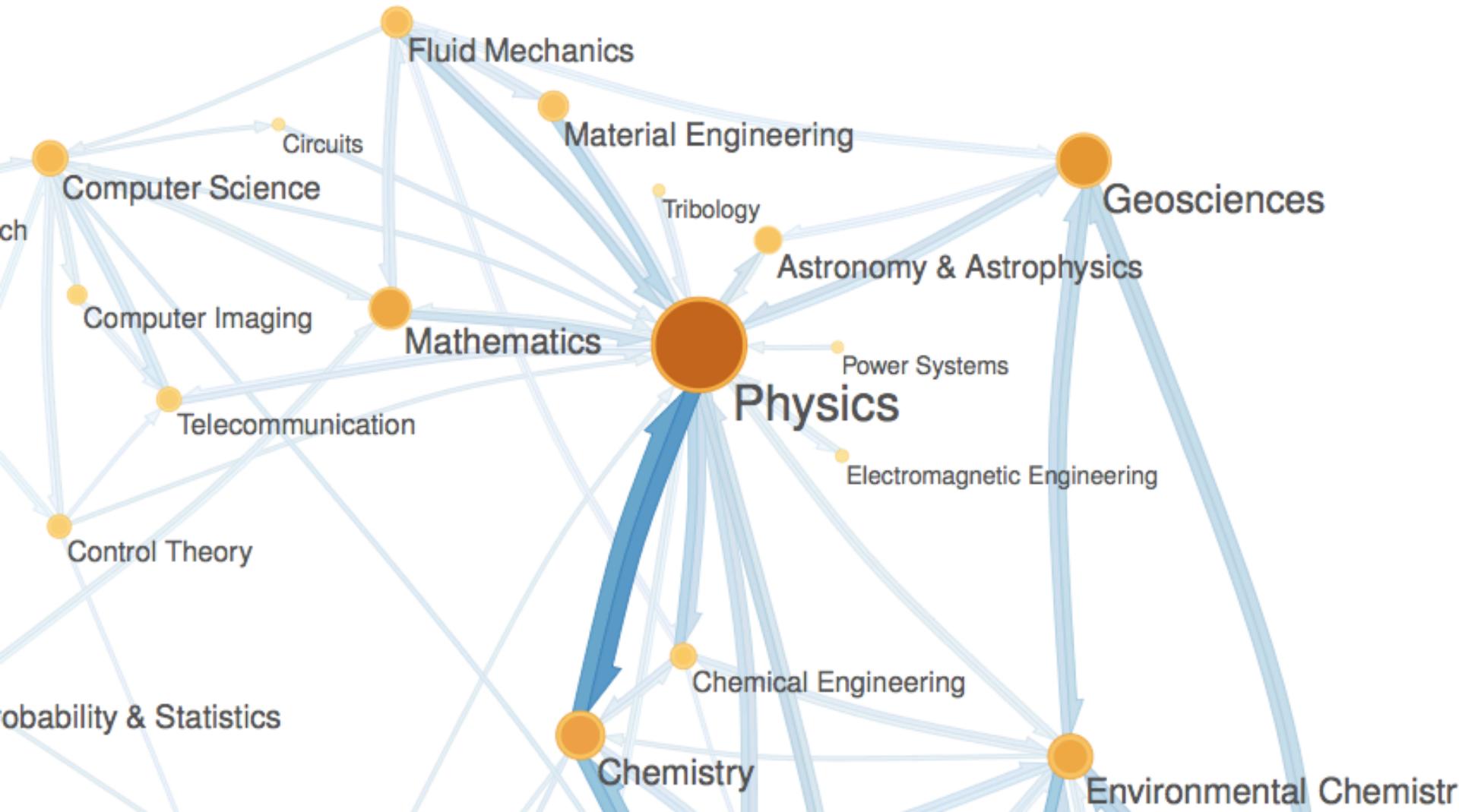
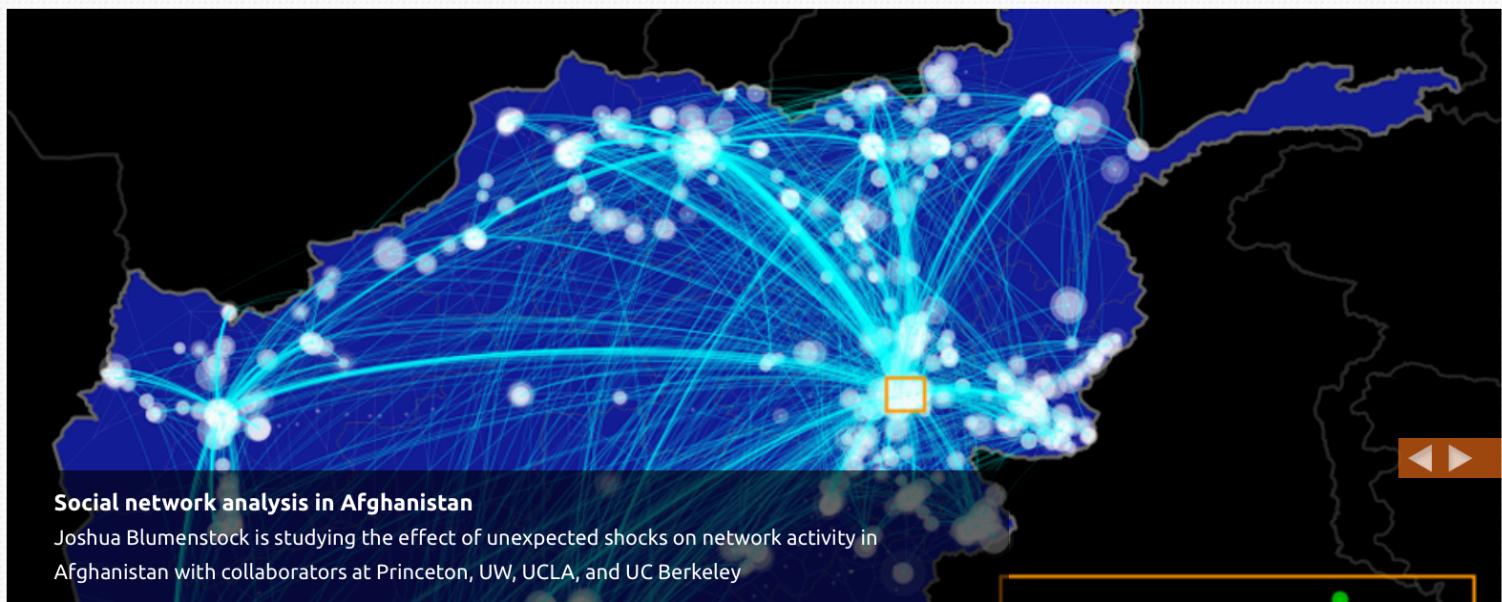


Mapping Knowledge Networks

Jevin West, Information School, University of Washington

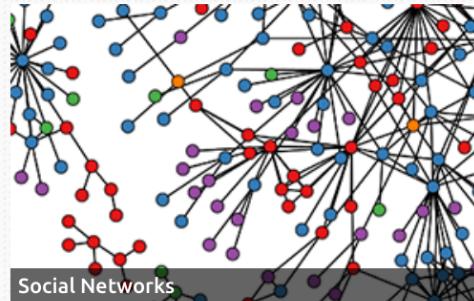




Research Focus Areas



Data for Development



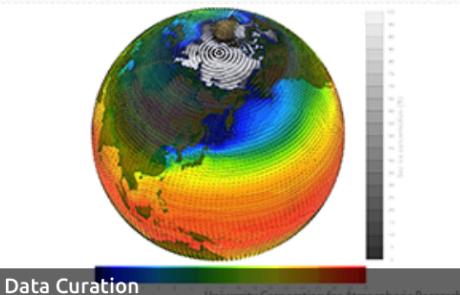
Social Networks



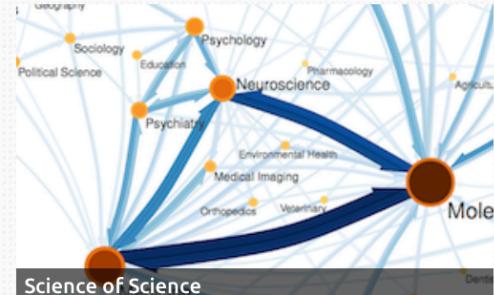
Data Visualization



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.

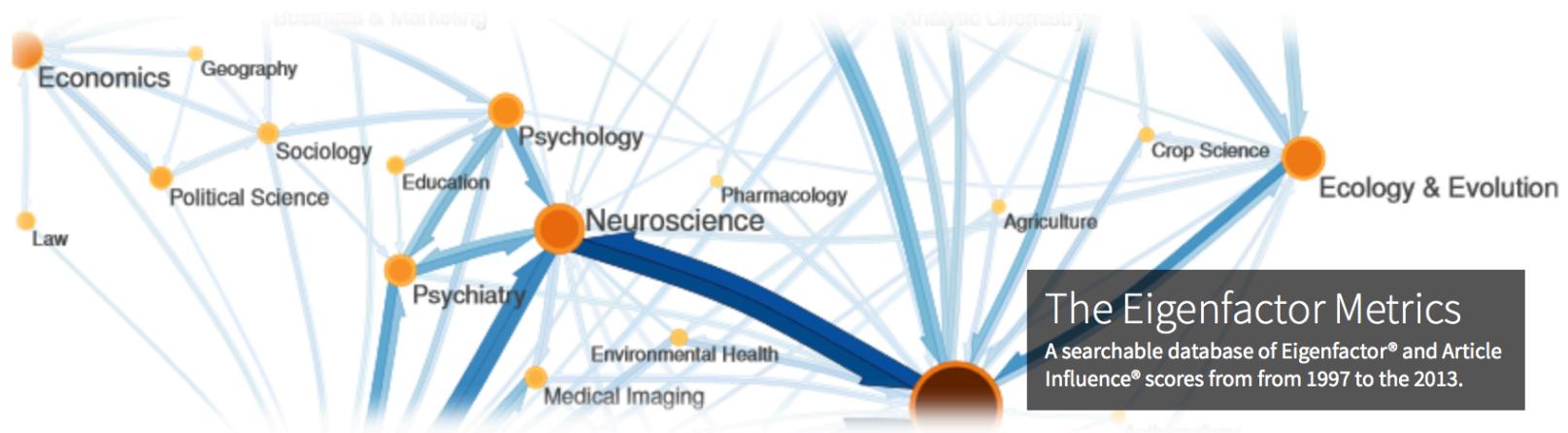
Latest eScience News

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

2 hours 4 min ago

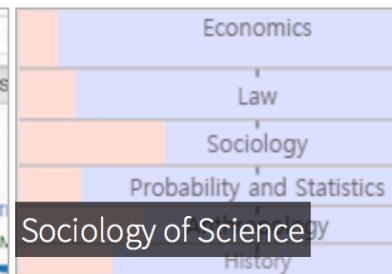
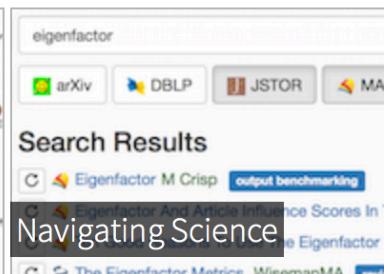
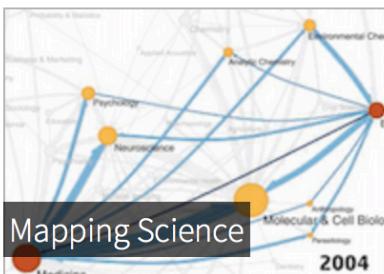
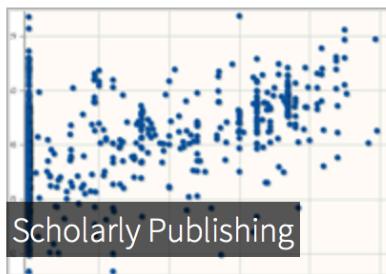
[Ben Marwick On How Computers Broke Science](#)

Search



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

RESEARCH AREAS



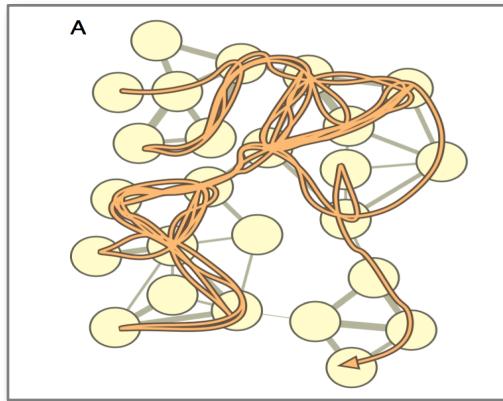
NEWS

23

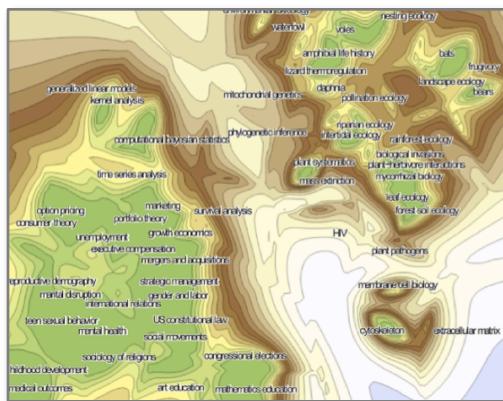
Nov. JEVIN WEST ON MEGAJOURNALS IN THE *CHRONICLE OF HIGHER EDUCATION*
Jevin West discusses the rise of the megajournal and our open access cost effectiveness tool in the *Chronicle of Higher Education*.

