emergence of scientific disciplines

Jevin West, Information School, University of Washington





Research Focus Areas



News and Updates

28

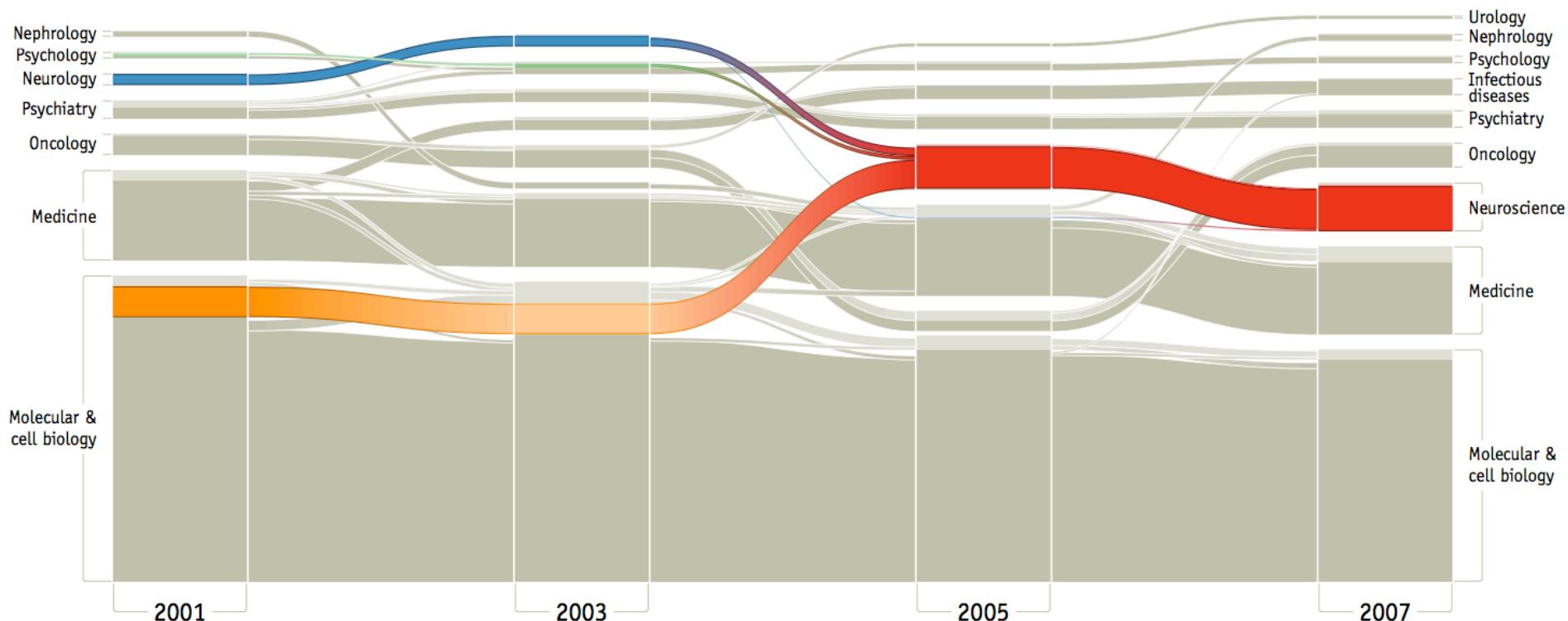
Blumenstock at Population Association of America

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

How do we *map* the origins of scientific disciplines?

The Emergence of Neuroscience





The Scholarly Graph



PatentVector™



WIKIPEDIA
The Free Encyclopedia



PNAS





The Scholarly Graph



Tens of millions articles, patents, books



Billions of citation links

PatentVector™

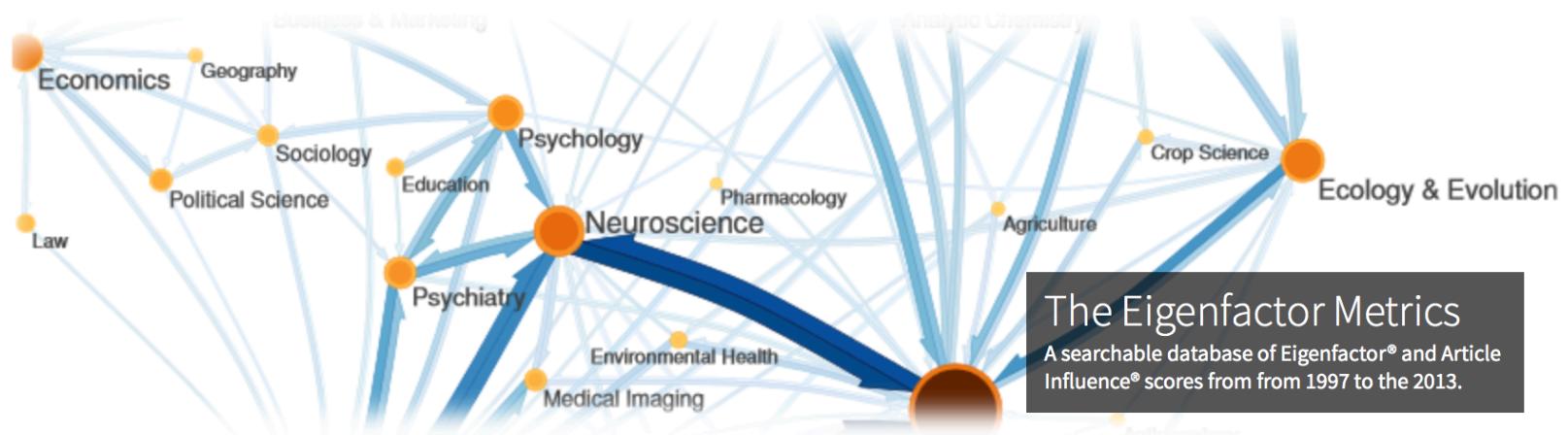


Years: 1600 - 2016

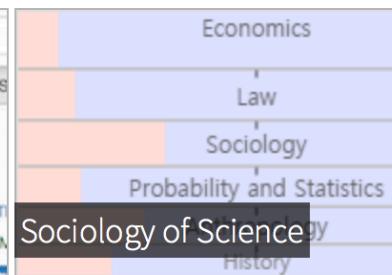
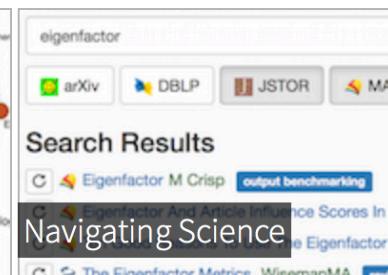
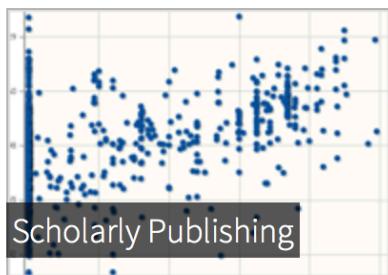


Science of Science

- What ideas, papers and scientists seeded new fields?
- Can we automate hypothesis generation? If so, how would this change science?
- Can we teach the computer to read the literature and understand figures?
- How can we improve scientific ‘navigation’ in the face of information overload?
- Can we predict new discoveries before they happen?
- How can we improve scientific institutions, reward mechanisms and funding processes?
- How can we facilitate interdisciplinary research and foster innovation rather than just piles of papers?



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

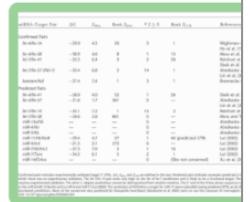
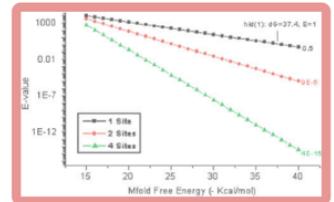
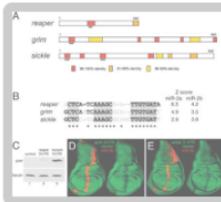
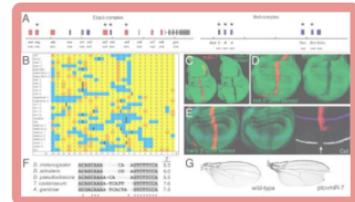
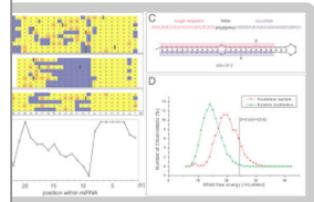
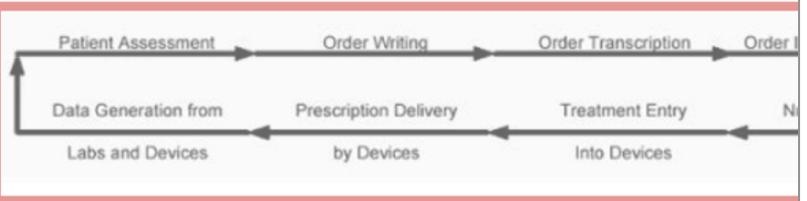
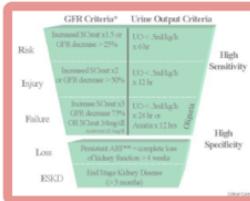
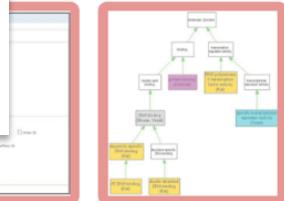
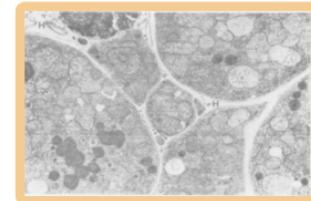
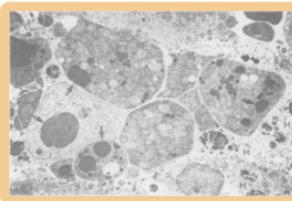
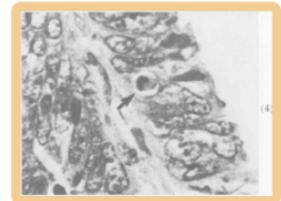
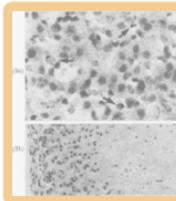
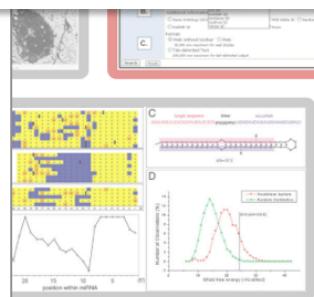
Next »