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 importance of individual articles. Using a mix of the article-level Eigenfactor algorithm and a deep learning model,

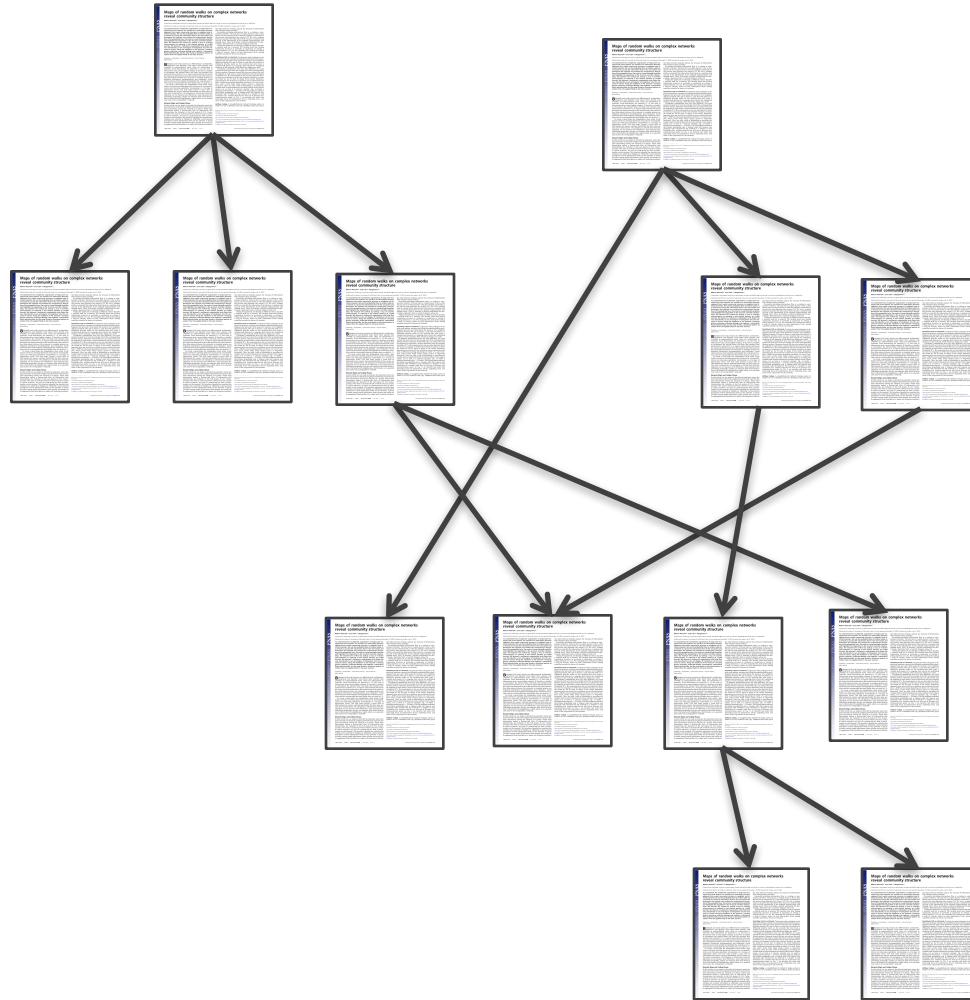


Science of Mapping



Mapping of Science

Citations form a vast network



de Solla Price, Science (1965)



The Scholarly Graph



PatentVector™



PNAS





The Scholarly Graph



Tens of millions articles, patents, books



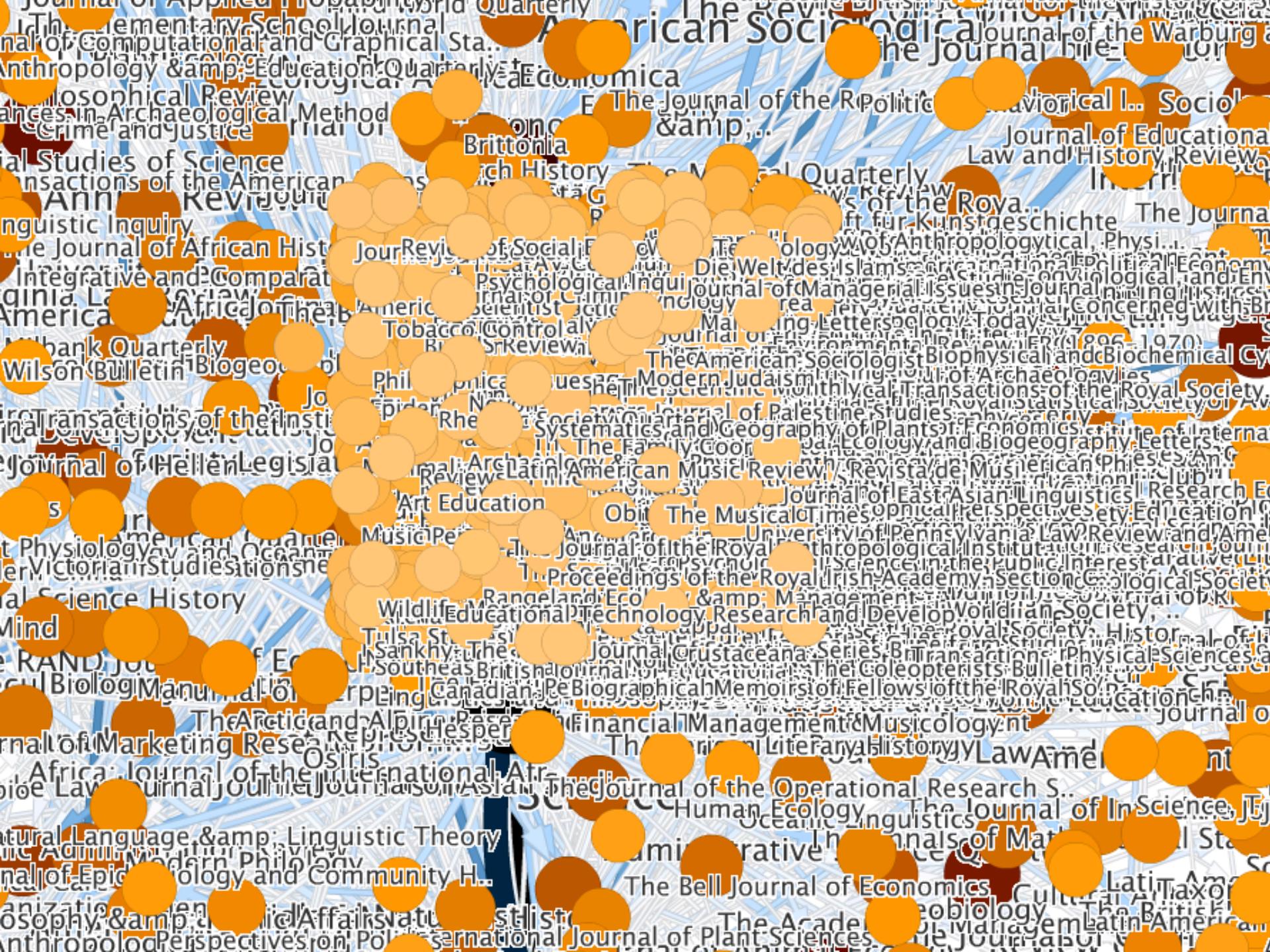
Billions of citation links

PatentVector™



Years: 1600 - 2016





Data

Compressing

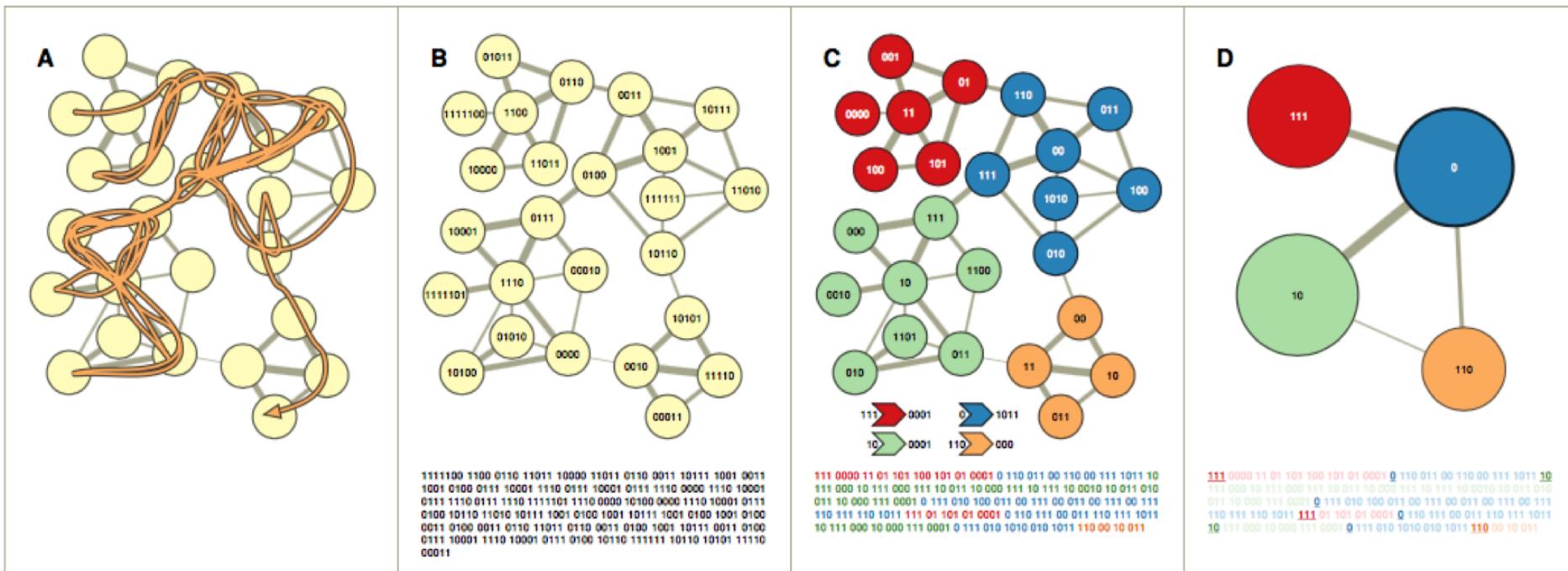


Finding patterns

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.