Papers related to

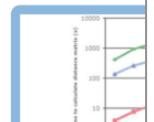
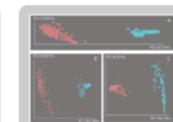
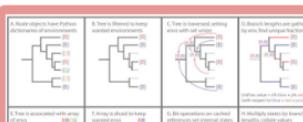
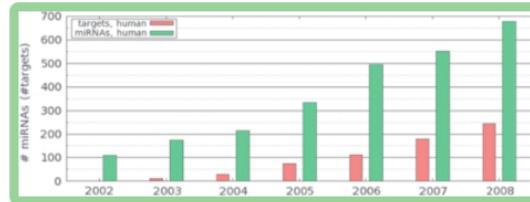
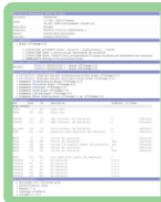
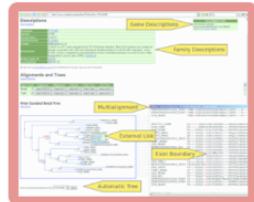
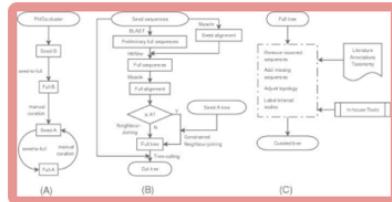
- Statistical Methods For Analyzing Speedup Learning Experiments O Etzioni satisfaction programs
- Automatically Configuring Constraint Satisfaction Programs: A Case Study S Minton 1995 satisfaction programs
- Abstraction Via Approximate Symmetry T Ellman 1992 satisfaction programs
- Integrating Heuristics 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

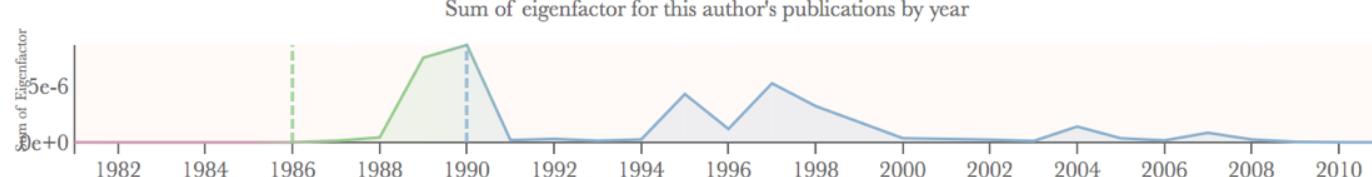
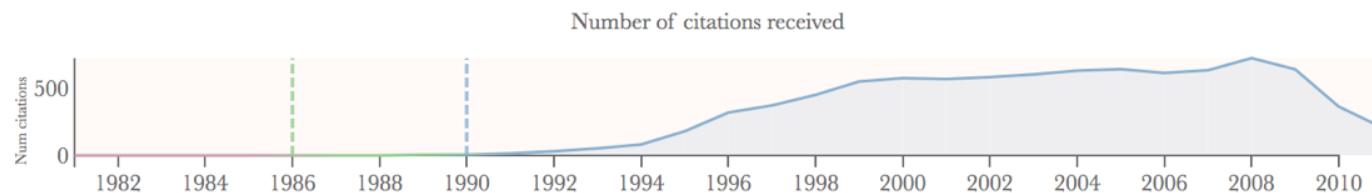
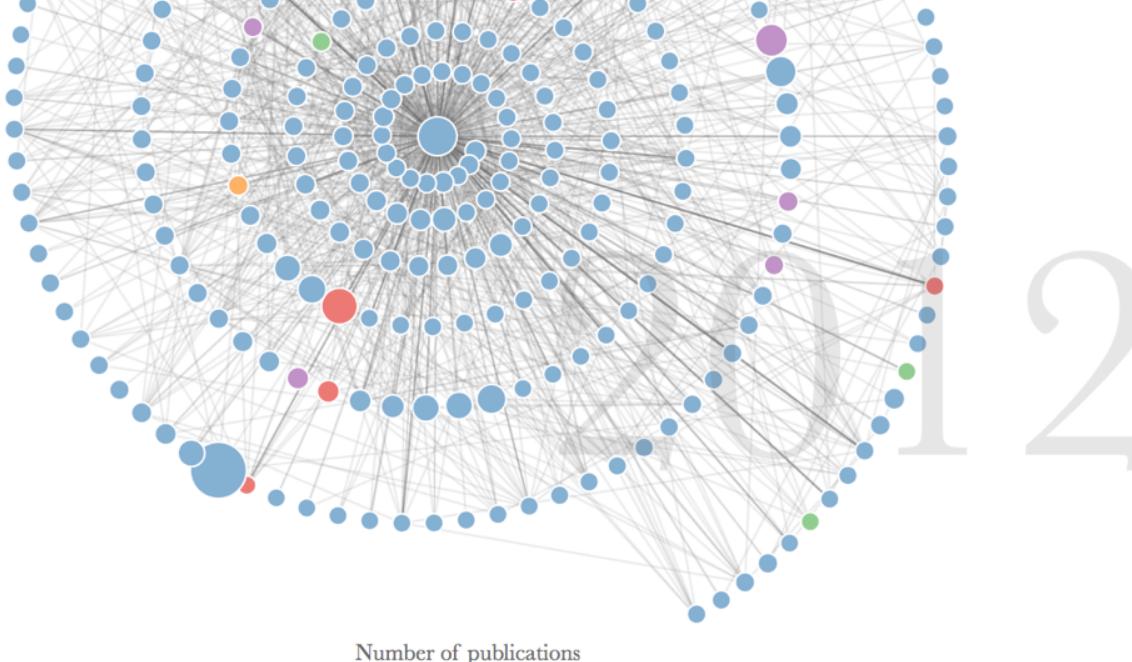


Enzyme	# Z >= 3	Z_{0.05}
(EC)2.4.1.41	2	7.31
(EC)3.1	2	6.53
(EC)2.6.1.42	1	5.75
(EC)1.1.1.35	1	4.80
(EC)2.4.4	1	4.51
(EC)2.8.3.5	1	4.23
(EC)1.1.1.35	1	3.9
(EC)1.1.1.31	1	3.81
(EC)1.2.1.27	1	3.72
		3.64
		3.61

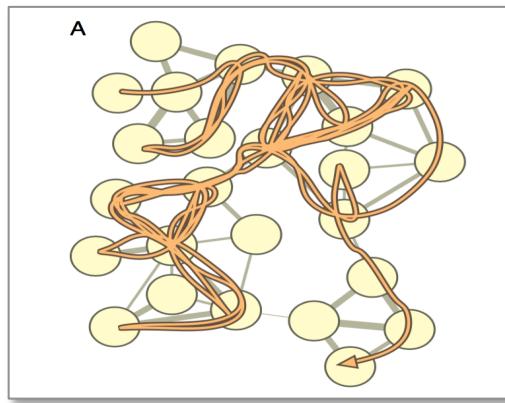




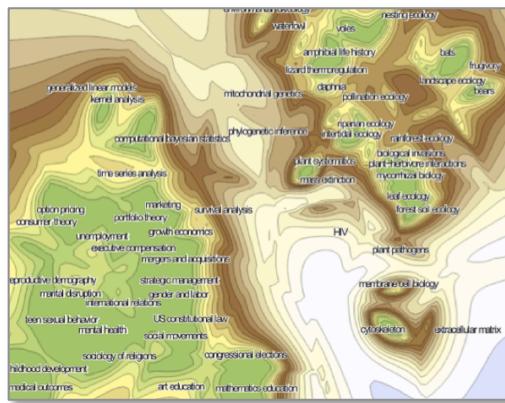
Jason Portenoy



scholar.eigenfactor.org

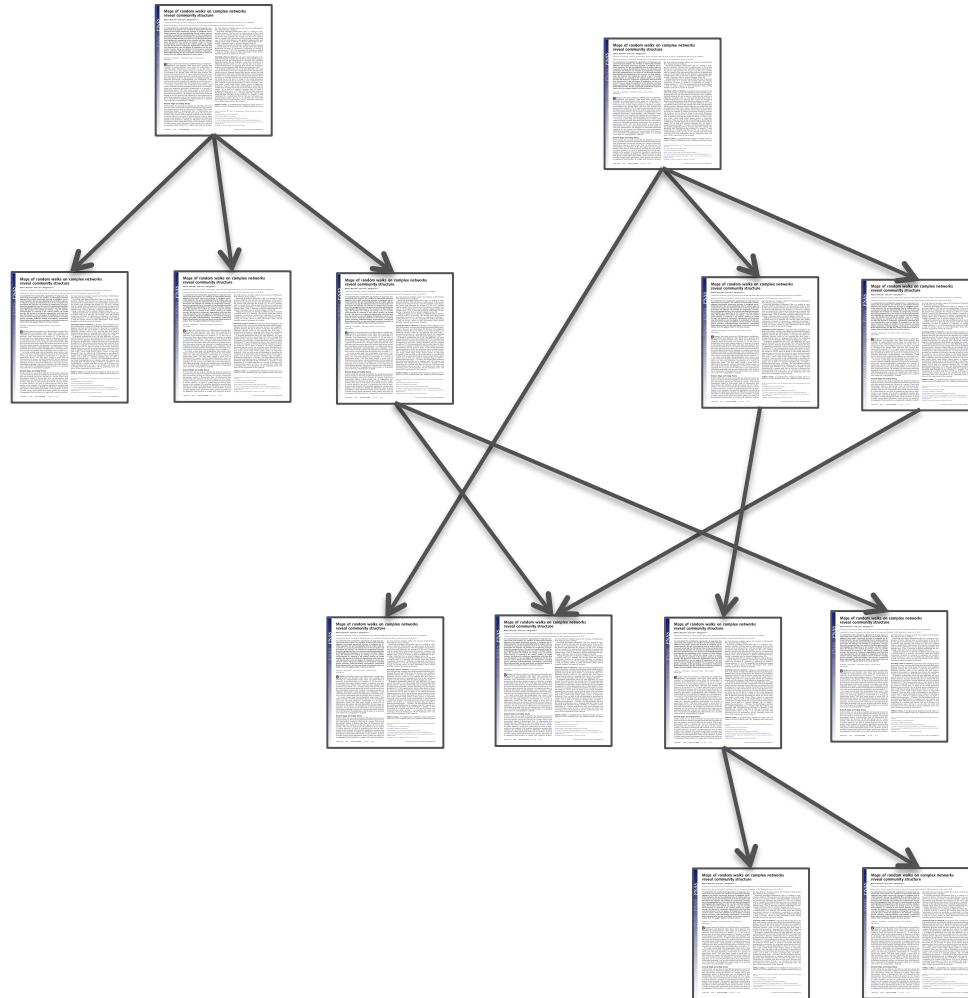


Science of Mapping

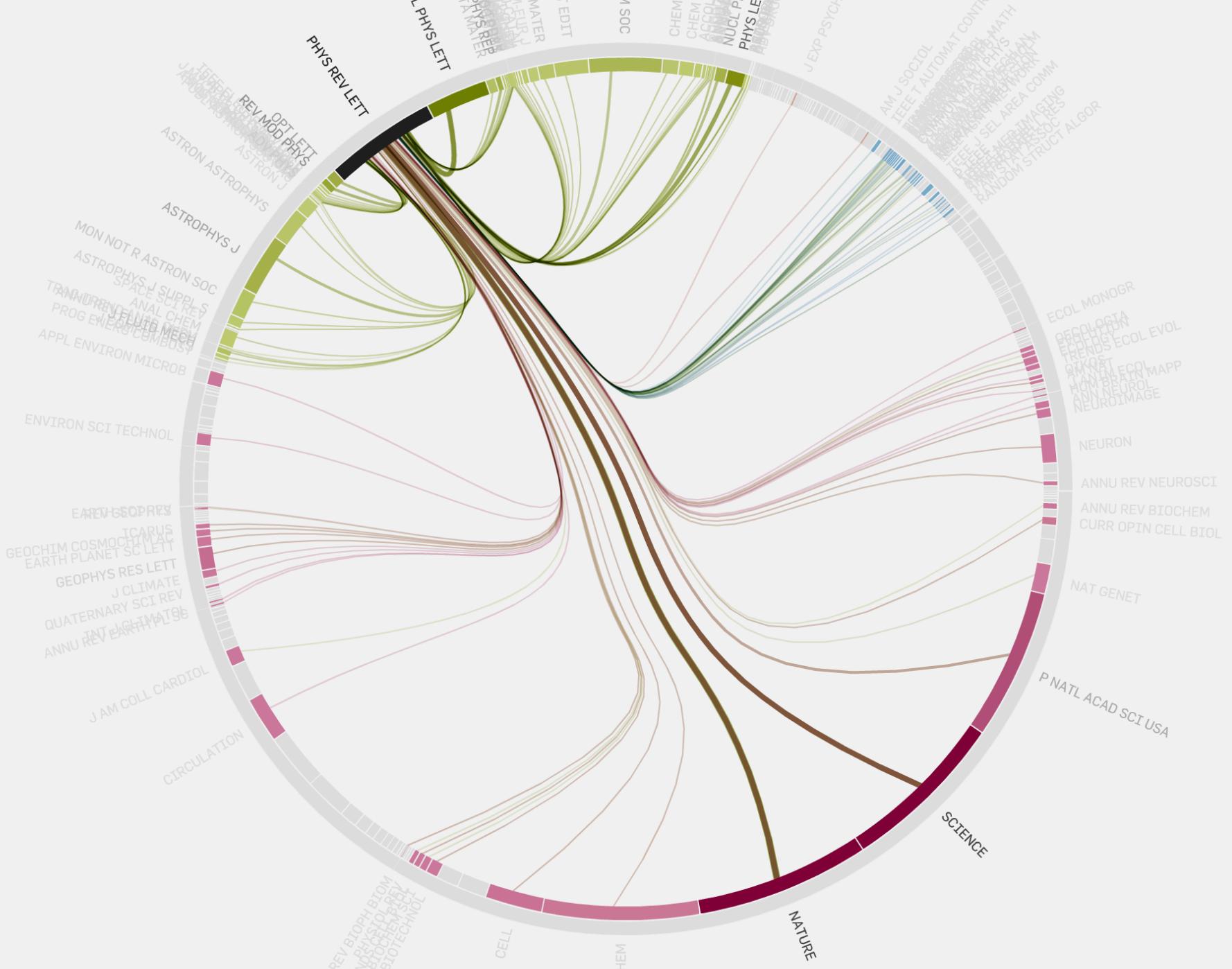


Mapping of Science

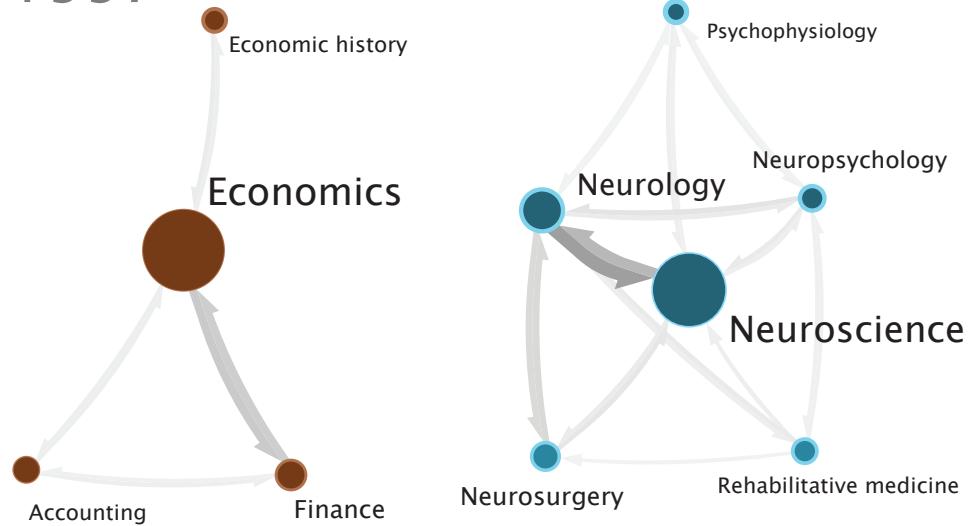
Citations form a vast network



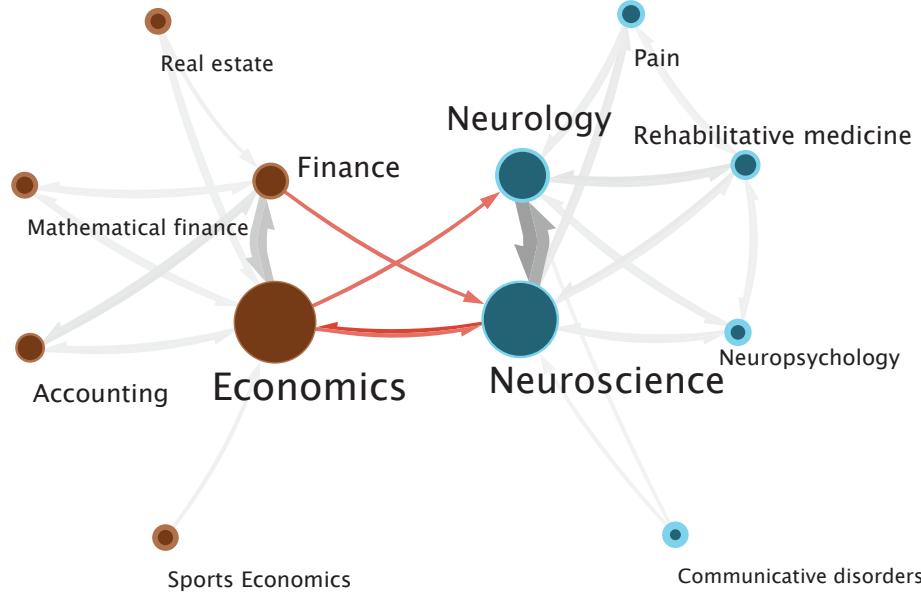
de Solla Price, Science (1965)



1997



2010



NEUROECONOMICS

Recent common ancestor

Of the 14 million pairs of papers, 7.1M of them have a common ancestor—a hit rate of about 50%. Among these 7.1M LCAs, here are the ten most frequent papers (and their frequencies):

1. (47,129) *Some Methods for Strengthening the Common χ^2 Tests* (Cochran, 1954)
2. (35,585) *The Evolution of Reciprocal Altruism* (Trivers, 1971)
3. (34,195) *On the Mathematical Foundations of Theoretical Statistics* (Fisher, 1922)
4. (34,093) *The Tragedy of the Commons* (Hardin, 1968)
5. (32,067) *Some Difficulties of the Determination Problem* (Harrison, 1933)
6. (29,458) *Diverse Doctrines of Evolution, Their Relation to the Practice of Science and of Life* (Jennings, 1927)
7. (28,149) *An Analysis of Transformations* (Box, 1964)
8. (26,000) *Fitting the Negative Binomial Distribution to Biological Data* (Bliss, 1953)
9. (25,410) *A Method for Cluster Analysis* (Edwards, 1965)
10. (24,611) *A Theory of the Allocation of Time* (Becker, 1965)

Visualizing Scholarly Influence Over Time

Influence of Pew Scholars

Roberta A. Gottlieb

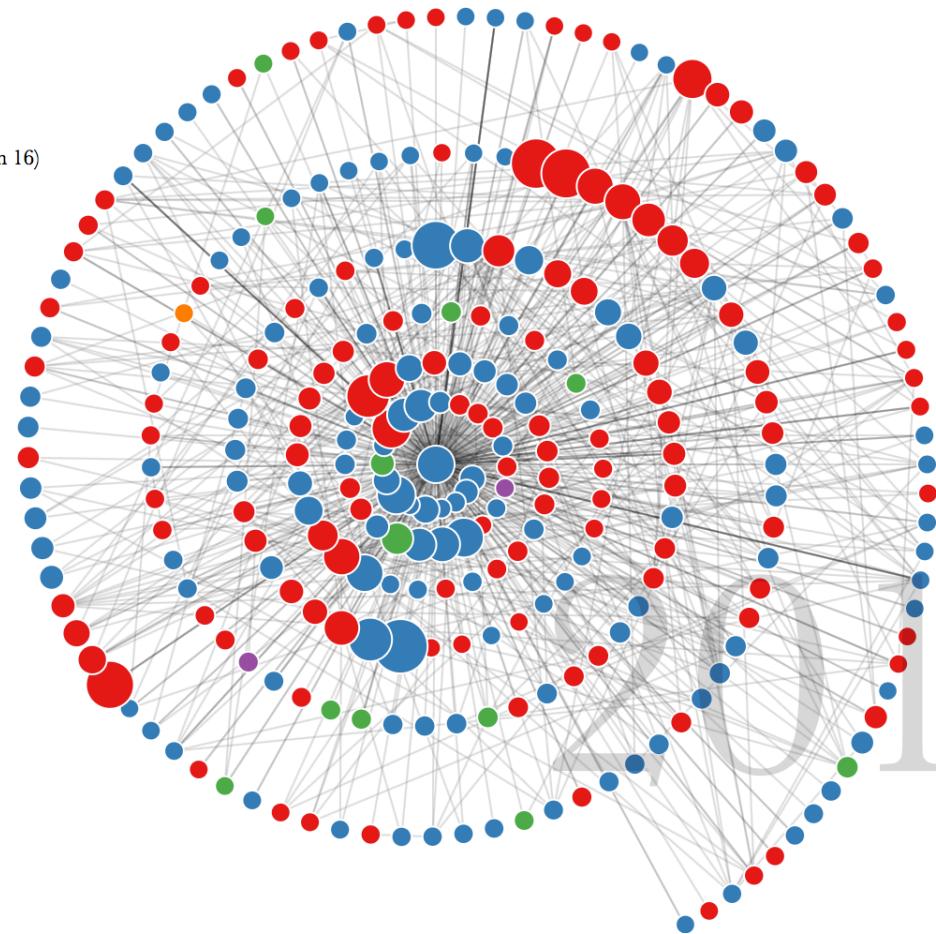
[Learn More](#)