Minimal Description Length (MDL) Statistics

Finding regularities in the dynamics on networks



Rosvall and Bergstrom (2008) PNAS

compressing \longleftrightarrow finding patterns

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



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

The map equation

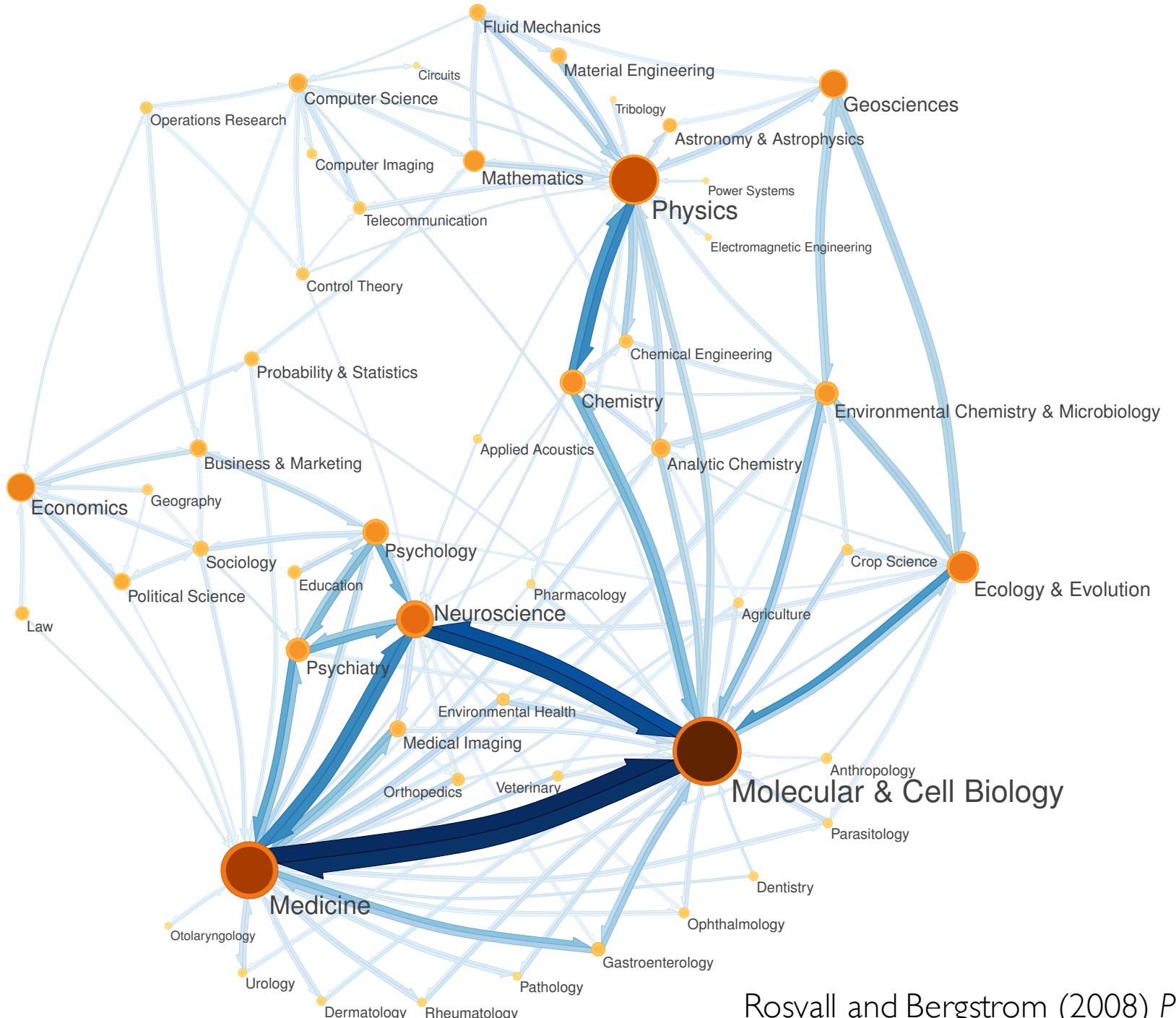
frequency of inter-module movements

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

frequency of movements within module i

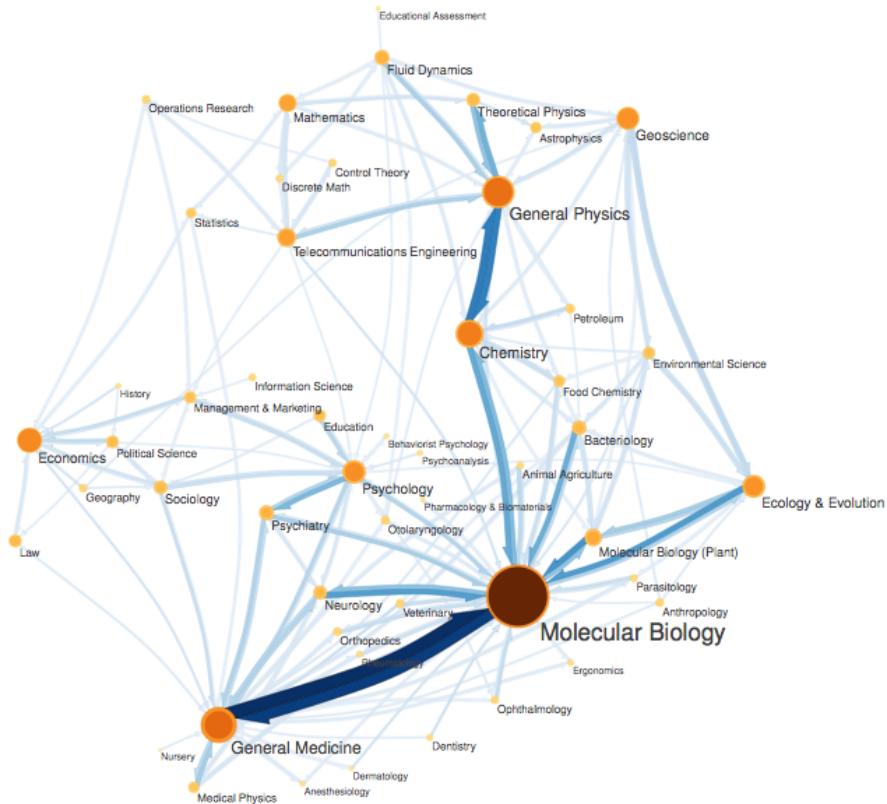
code length of module names

code length of node names in module i

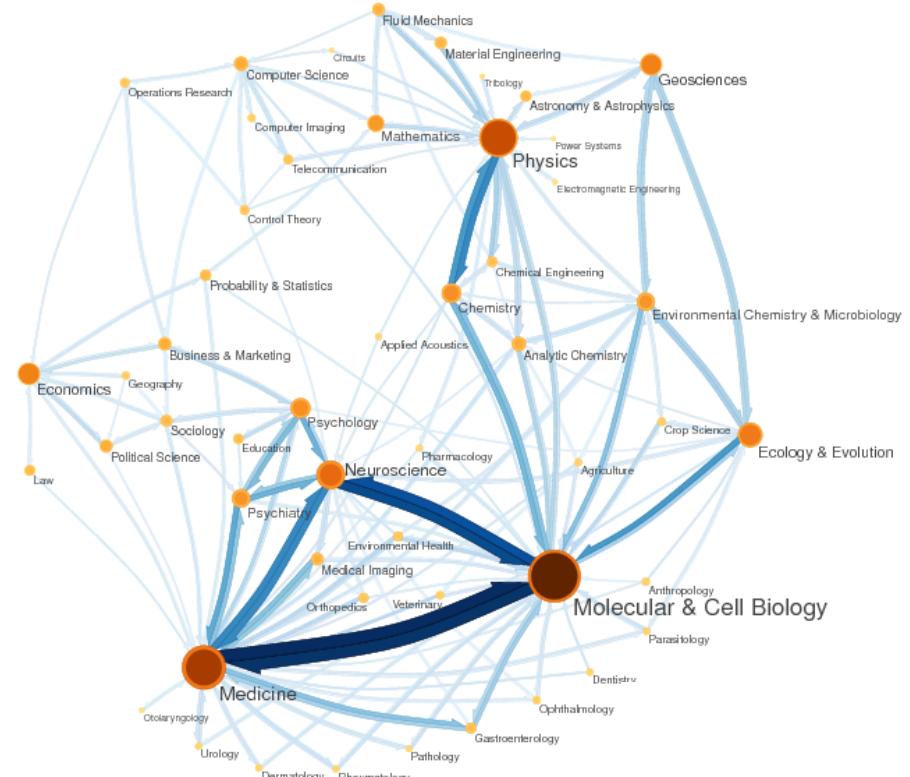


Rosvall and Bergstrom (2008) PNAS

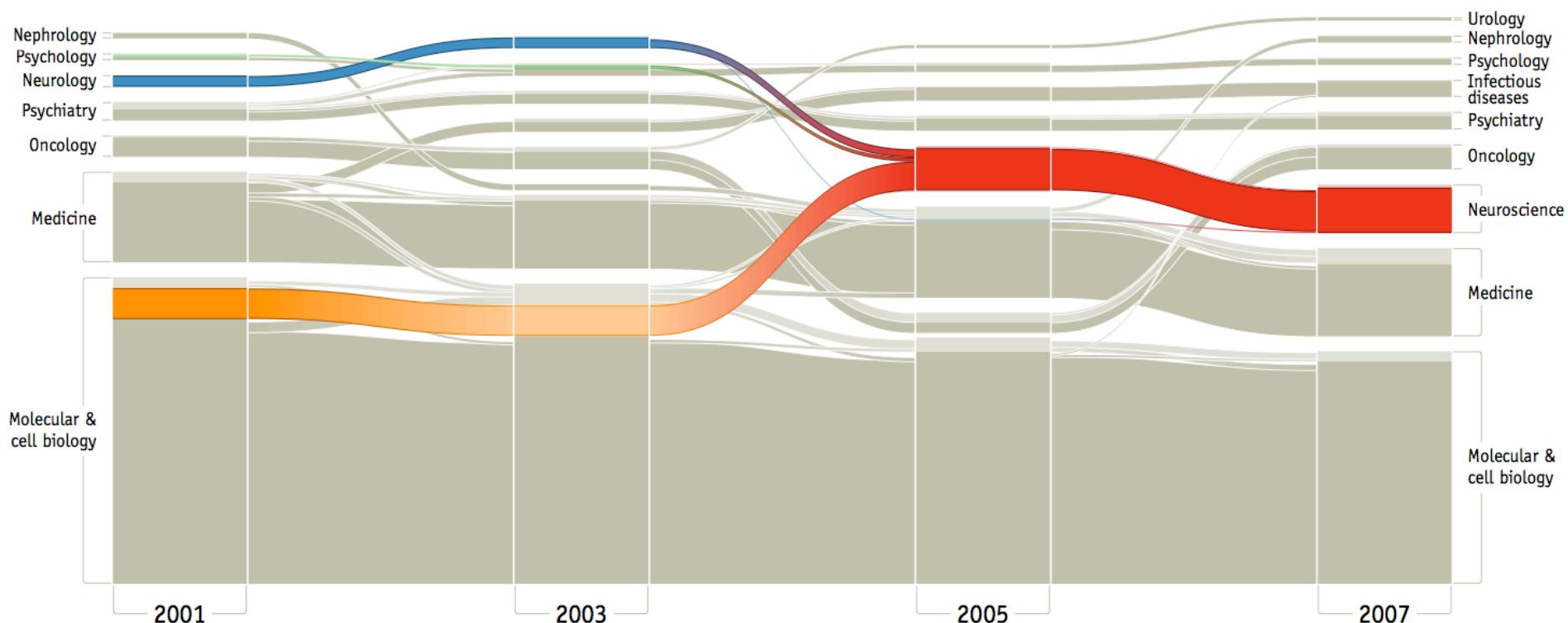
1995



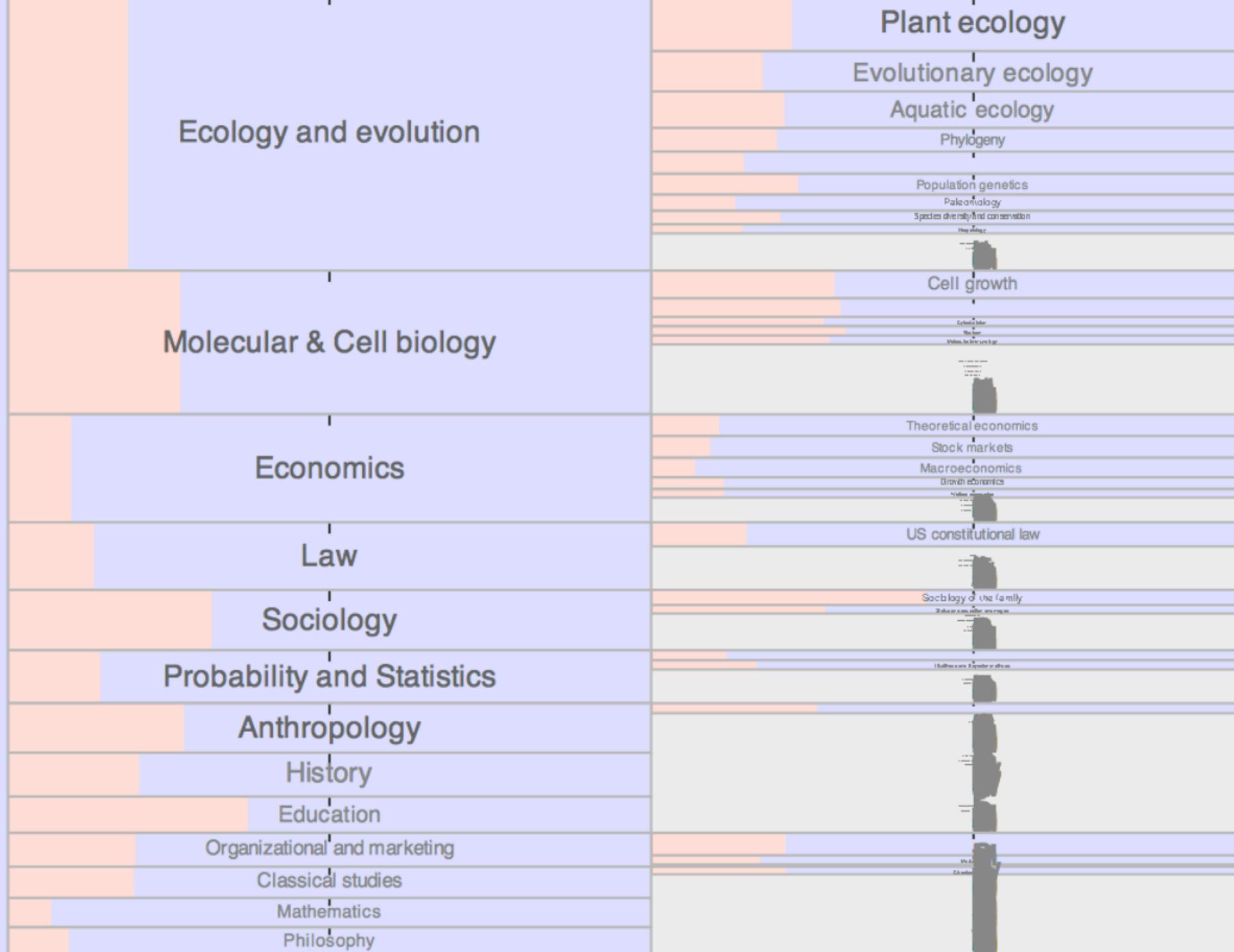
2004



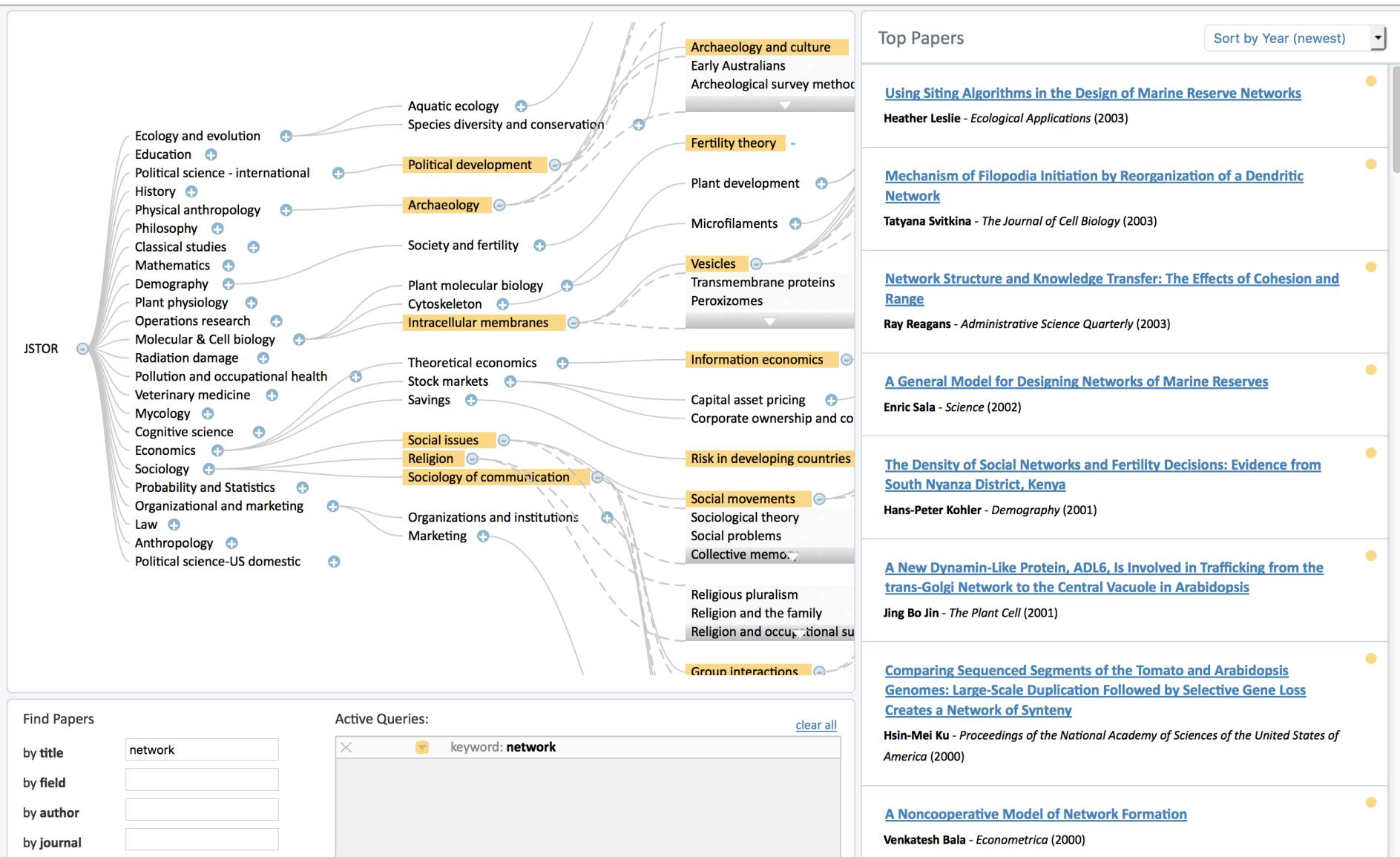
The Emergence of Neuroscience







“Network”