- █ 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)

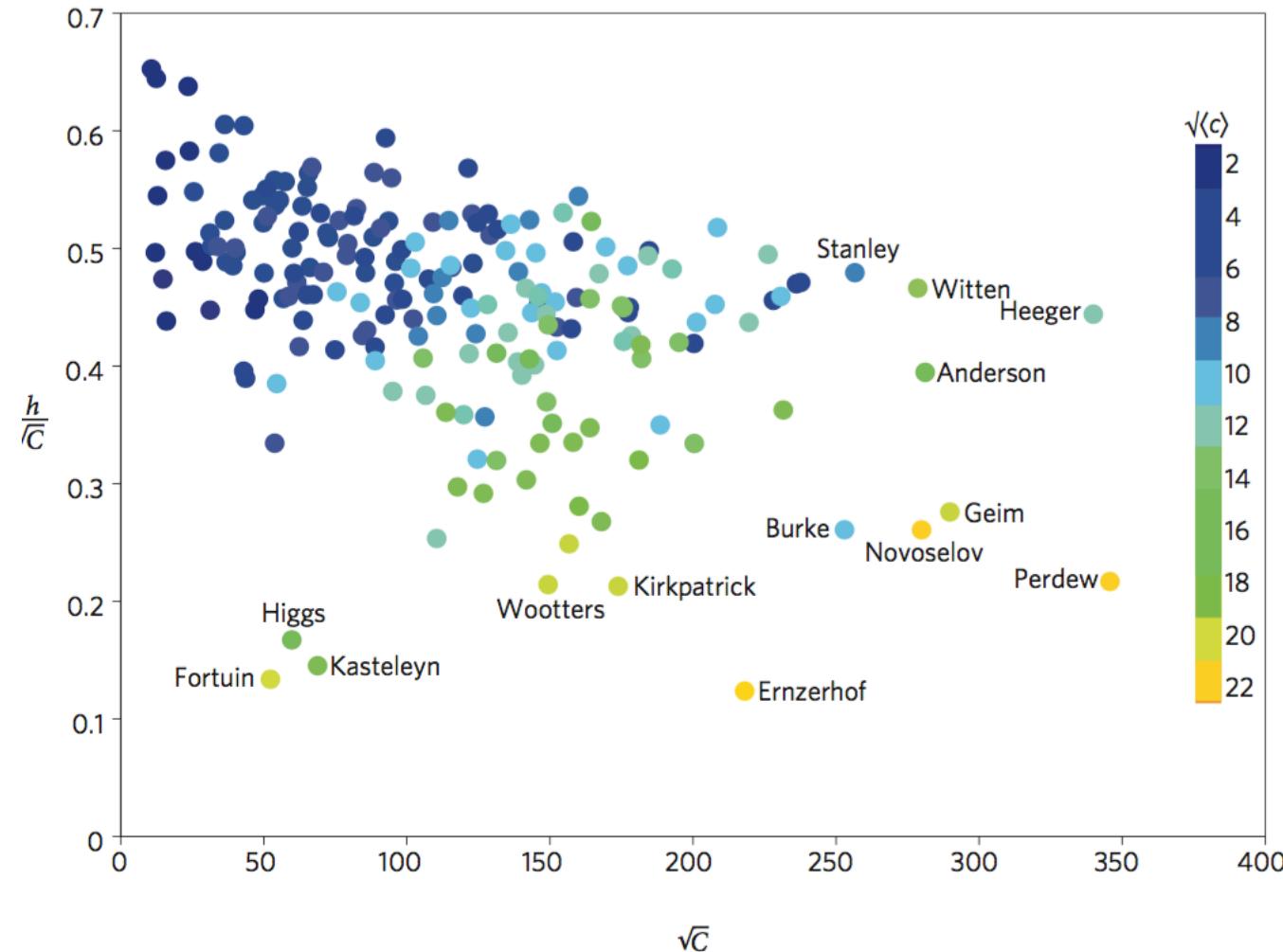
Roberta A.
Gottlieb



Pew Scholar
1997



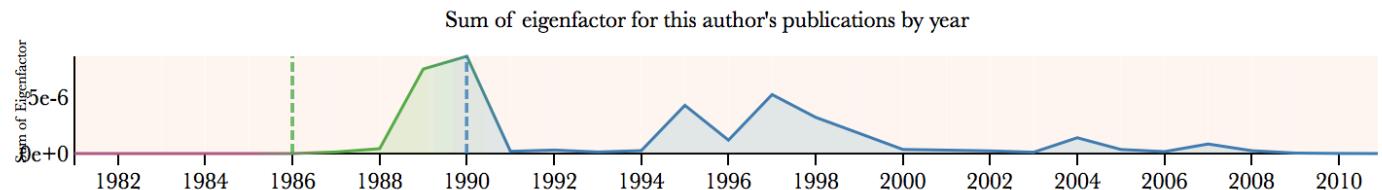
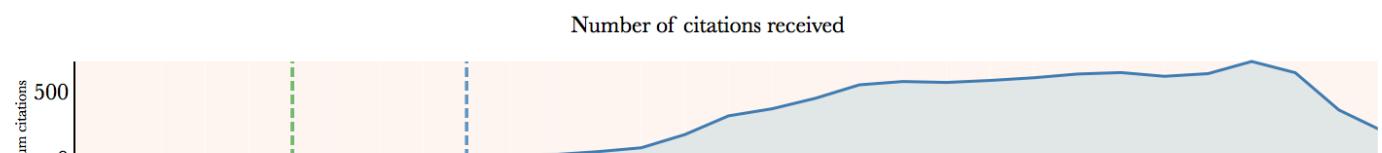
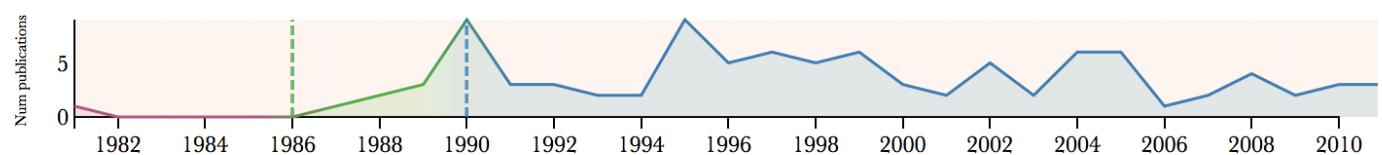
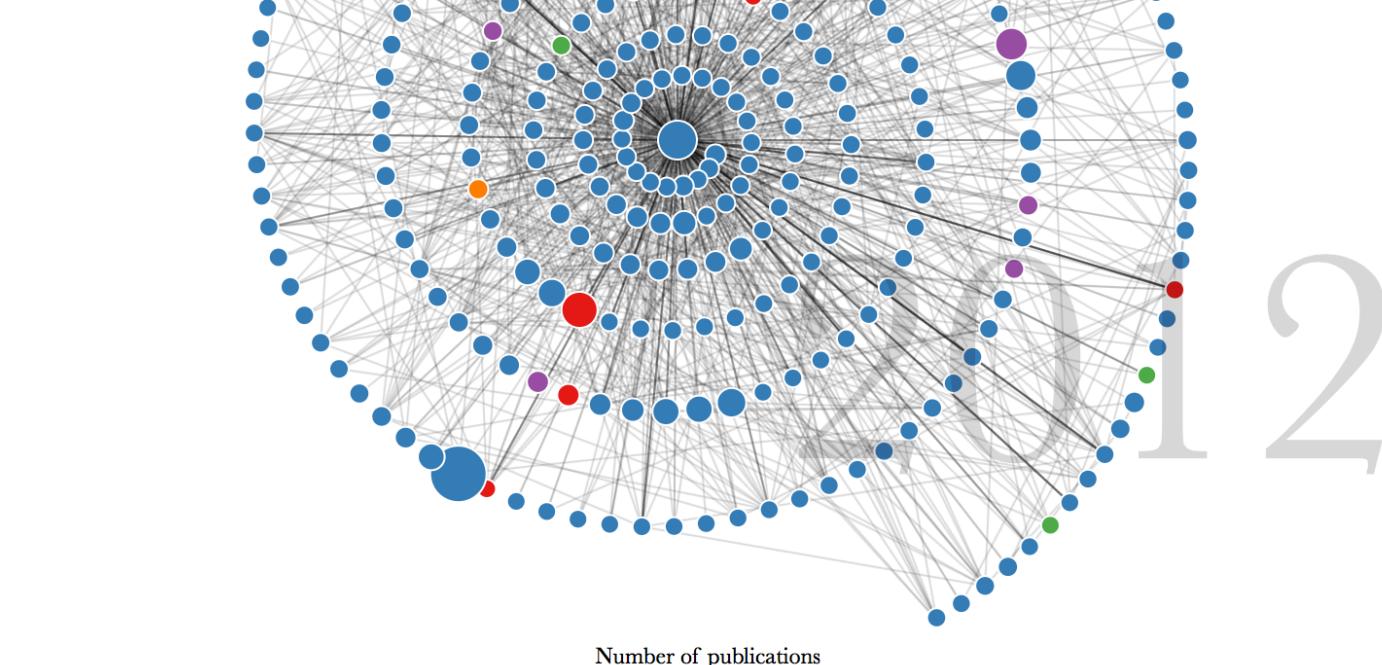
An evisceration of the H-index...