Next Steps

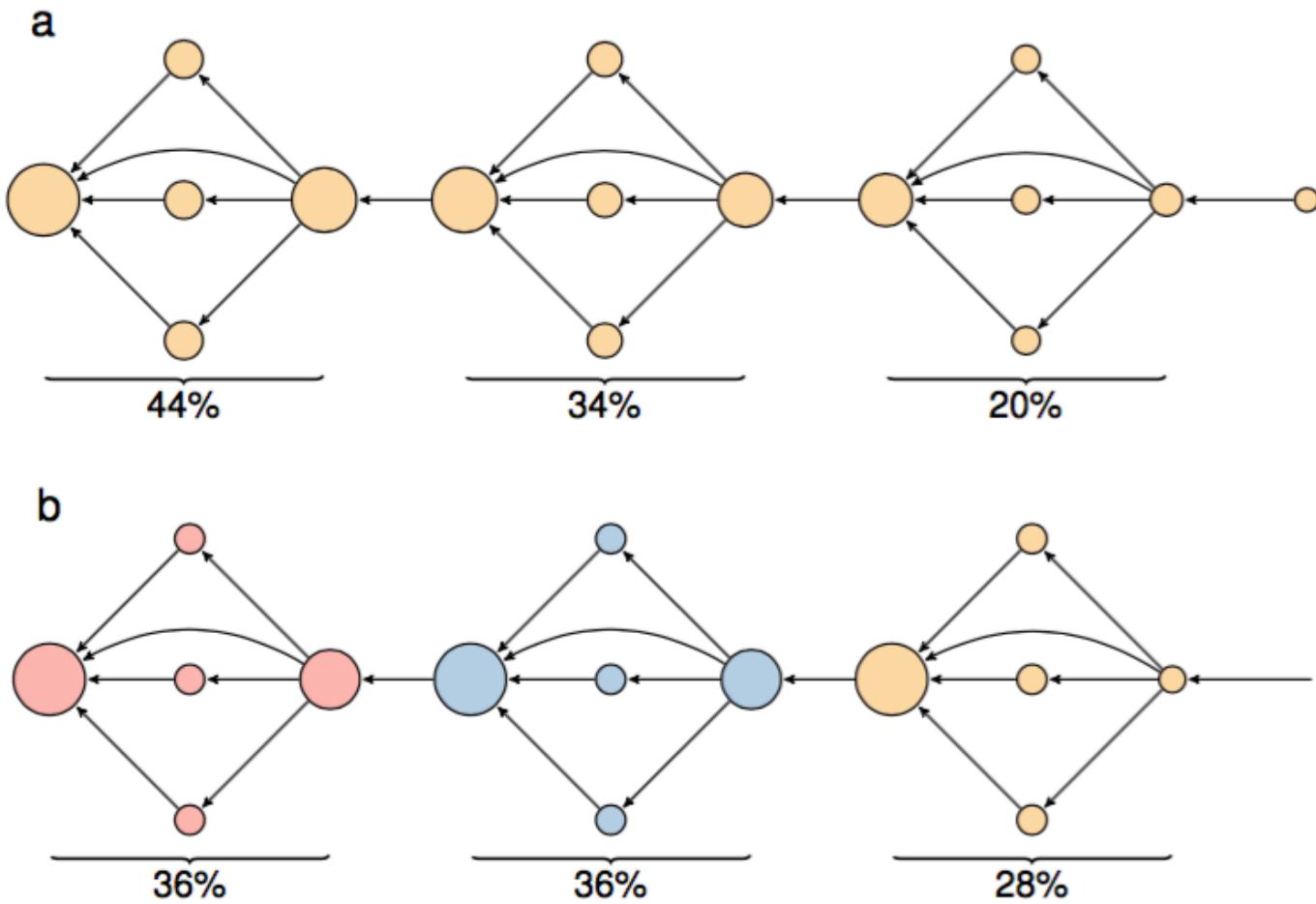
Higher Resolution Maps

Article-level Ranking

Author-centric Maps

Interactive Visualizations

Article-level Ranking and Mapping



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

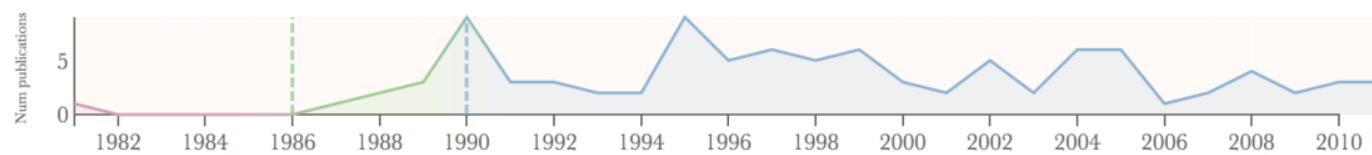
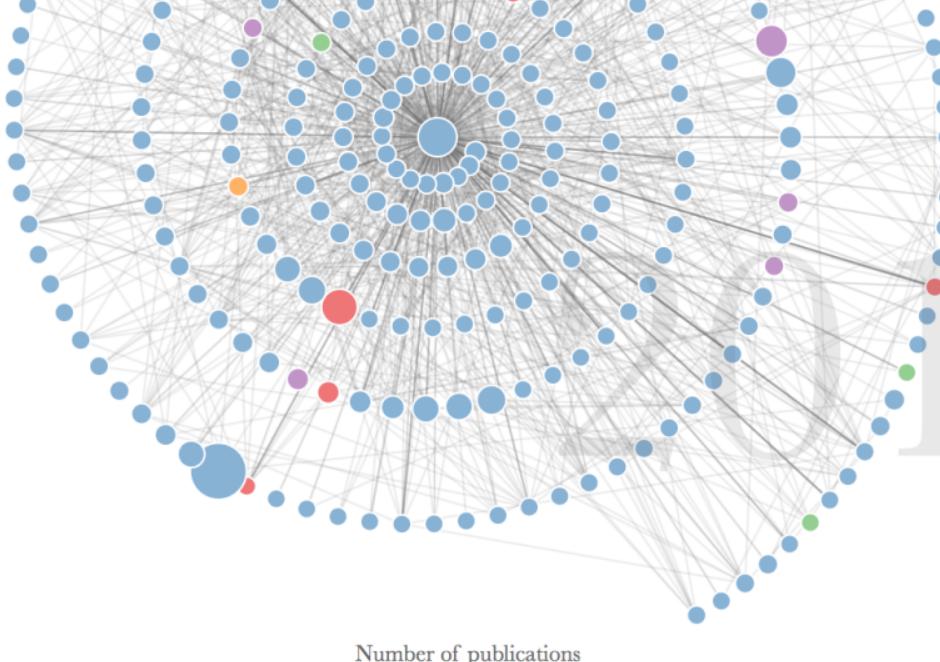
This data is available as a set of zipped text files stored in Microsoft Azure blob storage and available via HTTP. The file size (zipped) is ~30GB and may be downloaded [here](#).

The Challenge

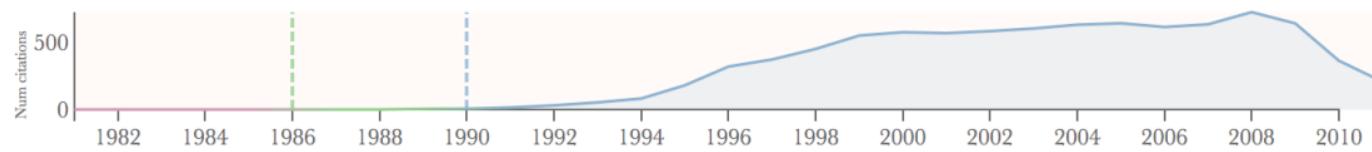
The goal of the Ranker Challenge is to assess the query-independent importance of scholarly articles, using data from the Microsoft Academic Graph.



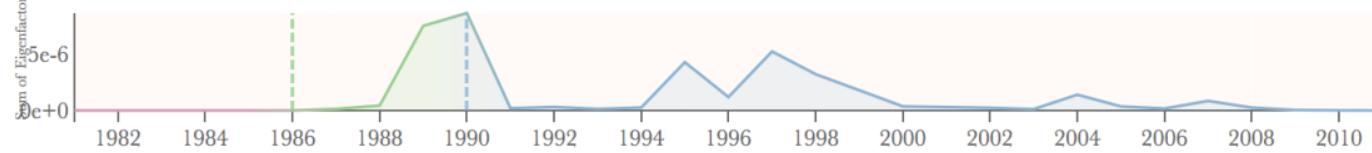
Jason Portenoy



Number of citations received



Sum of eigenfactor for this author's publications by year



scholar.eigenfactor.org

Visualizing Scholarly Influence Over Time

Influence of Pew Scholars

Roberta A. Gottlieb

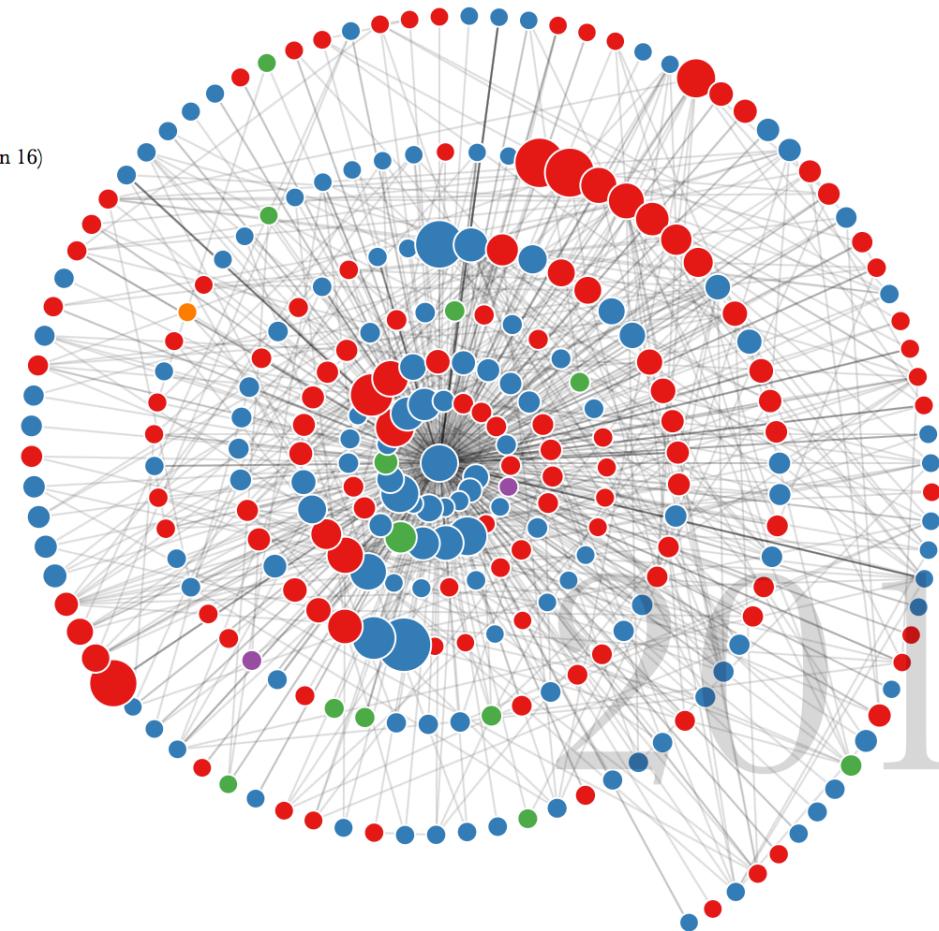
[Learn More](#)

- █ 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



Visualizing Scholarly Influence Over Time

Influence of Pew Scholars

Mark W. Grinstaff

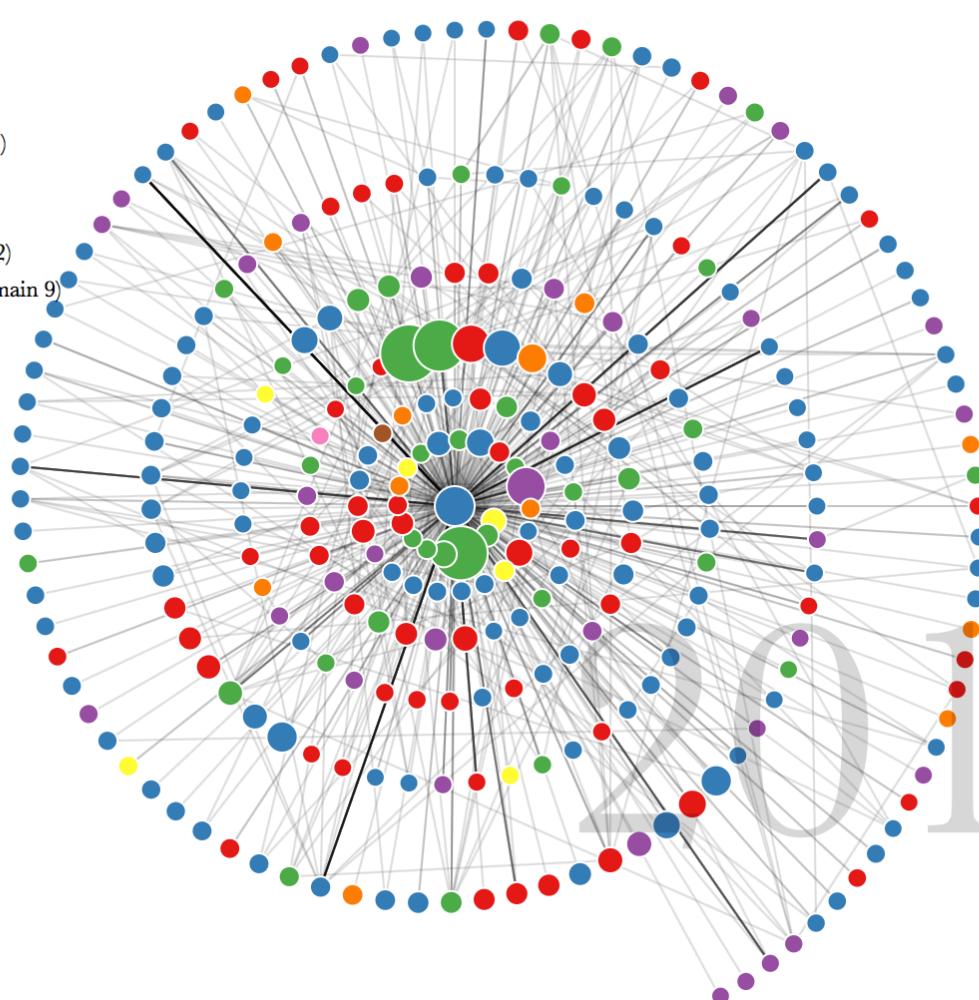
[Learn More](#)