Philip A.
Hieter



Pew Scholar
1986



Visualizing Scholarly Influence Over Time

Influence of Pew Scholars

Mark W. Grinstaff

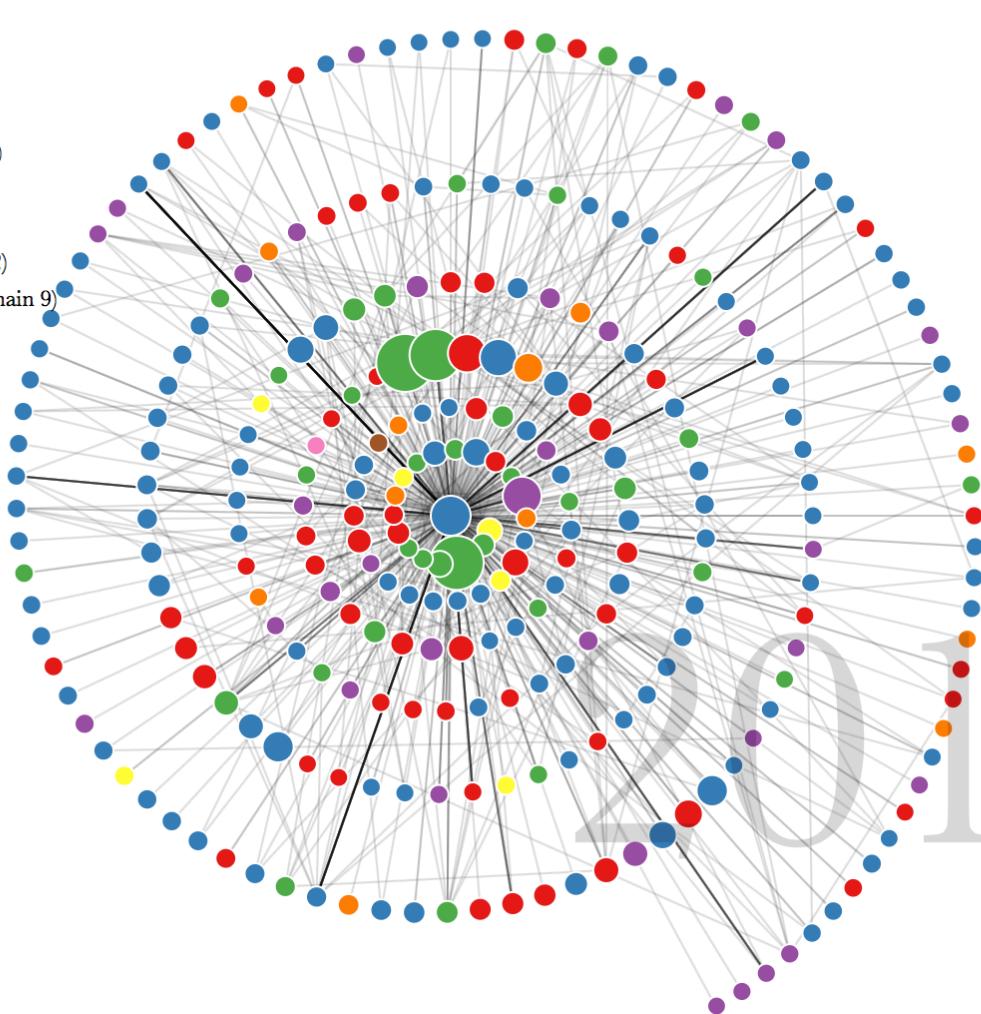
[Learn More](#)

- 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

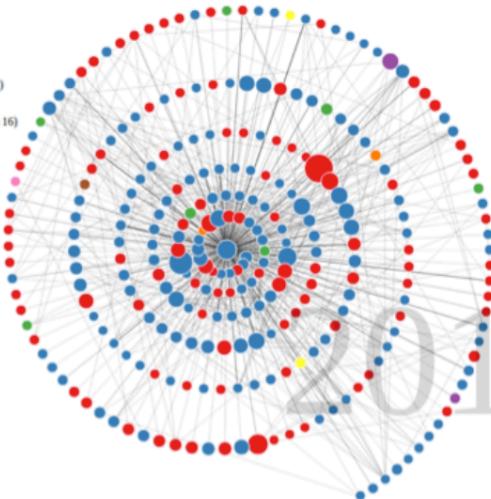


Pew Scholar
1999



Comparing Authors

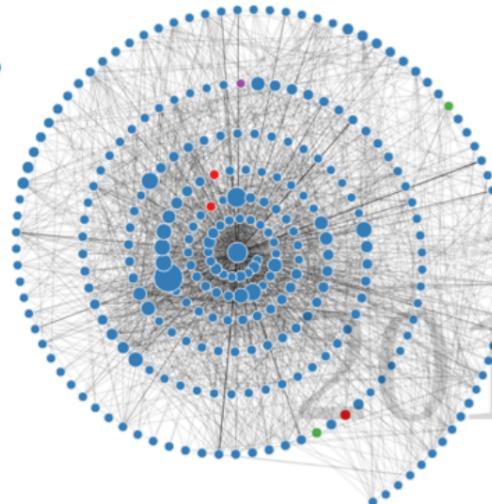
- Papers in category "Medicine" (domain 6)
- Papers in category "Biology" (domain 4)
- Papers in category "Chemistry" (domain 5)
- Papers in category "Engineering" (domain 8)
- Papers in category "Material Science" (domain 12)
- Papers in category "Physics" (domain 19)
- Papers in category "Agriculture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

- Papers in category "Biology" (domain 4)
- Papers in category "Medicine" (domain 6)
- Papers in category "Chemistry" (domain 5)
- Papers in category "Social Science" (domain 22)

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

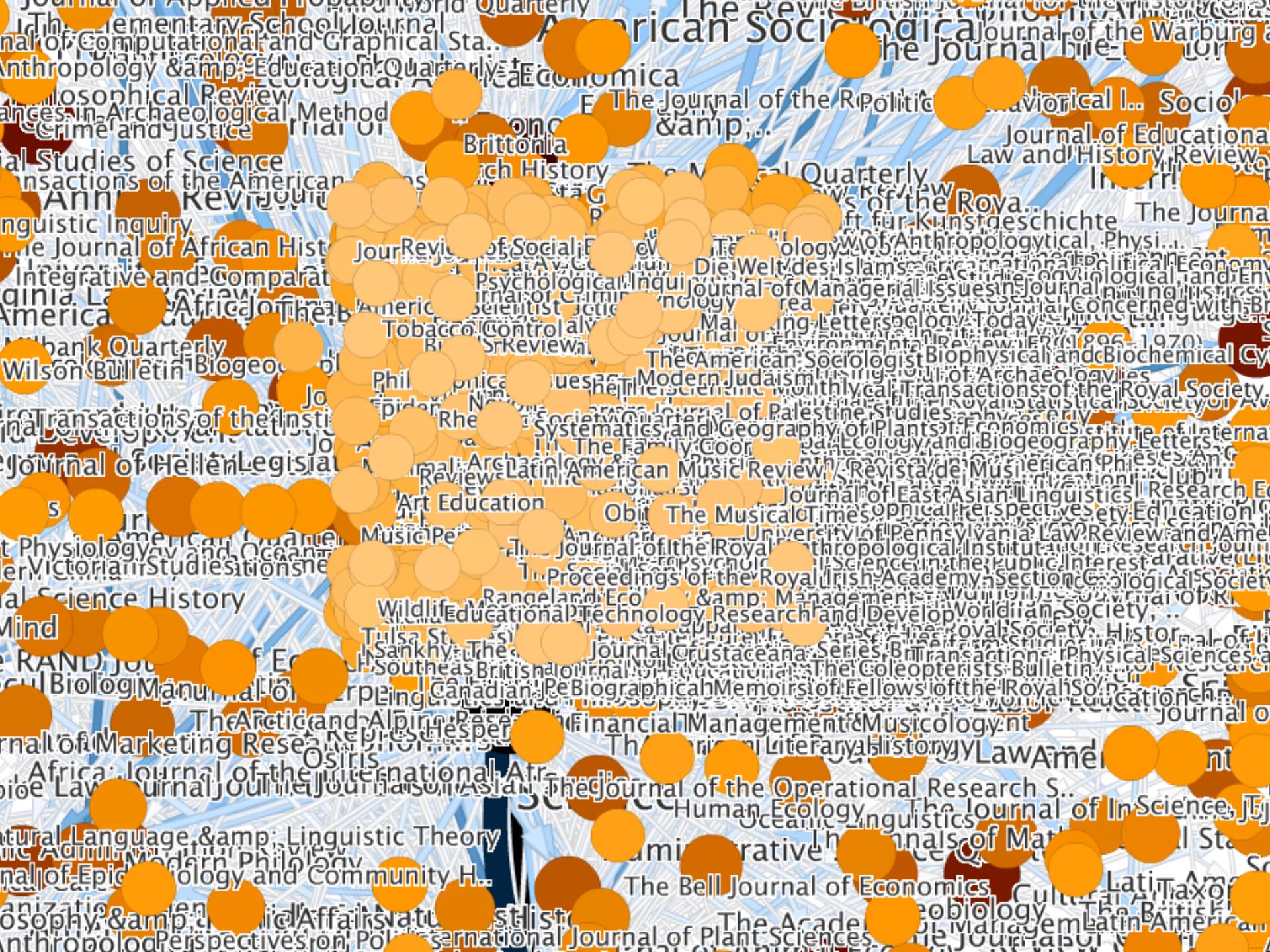


The map equation

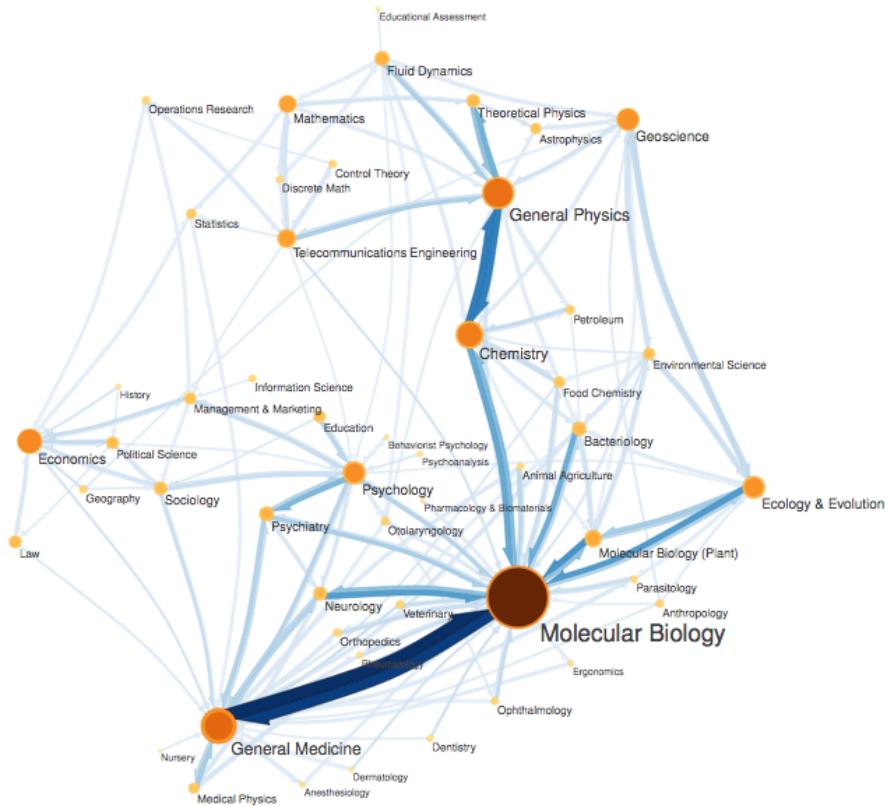
frequency of inter-module movements

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

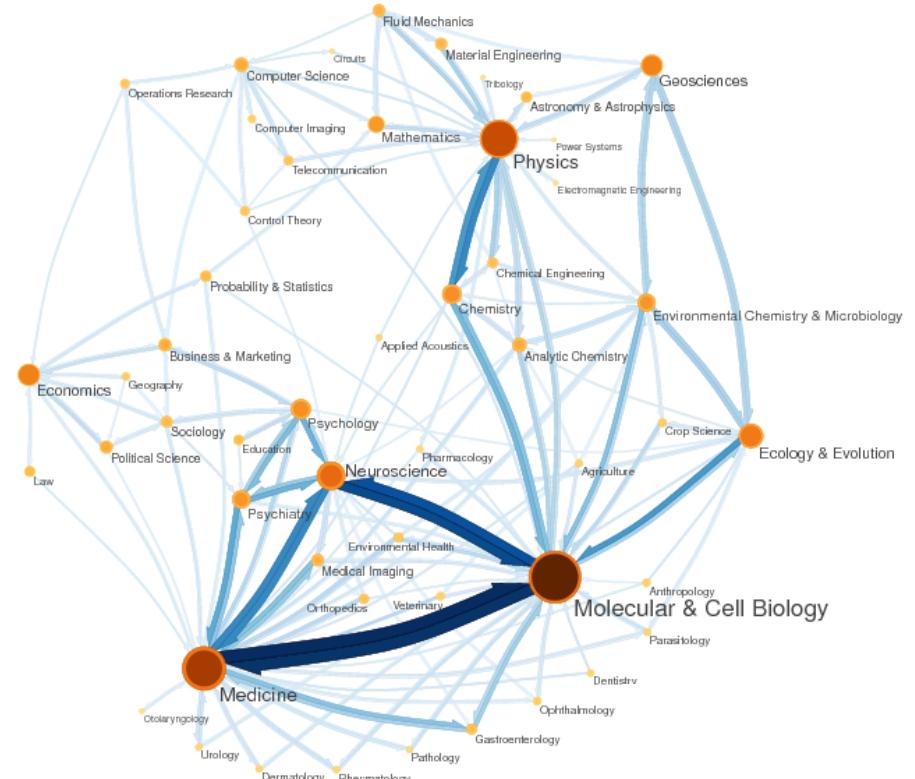
The diagram illustrates the components of the map equation. It features a central equation: $L(M) = q_{\sim} H(Q) + \sum_{i=1}^m p_i^{\circ} H(P^i)$. To the left of the first term, a blue vertical bar is labeled "frequency of inter-module movements". To the right of the second term, a red vertical bar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