- 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

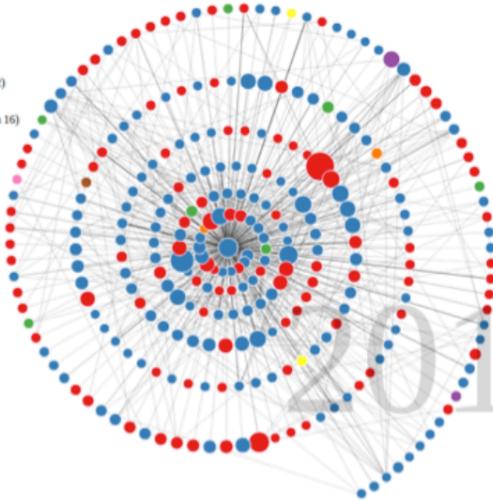


Pew Scholar
1999



Comparing Authors

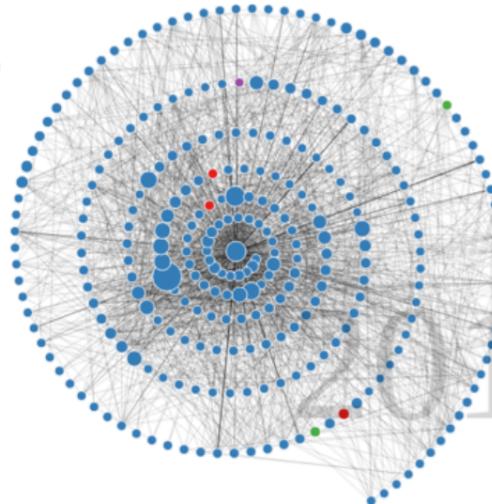
- Papers in category "Medicine" (domain 6)
- Papers in category "Biology" (domain 4)
- Papers in category "Chemistry" (domain 5)
- Papers in category "Engineering" (domain 8)
- Papers in category "Material Science" (domain 12)
- Papers in category "Physics" (domain 19)
- Papers in category "Agriculture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



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

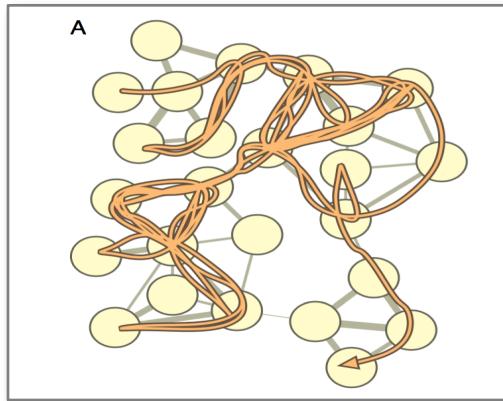


- Papers in category "Biology" (domain 4)
- Papers in category "Medicine" (domain 6)
- Papers in category "Chemistry" (domain 5)
- Papers in category "Social Science" (domain 22)

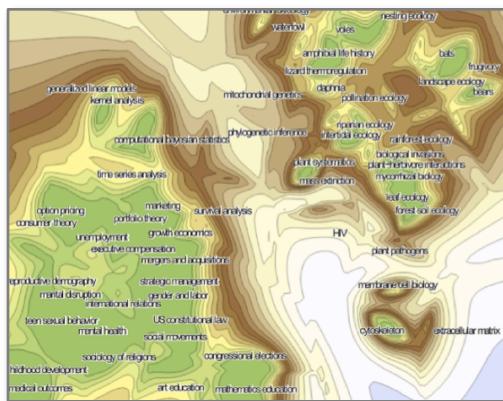


A more sparse network indicates fewer citations between papers shown in the network. This could be a result of the central scholar having impact across a wider set of academic communities.





Science of Mapping



Mapping of Science



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



Influential Articles



1960s 1970s 1980s 1990s 2000s 2010s 2020

- The Canonical Distribution of Commonness and ...
- An Equilibrium Theory of Insular Zoogeography
 - Turnover Rates in Insular Biogeography: ...
 - The Statistics and Biology of the ...