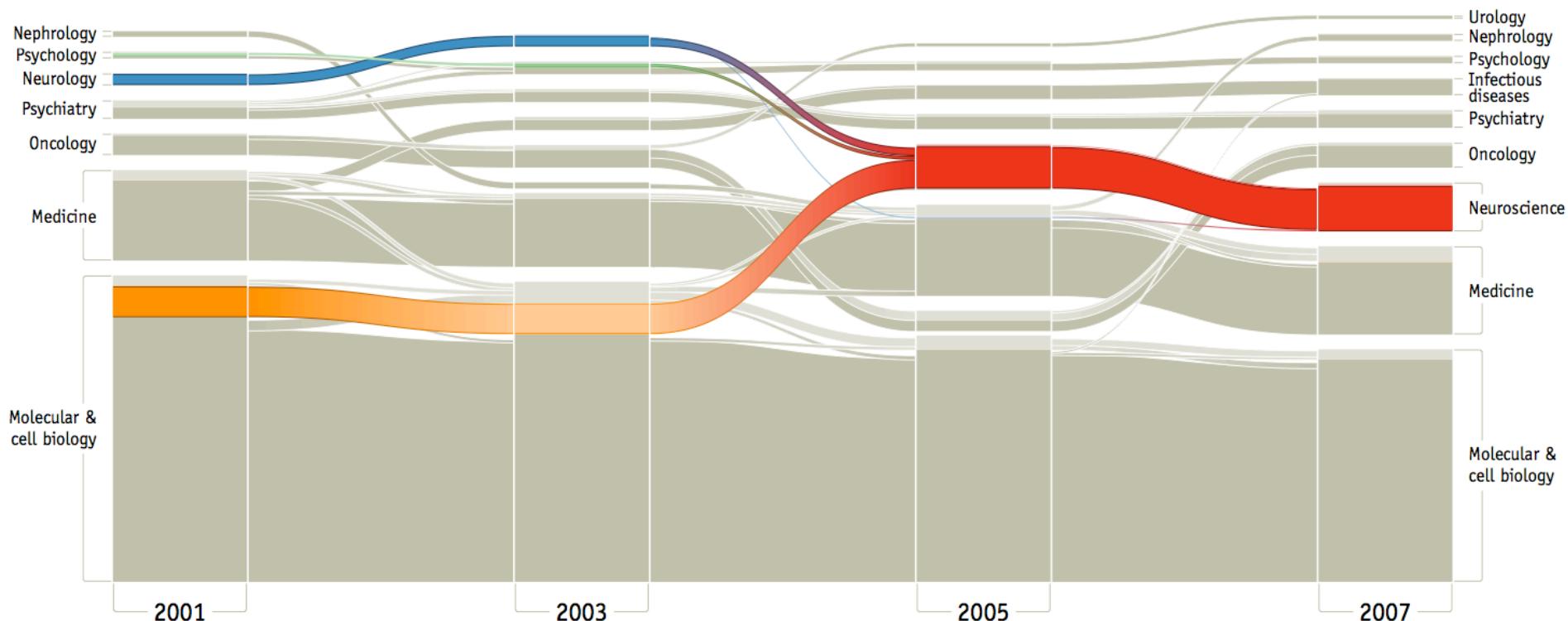
1995



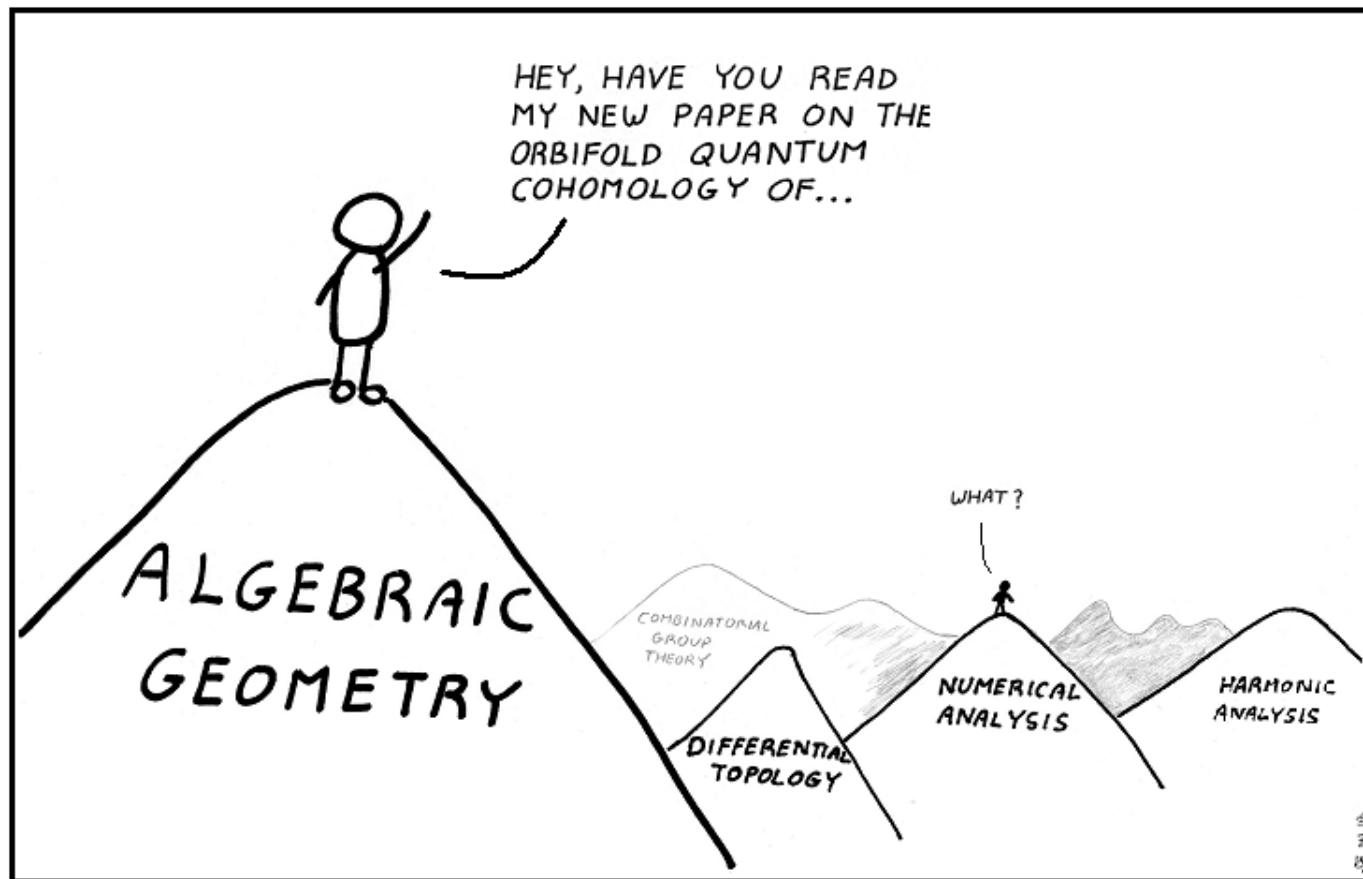
2004



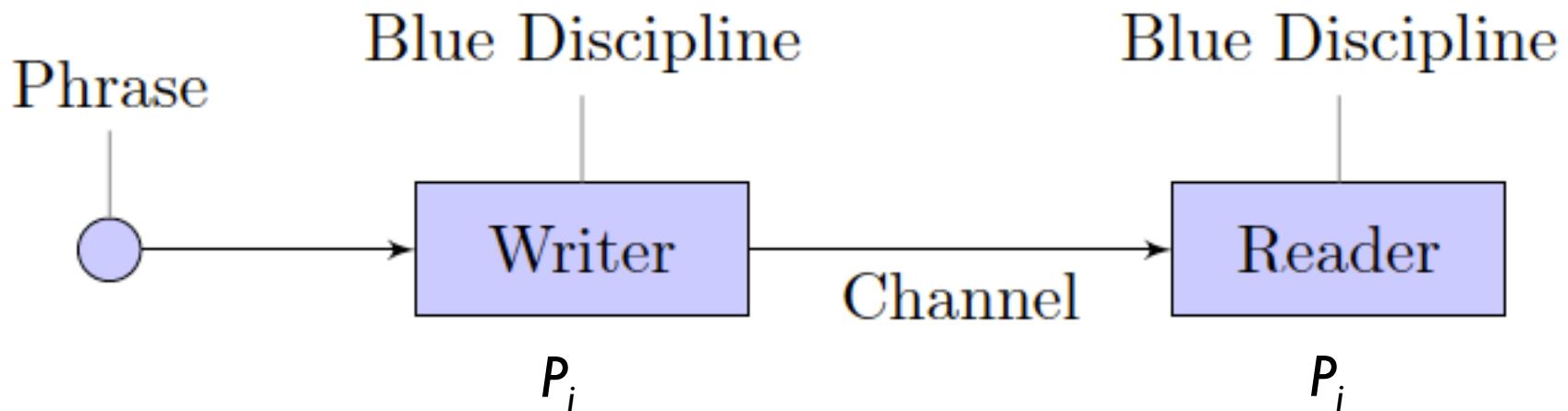
The Emergence of Neuroscience



The jargon barriers of science



The Landscape of Modern Mathematics



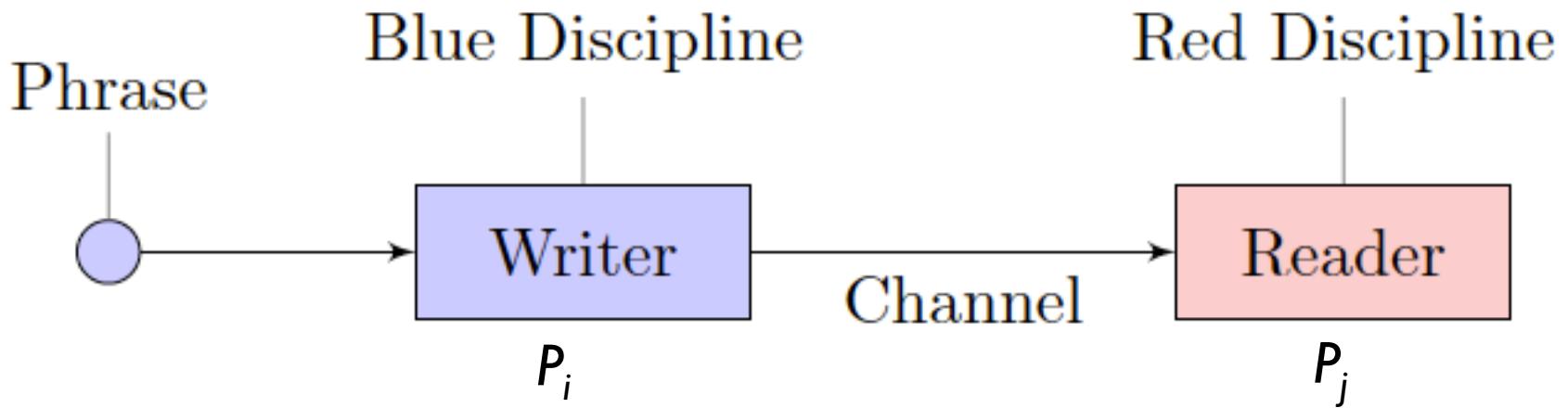
$X \sim$ space of all phrases

$P_i \sim$ 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: $Q(p_i||p_j) = - \sum_{x \in \mathcal{X}} p_i(x) \log_2 p_j(x)$

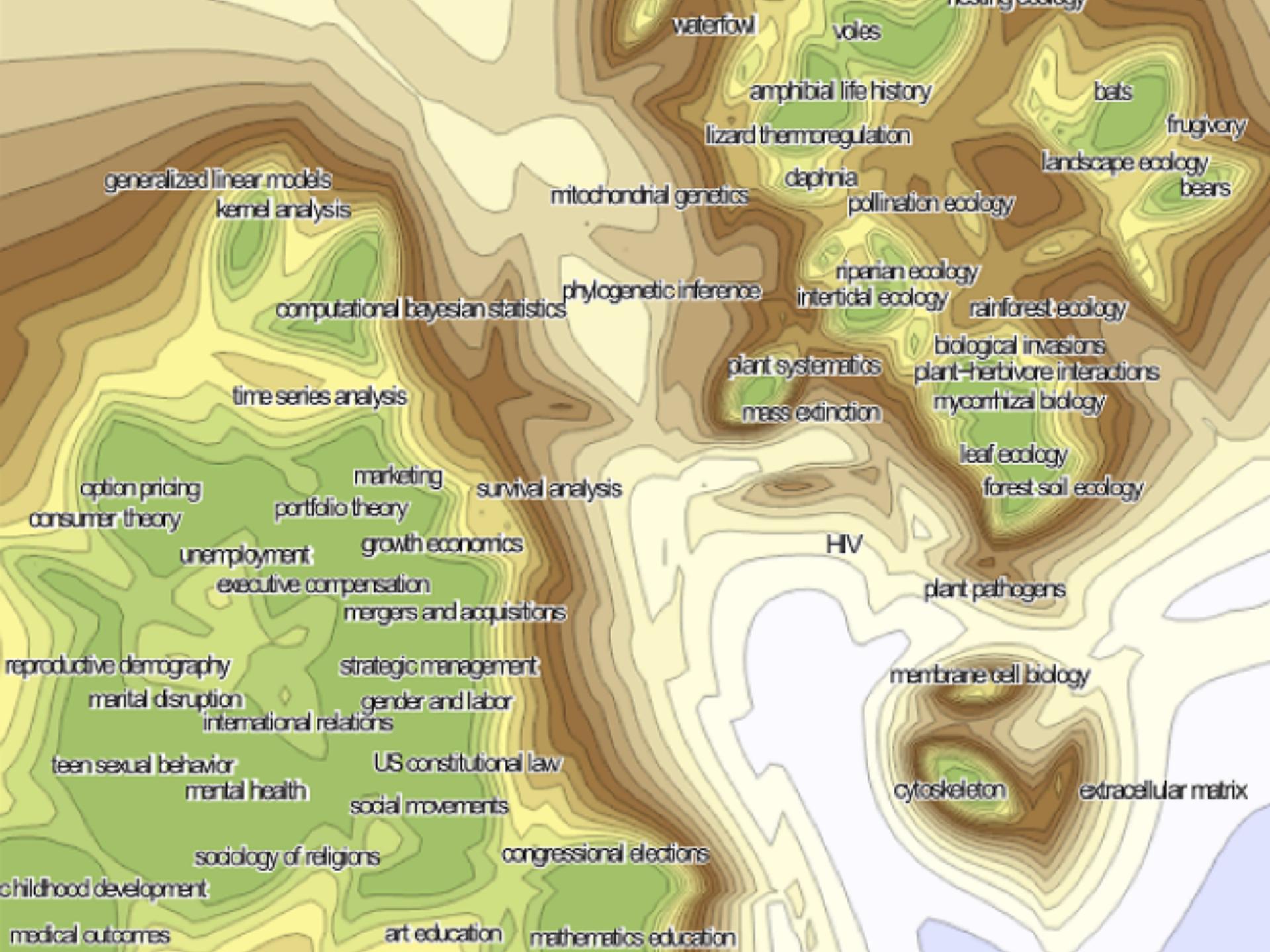
efficiency of communication

$$\downarrow$$

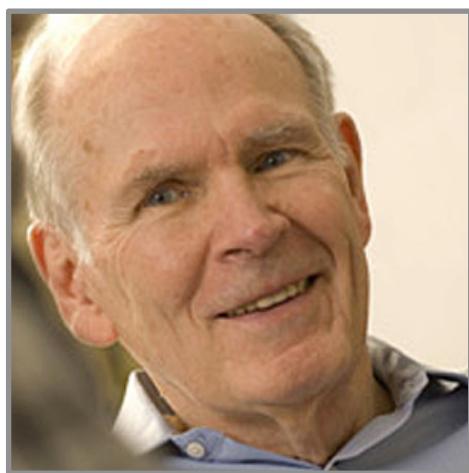
$$E_{ij} = \frac{H(X_i)}{Q(p_i||p_j)} = \frac{-\sum_{x \in \mathcal{X}} p_i(x) \log_2 p_i(x)}{-\sum_{x \in \mathcal{X}} p_i(x) \log_2 p_j(x)}$$

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

↑
 cultural hole



Translational Lag



A Life History of Innovation

inventors

Nature, PNAS, Cell



1980

UW licenses to WRF



lawyers

1987

Patent Application



examiners

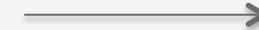
Patent Granted



1995

managers

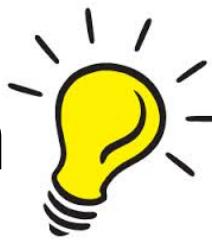
\$240 million in revenue



2014

\$31 million WRF -> UW



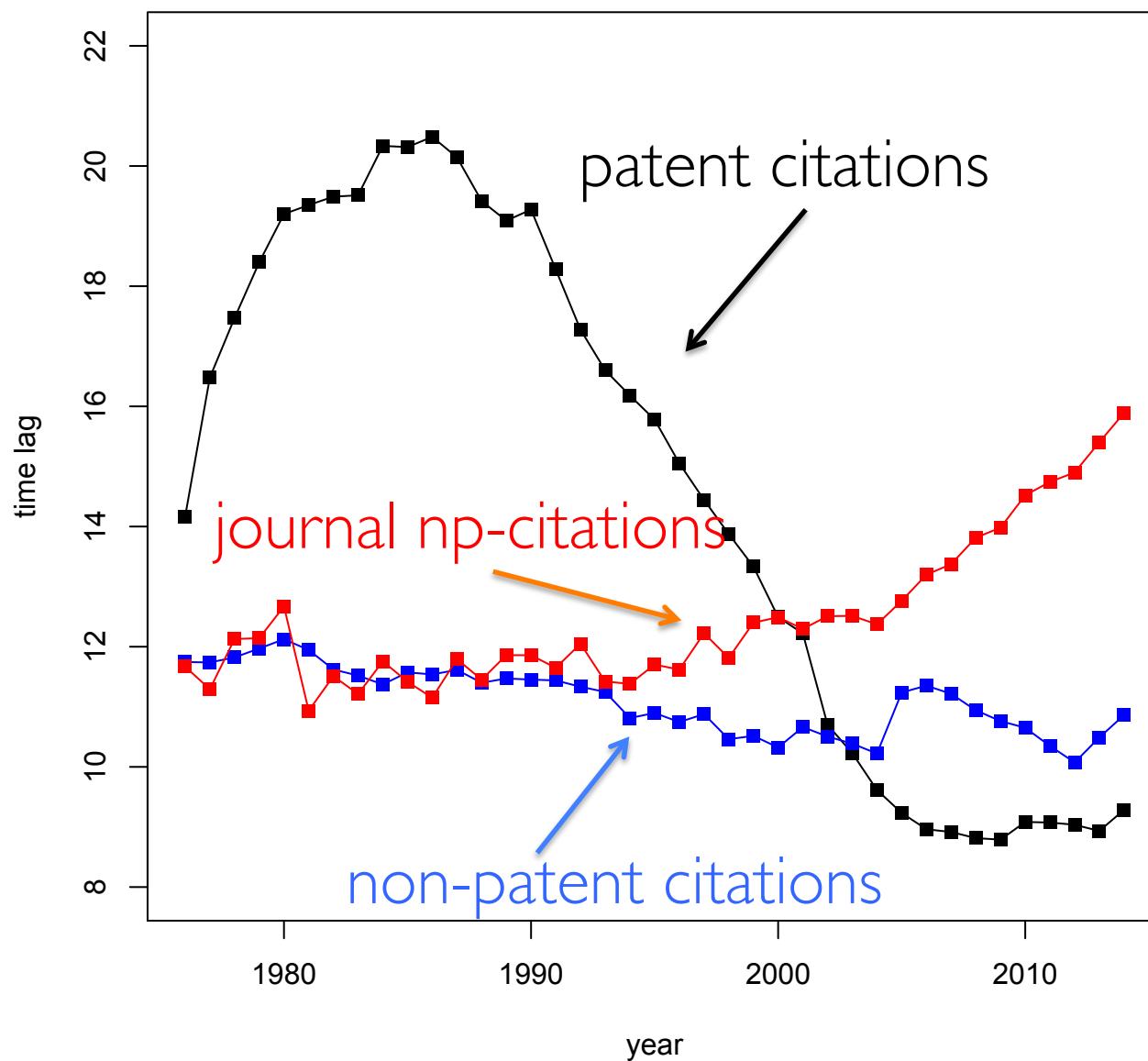
What is the time from  to patent?

Methods

- Extracted 111,038,761 citations (1976 – 2014)
 - 92,680,292 patent citations
 - 18,358,469 non-patent citations
 - 761,759 non-patent, journal citations
 - 20,470,405 examiner citations (2001 – 2014)
- Extracted years from patents and non-patents

Svensson, K. et al. *Evolution of subspecies of*. *Journal of Bacteriology*, Jun. 2005, vol. 187, No. 11, pp.3903-3908
- Only use first, unique instances of citing patent

Are patents becoming more myopic?



viziometrics.org

The screenshot shows a search interface with the following elements:

- Header:** VizioMetrics, About, Search, Crowdourcing.
- Search Bar:** Keywords or Cluster, Result Ordered by Random, with a dropdown for "Random".
- Filter Buttons:** Composite (checked), Equation, Diagram, Photo, Table, Visualization.
- Result Grid:** A 6x6 grid of cards, each representing a different type of visualization or document related to the search terms. Some cards are highlighted with colored borders (red, green, blue).
- Bottom Text:** A project of the eScience Institute at the University of Washington.

(b)

VizioMetrics About Search Crowdsourcing

Random Keywords or Cluster, Result Ordered by Random

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

The screenshot displays a search interface with filters for 'Random' and 'Keywords or Cluster, Result Ordered by Random'. Below are several data cards:

- A card titled 'On the Number of New World Founders: A Population Genetic Portrait of the Peopling of the Americas' by Clark Andy G. and Hey Jody, featuring a green border.
- A card titled 'Population Size Change Rate Years' with a green border.
- A card titled 'Population Size Change Rate Years' with a grey border.
- A card titled 'Population Size Change Rate Years' with a grey border.
- A card titled 'Malaria in American troops in the South and Southwest Pacific in World War II.' with a red border, showing three thumbnail images of historical documents.

The interface includes a search bar, a progress bar, and a sidebar with icons for file operations.

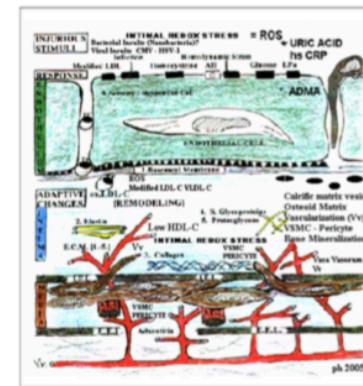
A project of the eScience Institute at the University of Washington

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 Ramesh, Sowers James R, Kolb Lisa, Tyagi Suresh C and Hayden Melvin
Cardiovascular Diabetology 2005

Abstract

| Hide Abstract | View Paper | View Cluster

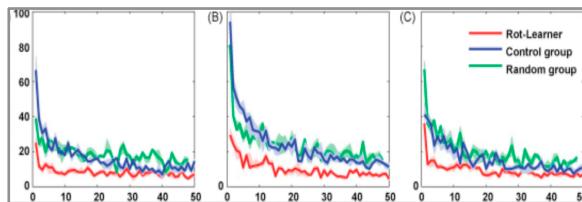
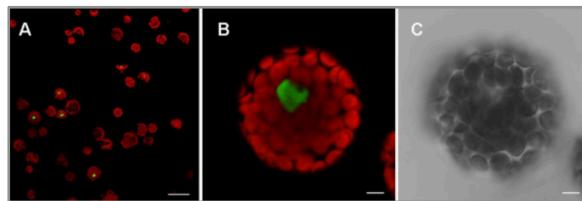
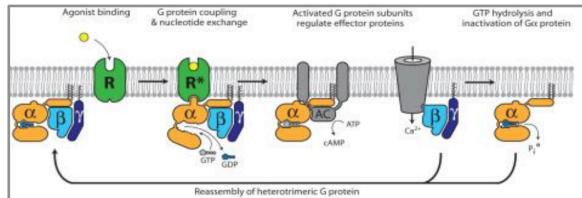


The central role of the endothelium in VOC and atherosclerosis. These images portray the endothelium as the first line of defense against multiple injurious stimuli. Most of the injurious stimuli are associated with the vascular wall. In addition to the vascular wall damage associated with the endothelium, the vascular wall is also damaged by the presence of VOCs. As mentioned earlier, VOCs are important in the initiation of atherosclerosis. In addition to the vascular wall damage associated with the endothelium, the vascular wall is also damaged by the presence of VOCs. As mentioned earlier, VOCs are important in the initiation of atherosclerosis.

[Hide Other Figures](#) | [Show Related Figures](#)



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



	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)

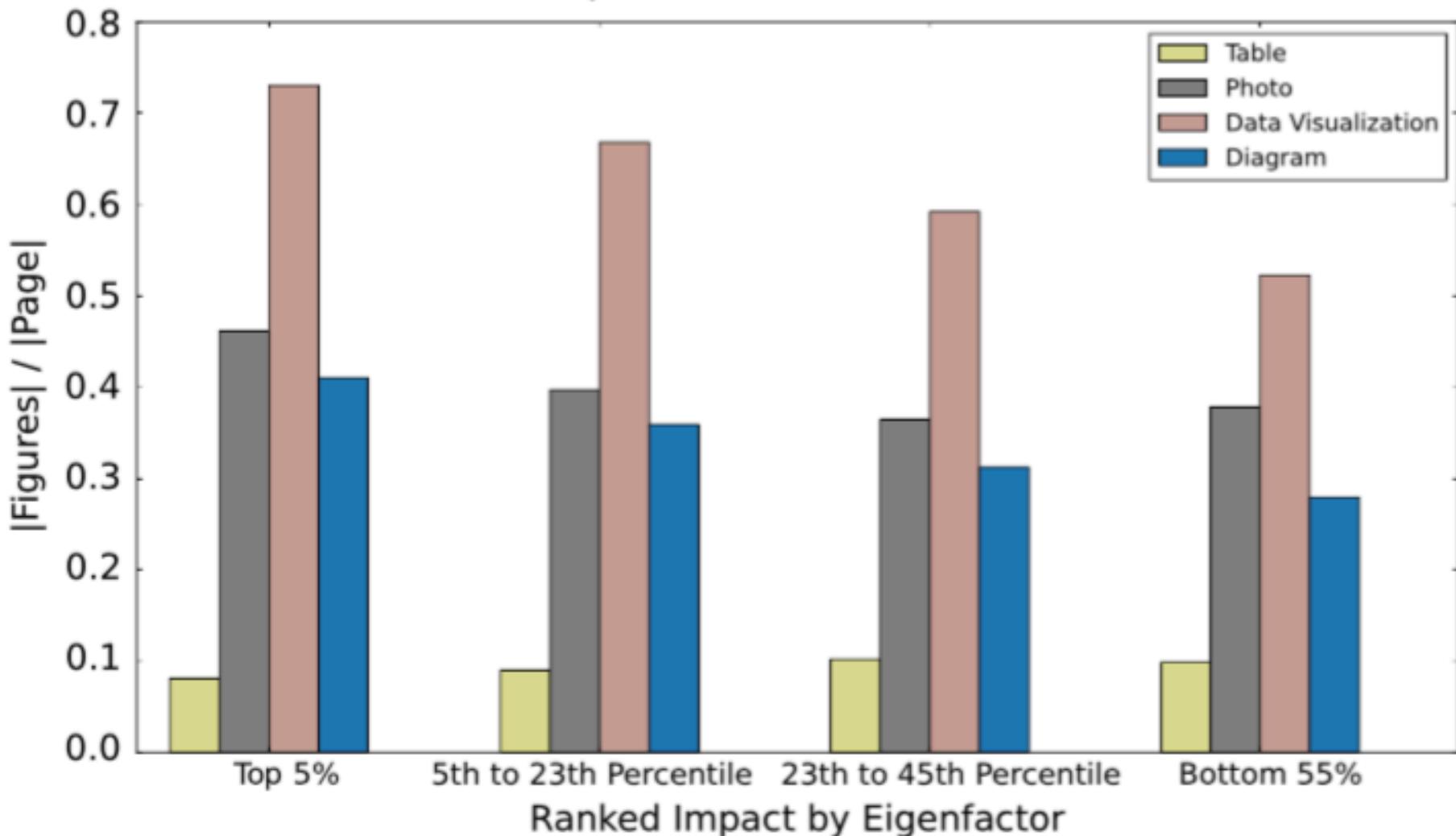
Schematics (769)

Photos (782)

Plots (890)

Tables (436)

Impact versus Figure Density



Summary

- Assembled scientific knowledge into a machine readable form
- Interrogating science from the meta-view
- Trying to better understand the origins of scientific ideas
- Building tools that facilitate discovery at the interdisciplinary boundaries of science



The Scholarly Graph



What else should we ask of this data?

PatentVector™



PNAS



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