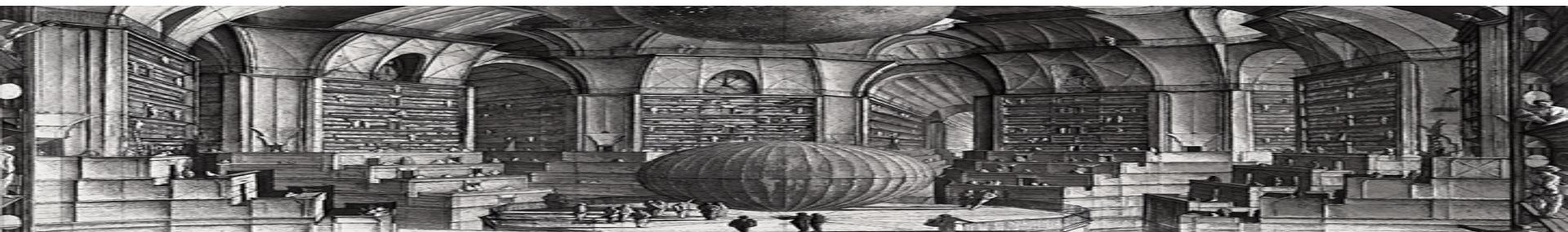
Related Topics



Species

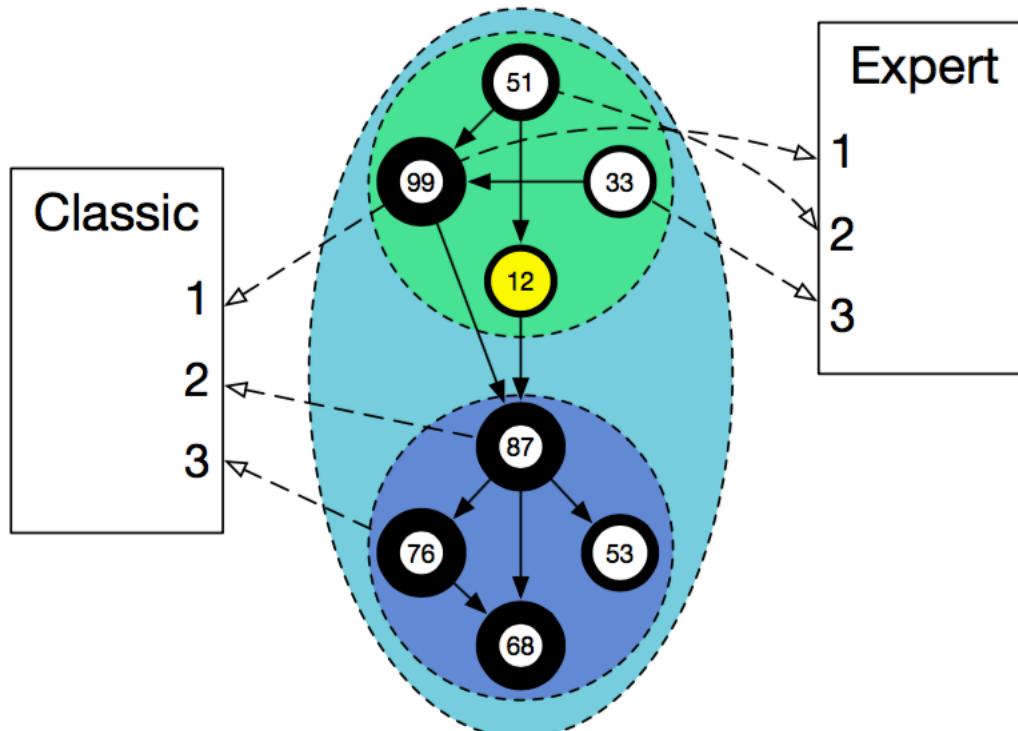
Habitat conservation

Explore the recommendations
babel.eigenfactor.org

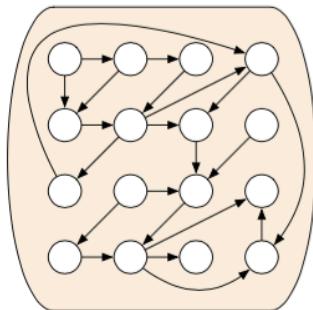


jevinw@uw.edu

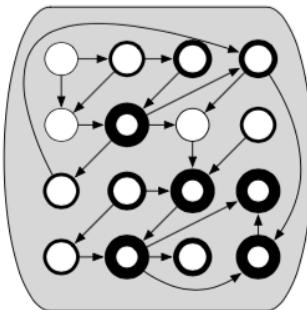
Recommend



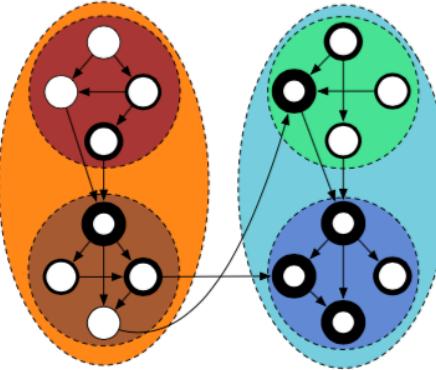
Assemble



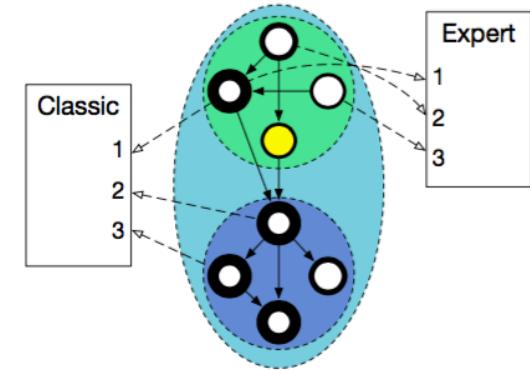
Rank



Cluster

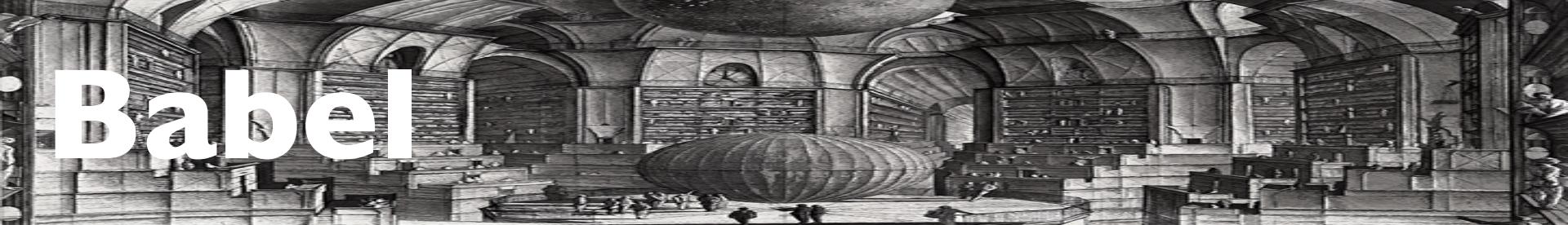


Recommend



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

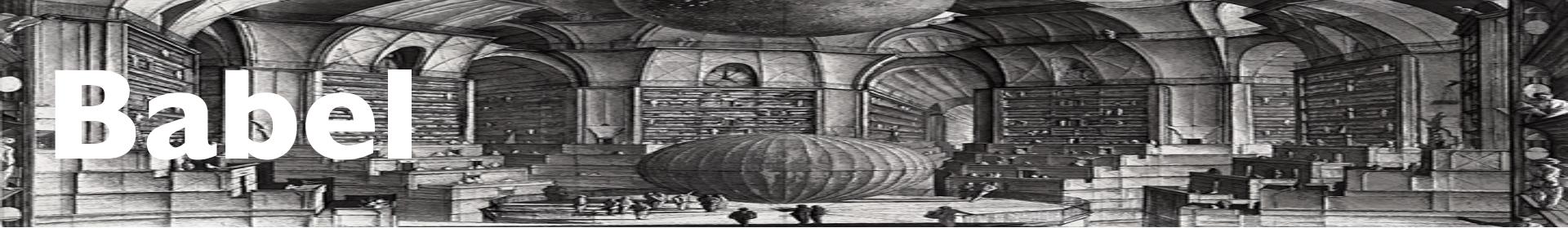
Babel



- 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

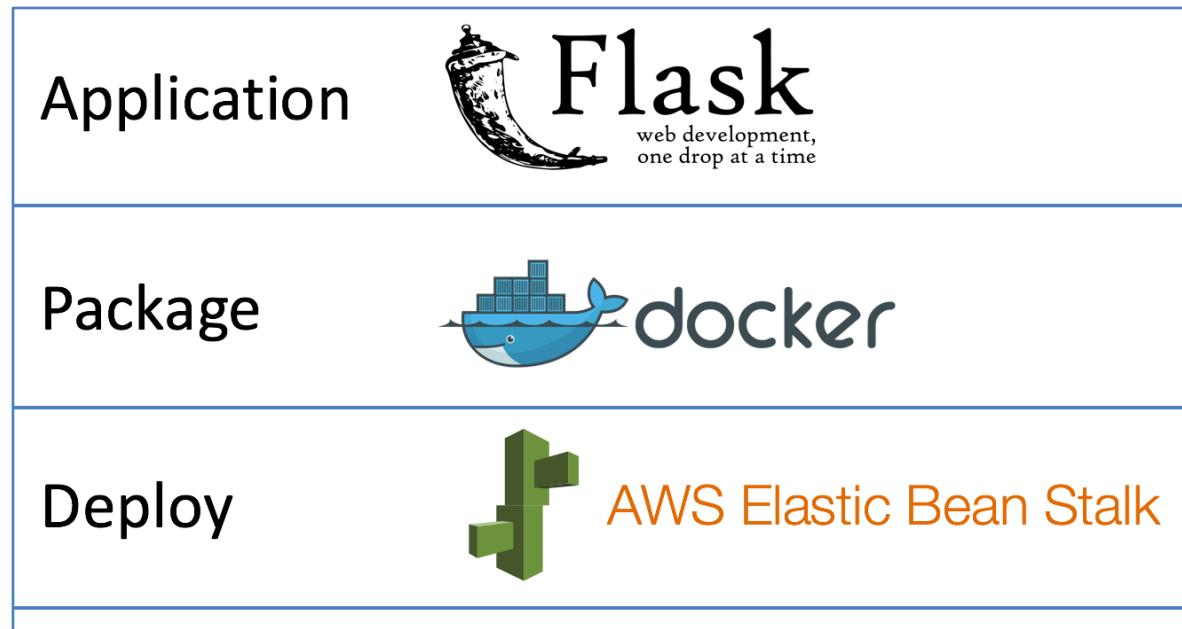
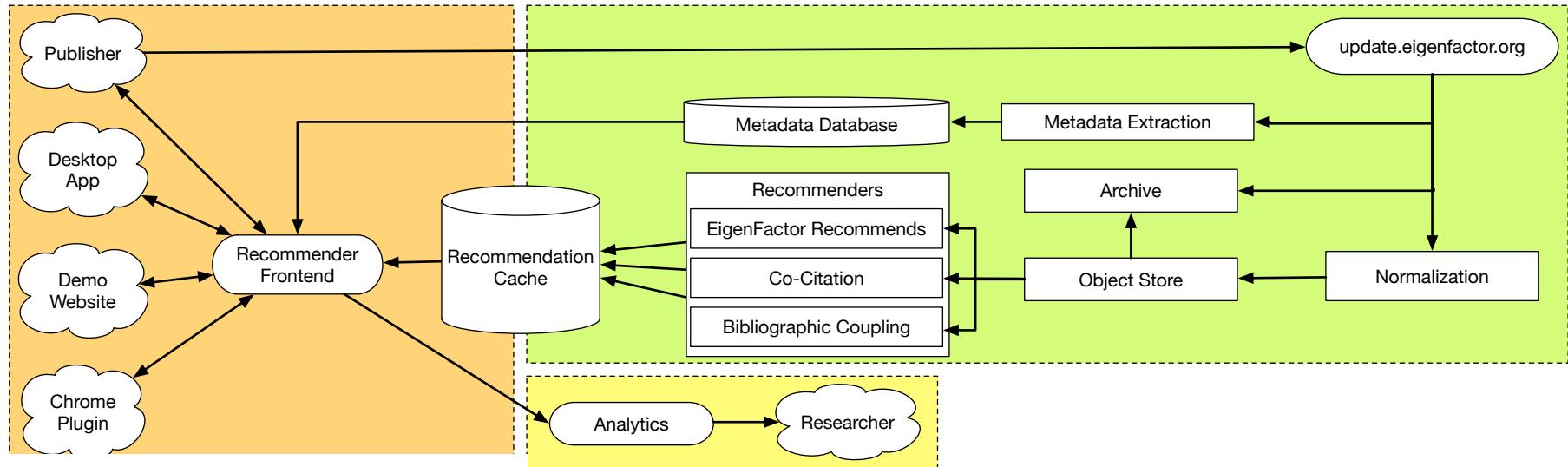
babel.eigenfactor.org

Babel



- Promote creation of better tools for scholarly authoring
- Allows 3rd parties to easily integrate with publishers, researchers
 - Examples: Zotero Plugin, Google Scholar Browser Plugin
- Easily test out new recommendations from other groups
 - Just change your end point!
- Publishers can host their own endpoints, meet their reliability requirements

Babel Architecture





Open API

File ▾ Preferences ▾ Generate Server ▾ Generate Client ▾ Help ▾ ✓ Processed with no error

```
268     ...
269     | '/recommendation/{publisher}/{paper_id}':
270     |   get:
271     |     description: Generates recommendations for `paper_id`
272     |     operationId: application.get_recommendation
273     |     parameters:
274     |       - $ref: '#/parameters/paper_id'
275     |       - $ref: '#/parameters/publisher'
276     |       - $ref: '#/parameters/client_id'
277     |       - default: 5
278     |       description: Maximum number of recommendations to return
279     |       in: query
280     |       maximum: 10
281     |       minimum: 1
282     |       name: limit
283     |       required: false
284     |       type: integer
285     |       - description: Algorithm to generate recommendations with. If not provided a
286     |         random algorithm will be used.
287     |       enum:
288     |         - ef_expert
289     |         - ef_classic
290     |       in: query
291     |       name: algorithm
292     |       required: false
293     |       type: string
294     |     produces:
295     |       - application/json
296     |     responses:
297     |       '200':
298     |         description: Successful response
299     |         schema:
300     |           $ref: '#/definitions/Recommendations'
301     /search:
302     get:
303     |       description: Searches metadata for known papers, authors or labels'
304     |       operationId: application.search
305     |       parameters:
306     |         - description: Search query
307     |           in: query
308     |           name: q
309     |           required: true
310     |           type: string
311     |         - collectionFormat: multi
312     |           description: Publishers to restrict search to.
313     |           in: query
314
```

/recommendation/{publisher}/{paper_id}

GET /recommendation/{publisher}/{paper_id}

Description

Generates recommendations for `paper_id`

Parameters

| Name | Located in | Description | Required | Schema |
|-----------|------------|--|----------|-----------|
| paper_id | path | Publisher assigned identifier of a paper | Yes | ⇒ string |
| publisher | path | Publisher to perform this operation on | Yes | ⇒ string |
| client_id | query | Identifier provided by the platform to clients to track client usage. | No | ⇒ string |
| limit | query | Maximum number of recommendations to return | No | ⇒ integer |
| algorithm | query | Algorithm to generate recommendations with. If not provided a random algorithm will be used. | No | ⇒ string |

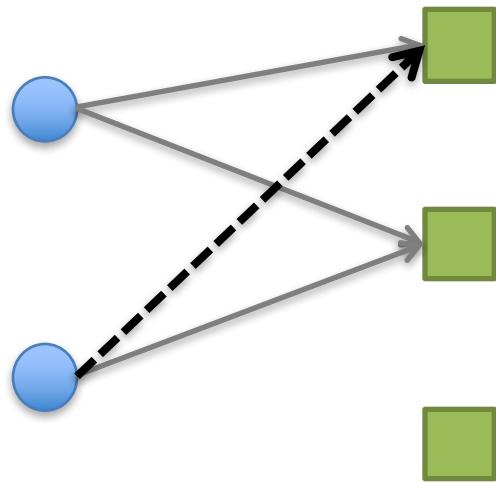
Responses

| Code | Description | Schema |
|------|---------------------|---|
| 200 | Successful response | ▼ Recommendations { ⇒ results: ⇒ [] ⇒ transaction_id: ⇒ string } |

Try this operation

A-B Testing

Usage-based

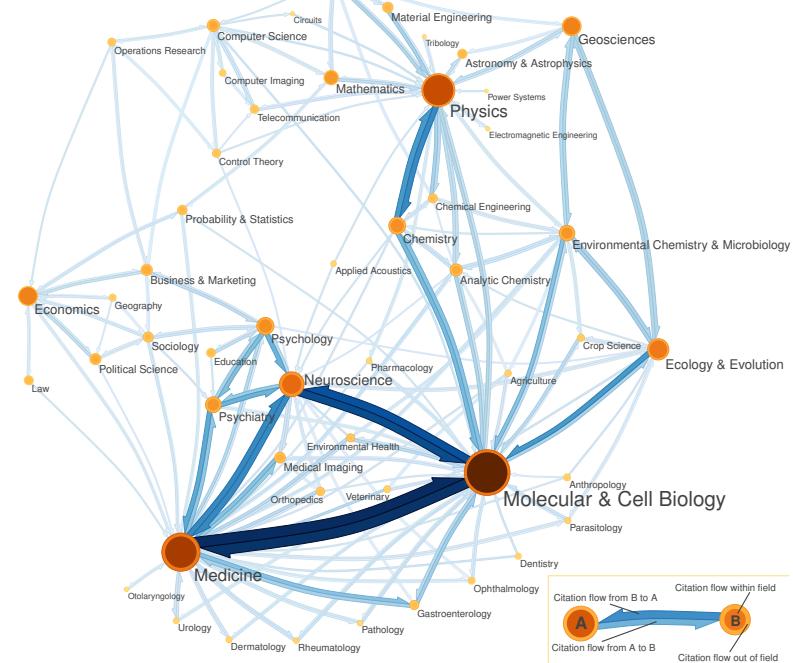


People

Papers

Utilizes download paper of similar readers

Citation-based



Utilizes hierarchical structure of citation graph and the relative position of papers

Browser Plugins



Eigenfactor - Google Scholar

https://scholar.google.com/scholar?q=Eigenfactor&btnG=&hl=en&as_sdt=0%2C48

Web Images More...

Google Eigenfactor

Scholar About 38,900 results (0.05 sec)

Articles **The Eigenfactor™ metrics** CT Bergstrom, JD West... - The Journal of ..., 2008 - Soc Neuroscience
Quantitative metrics are poor choices for assessing the research output of an individual scholar. Summing impact factors, counting citations, tallying an h-index, or looking at **Eigenfactor™ Scores** (described below)—none of these methods are adequate compared ...
Cited by 234 Related articles All 11 versions Cite Save

Case law My library

Any time Since 2015 Since 2014 Since 2011 Custom range...
P Yu, H Van de Sompel - Science, 1965 - eigenfactor.org
... Our aim at eigenfactor.org is develop ways of extracting this information. ... **Eigenfactor** algorithm modifies the basic eigenvector centrality algorithm to overcome these problems and to better handle certain peculiarities of journal citation data. ...
Cited by 2203 Related articles All 10 versions Cite Save More

Sort by relevance Sort by date
 include patents include citations
 Create alert

Eigenfactor: Does the principle of repeated improvement result in

Eigenfactor - Google Scholar

https://scholar.google.com/scholar?hl=en&q=Eigenfactor&btnG=&as_sdt=1%2C48&as_sdtp=

Web Images More...

Google Eigenfactor

Scholar About 38,900 results (0.02 sec)

Articles **The Eigenfactor™ metrics** CT Bergstrom, JD West... - The Journal of ..., 2008 - Soc Neuroscience
Quantitative metrics are poor choices for assessing the research output of an individual scholar. Summing Impact factors, counting citations, tallying an h-index, or looking at **Eigenfactor™ Scores** (described below)—none of these methods are adequate compared ...
Cited by 234 Related articles All 11 versions Cite Save

Case law My library

Any time Since 2015 Since 2014 Since 2011 Custom range...
P Yu, H Van de Sompel - Science, 1965 - eigenfactor.org
... Our aim at eigenfactor.org is develop ways of extracting this information. ... **Eigenfactor** algorithm modifies the basic eigenvector centrality algorithm to overcome these problems and to better handle certain peculiarities of journal citation data. ...
Cited by 2203 Related articles All 10 versions Cite Save More

Sort by relevance Sort by date
 include patents include citations
 Create alert

Eigenfactor: Does the principle of repeated improvement result in

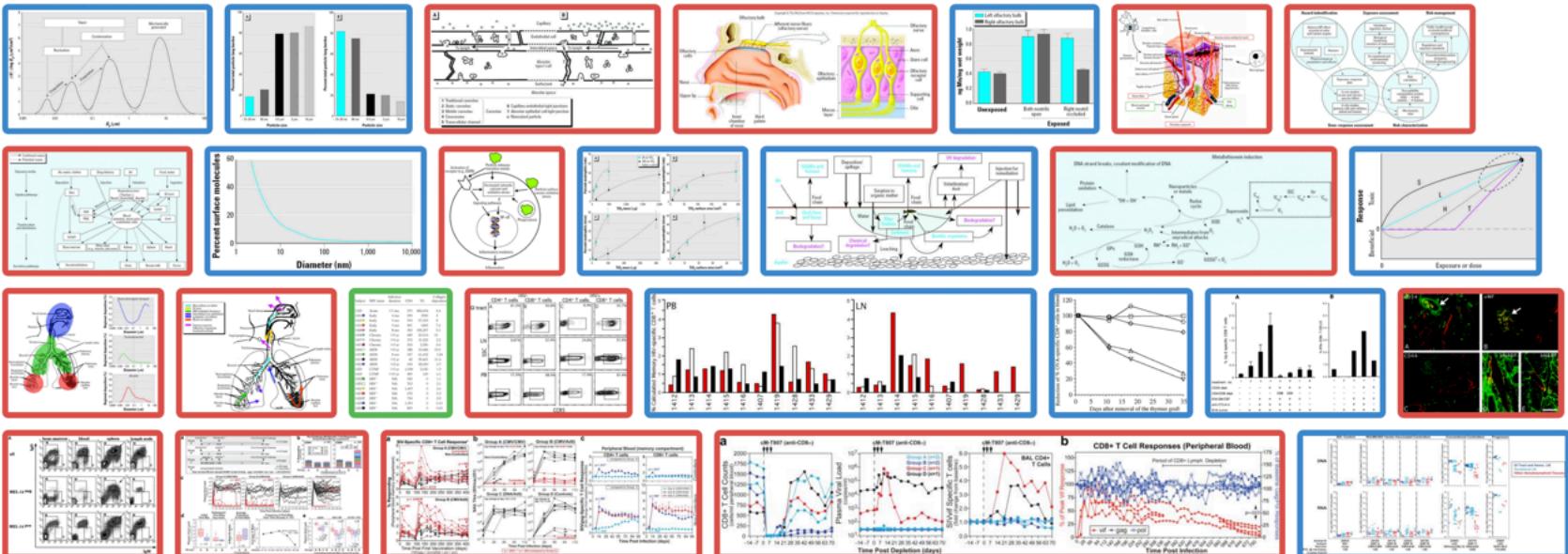
Figure-Centric Search Engine

 viziometrics.org

VizioMetrix About **Search** Crowdsourcing

Impact blood lymph

Composite Equation Diagram Photo Plot Table



A project of the eScience Institute at the University of Washington

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?

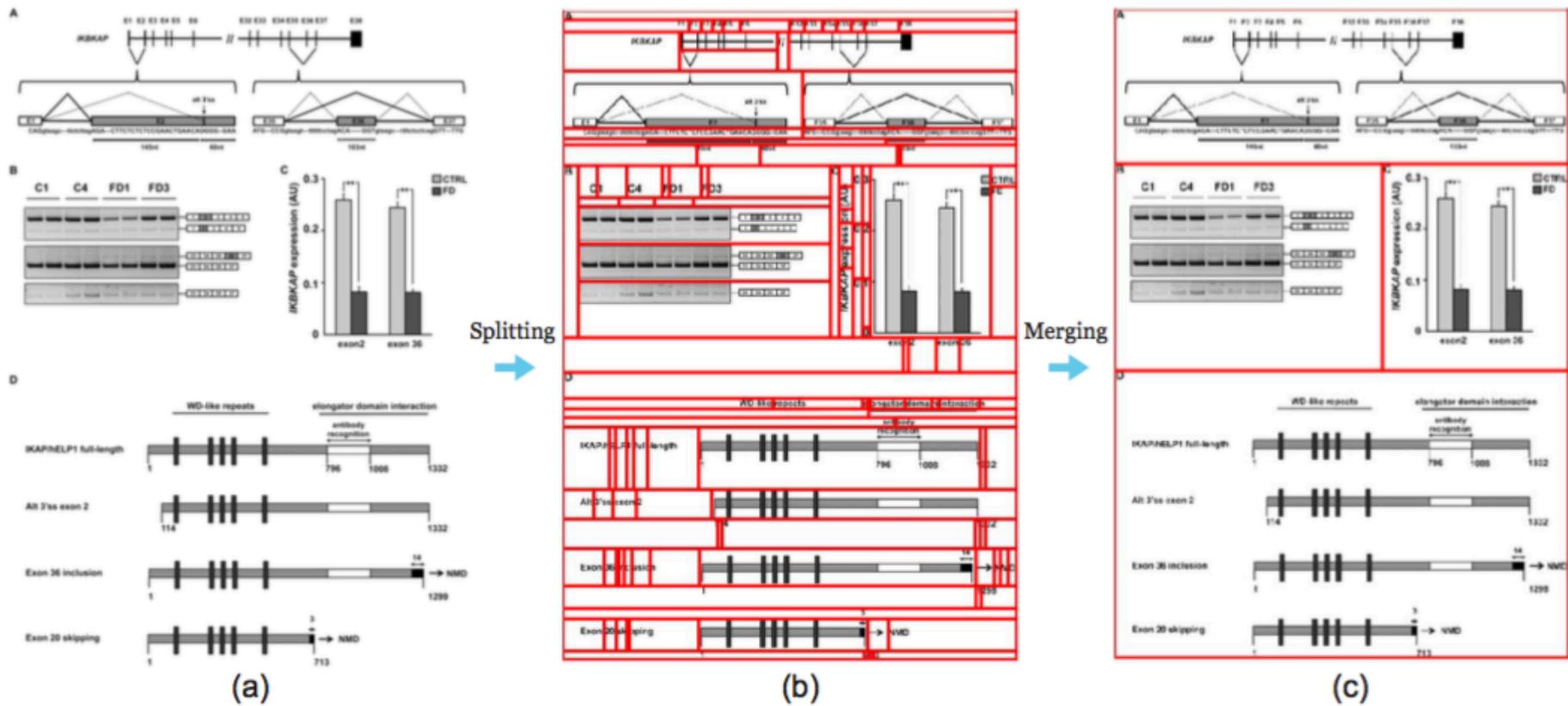
*How can we better utilize visual information
in the search and navigation process?*



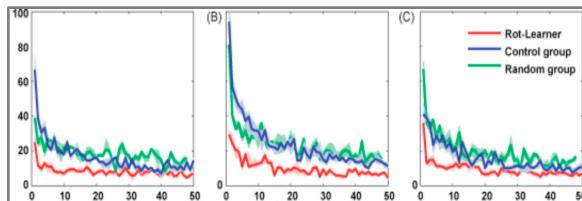
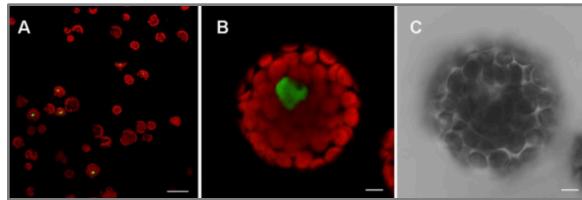
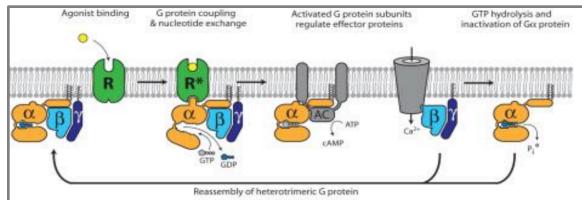
650,000 papers

5 million figures

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)

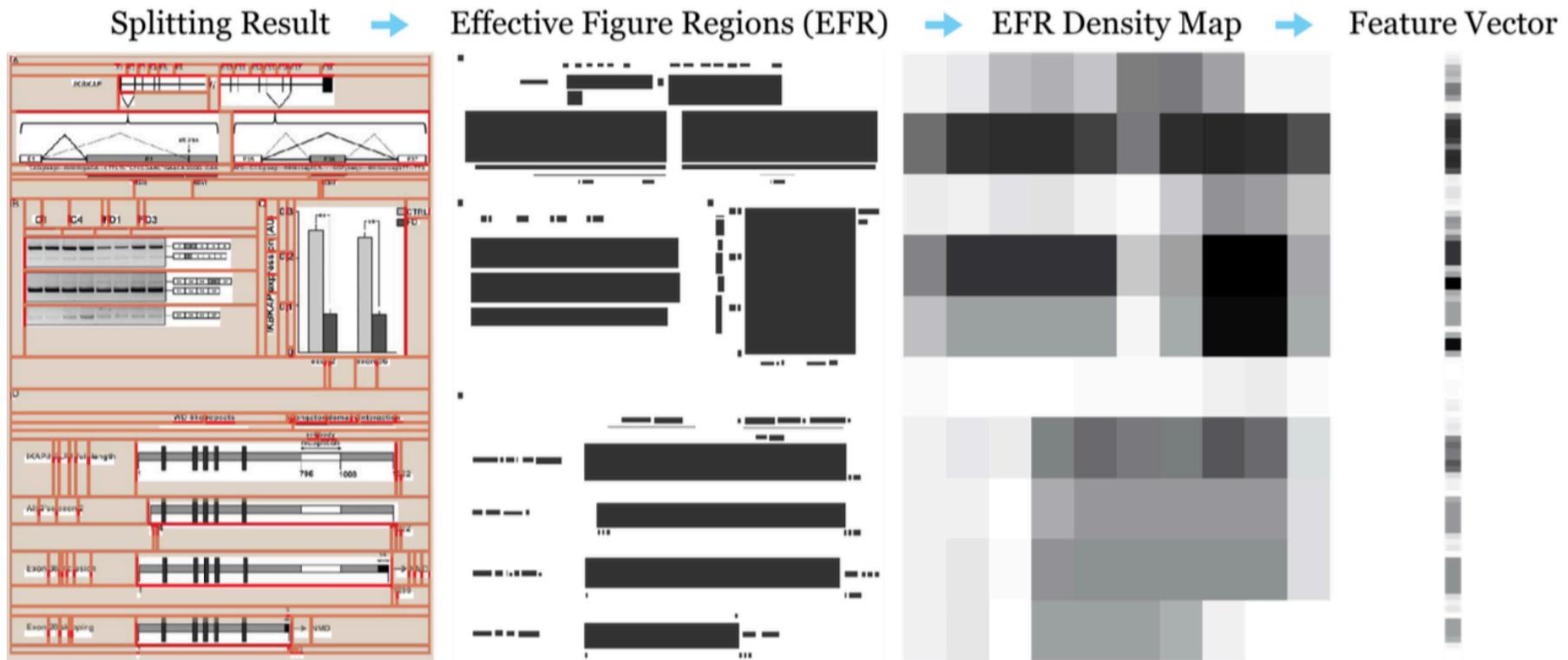
Schematics (769)

Photos (782)

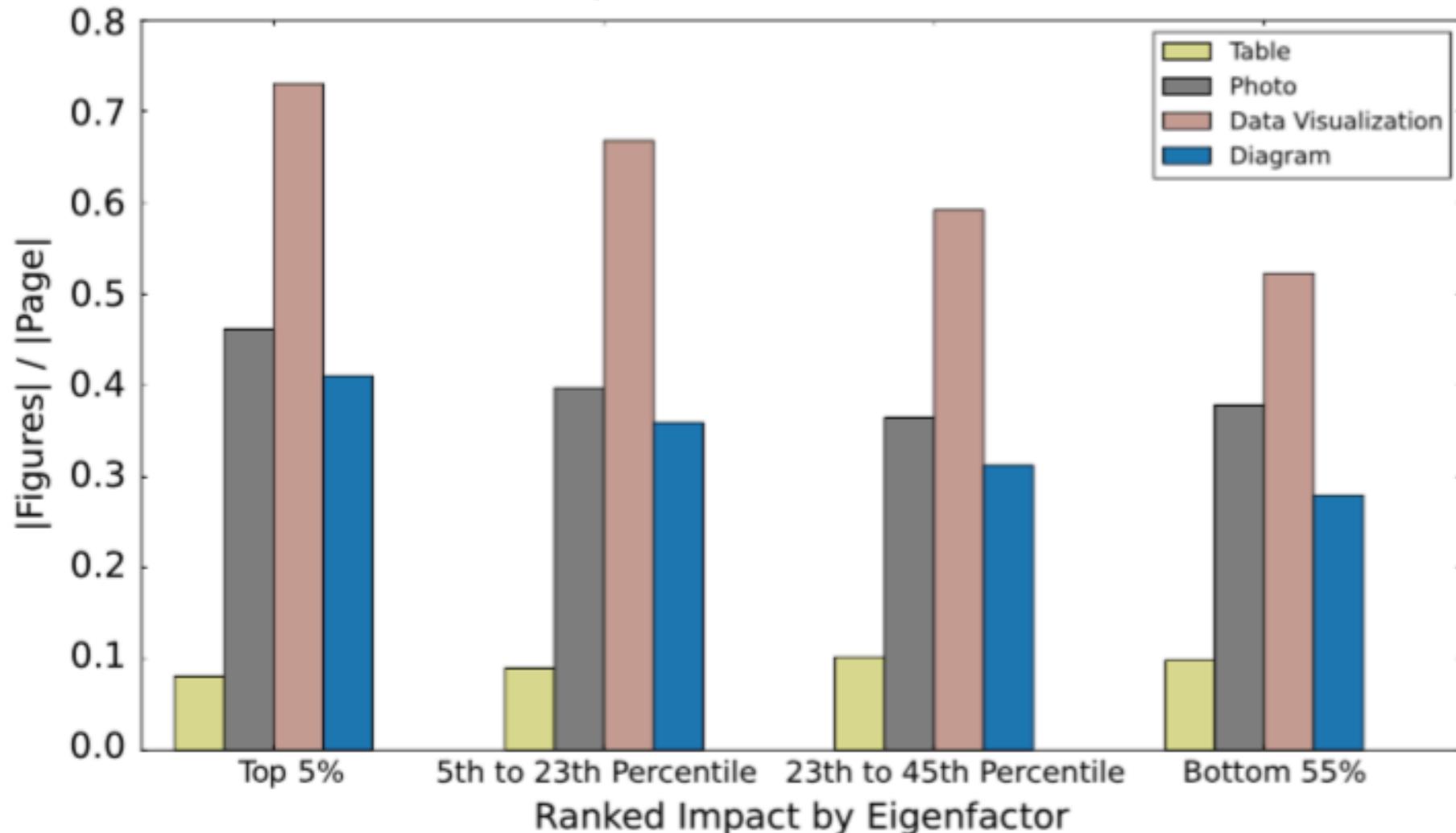
Plots (890)

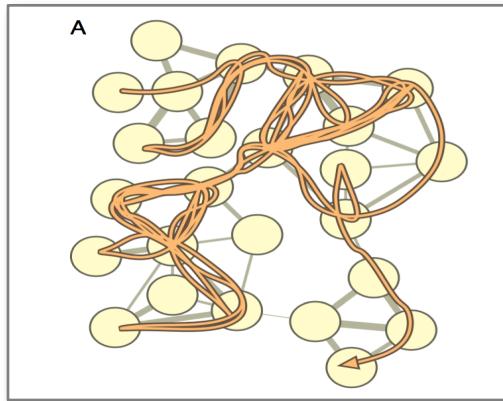
Tables (436)

Feature Extraction

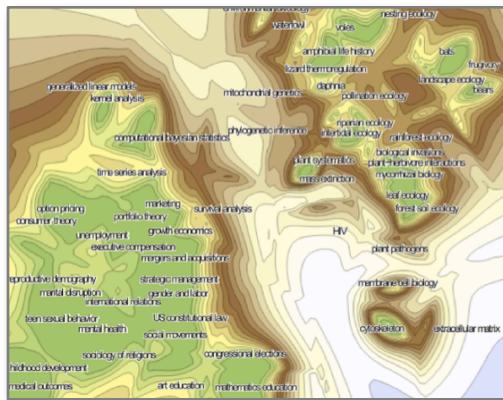


Impact versus Figure Density



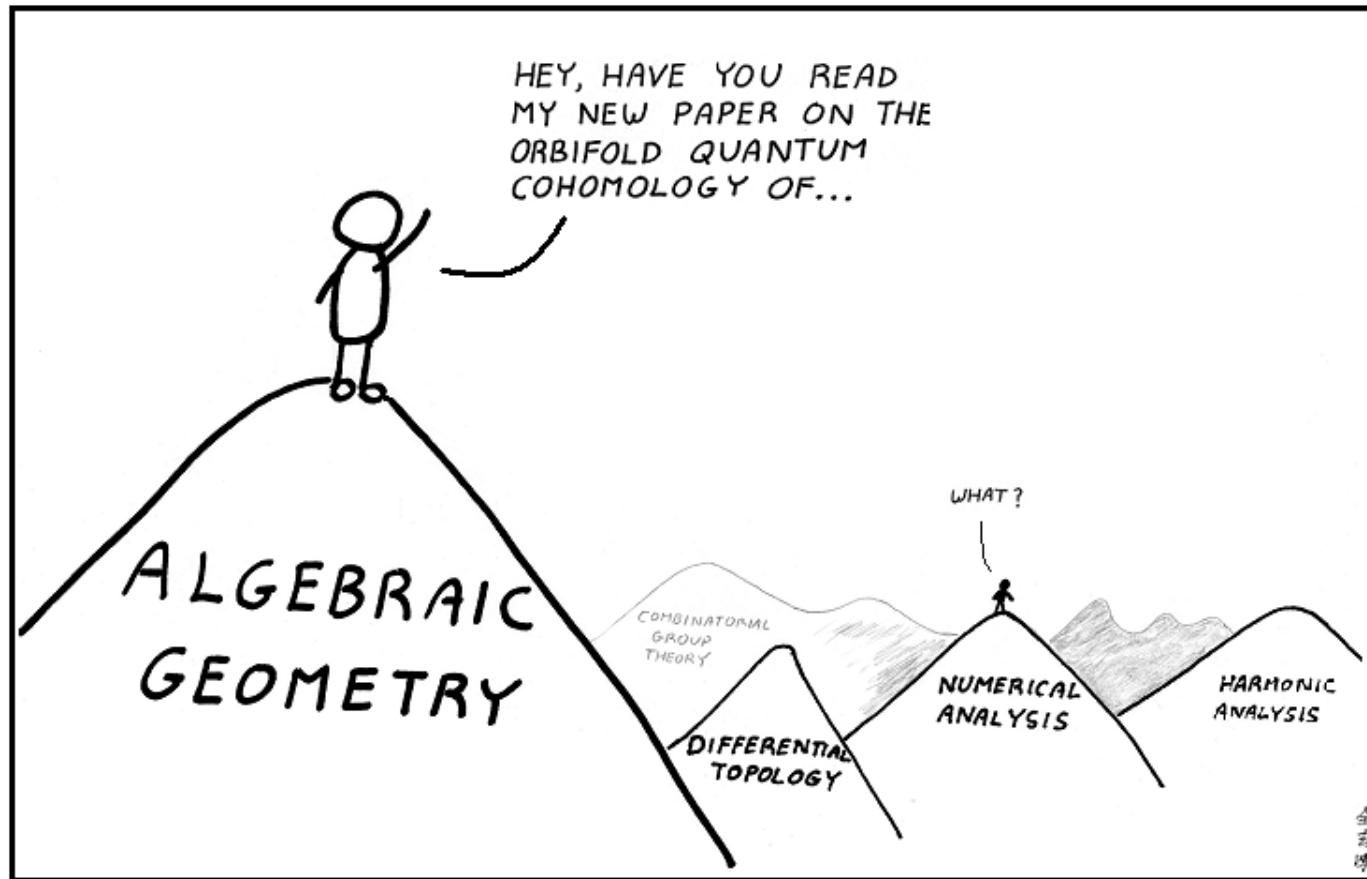


Science of Mapping

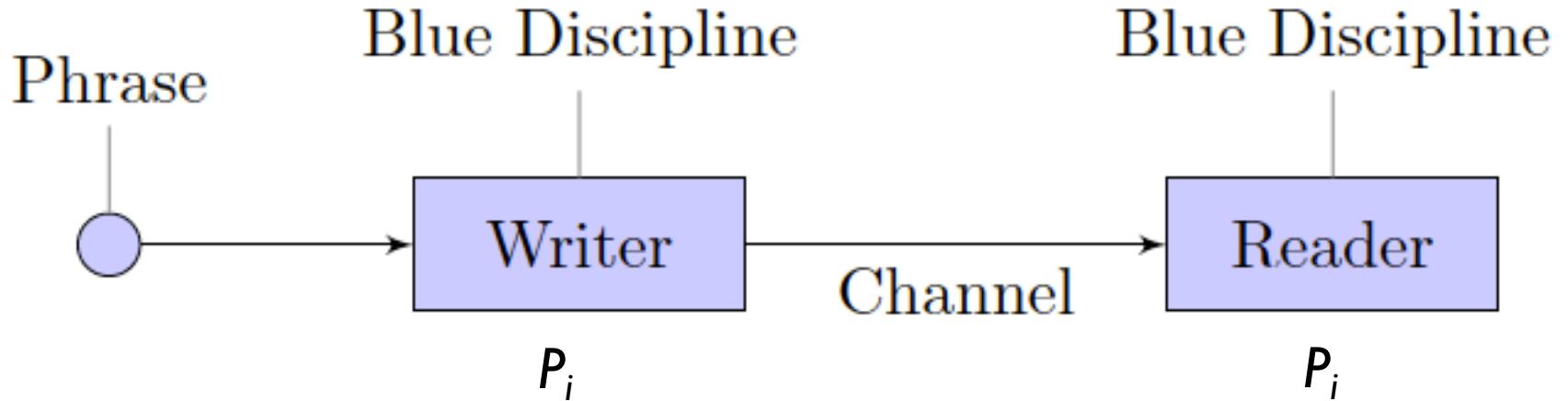


Mapping of Science

The jargon barriers of science



The Landscape of Modern Mathematics



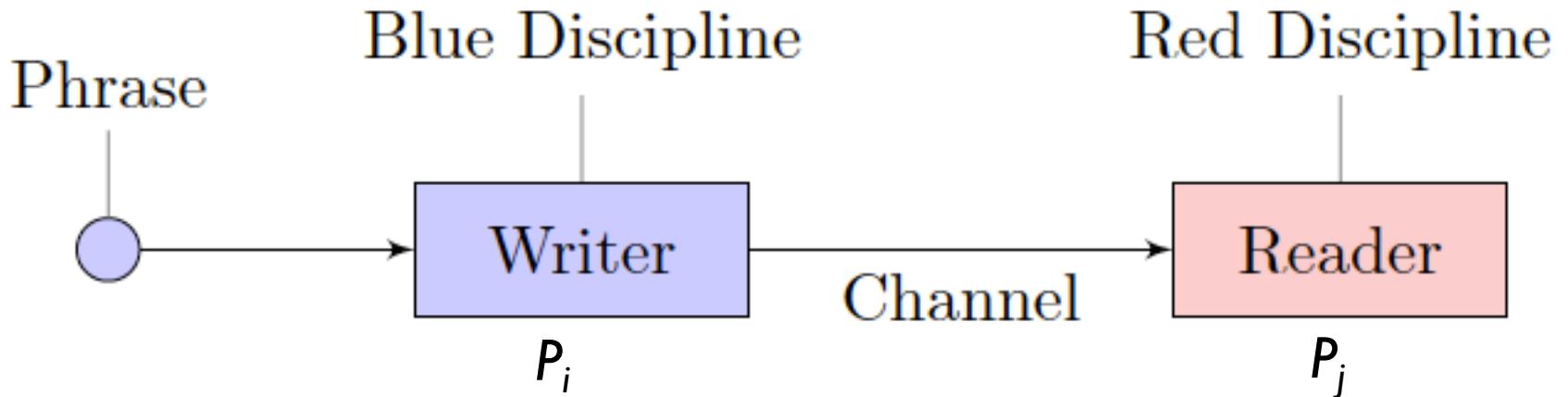
$X \sim$ space of all phrases

$P_i \sim$ 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

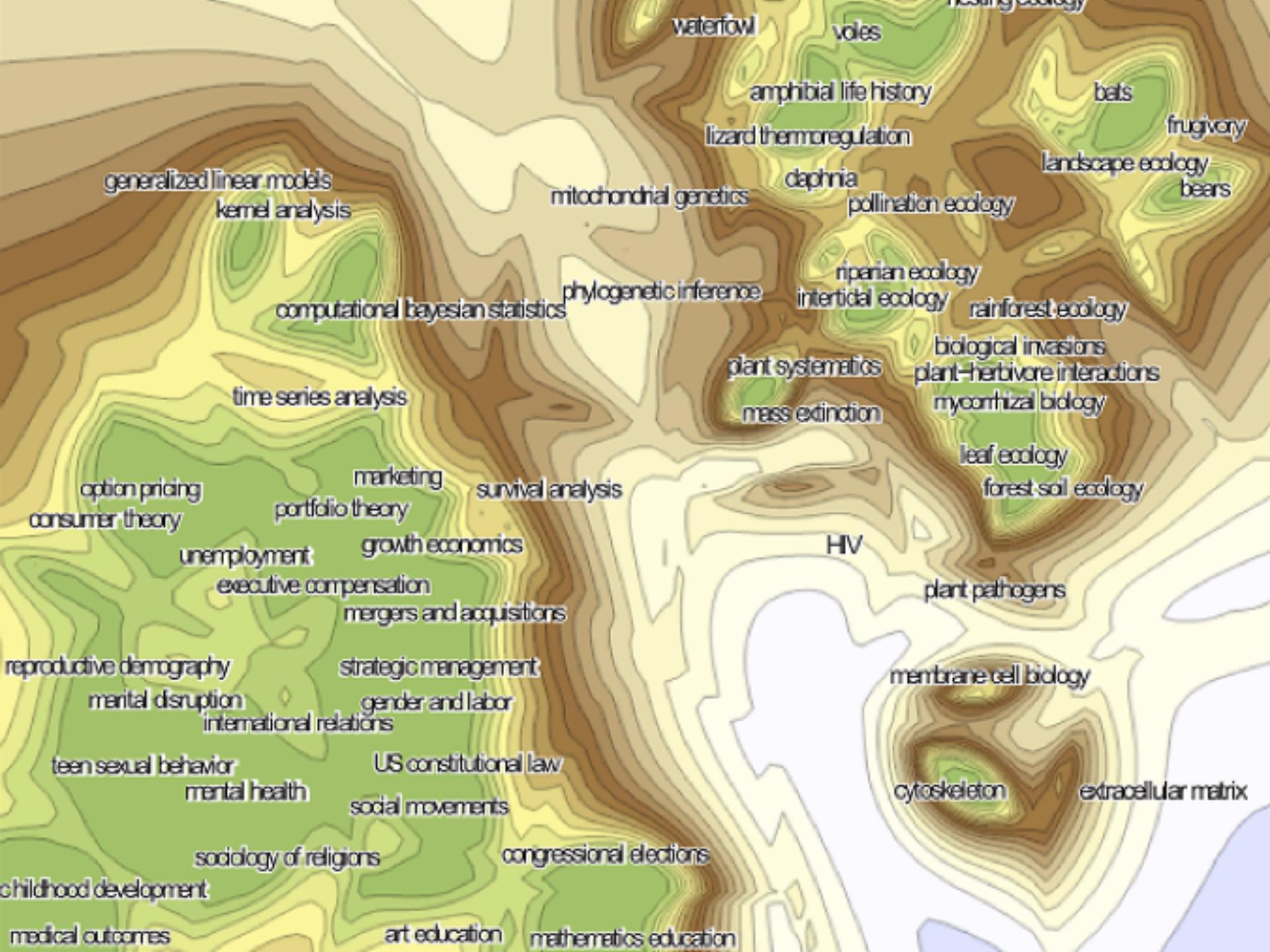
expected message length: $Q(p_i||p_j) = - \sum_{x \in \mathcal{X}} p_i(x) \log_2 p_j(x)$

efficiency of communication

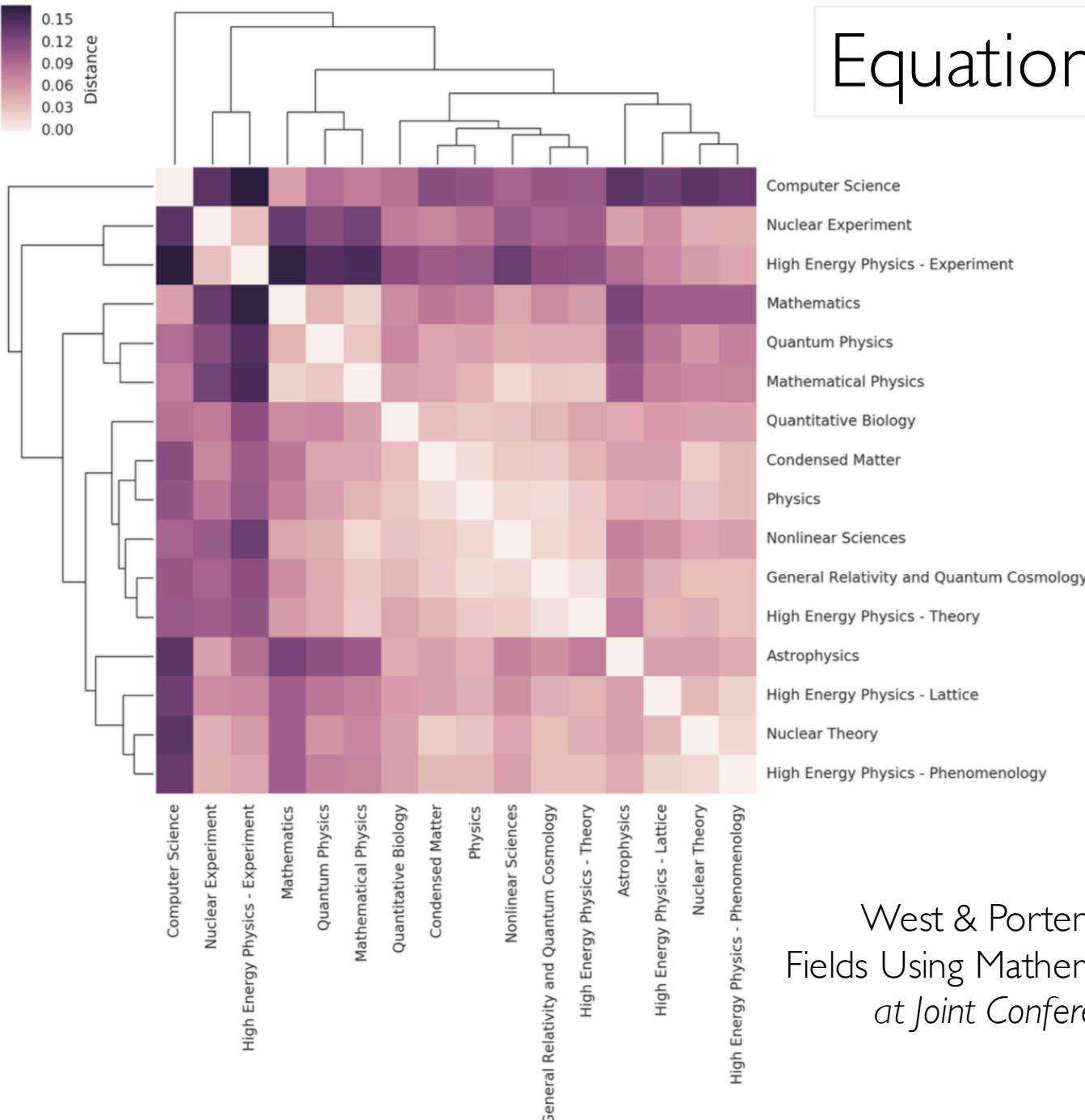
$$E_{ij} = \frac{H(X_i)}{Q(p_i||p_j)} = \frac{-\sum_{x \in \mathcal{X}} p_i(x) \log_2 p_i(x)}{-\sum_{x \in \mathcal{X}} p_i(x) \log_2 p_j(x)}$$

cultural hole

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



Equation Distance



West & Portenoy (2016) Delineating
Fields Using Mathematical Jargon. *BIRNDL*
at Joint Conference on Digital Libraries

Summary

- Study the *Science of Science*
- Assemble knowledge graph into machine readable formats
- Ask questions about the origin and evolution of ideas and fields, interdisciplinarity, impact assessment and sociology of science
- Building statistical and visualization tools that improve navigation, make relevant connections and facilitate knowledge discovery
- 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

Jevin West

jevinw@uw.edu

jevinwest.org

Eigenfactor.org

