

Analysis and Visualization of Science



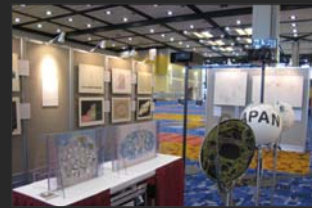
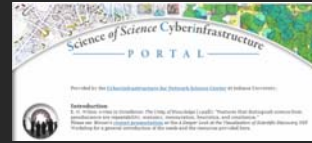
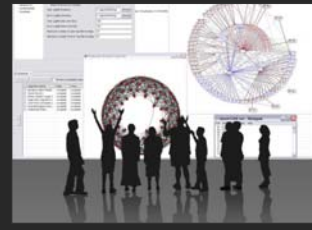
Angela Zoss, M.S.

Research Assistant, Cyberinfrastructure for Network Science Center
Doctoral Student, School of Library and Information Science

Indiana University, Bloomington, IN

amzoss@indiana.edu

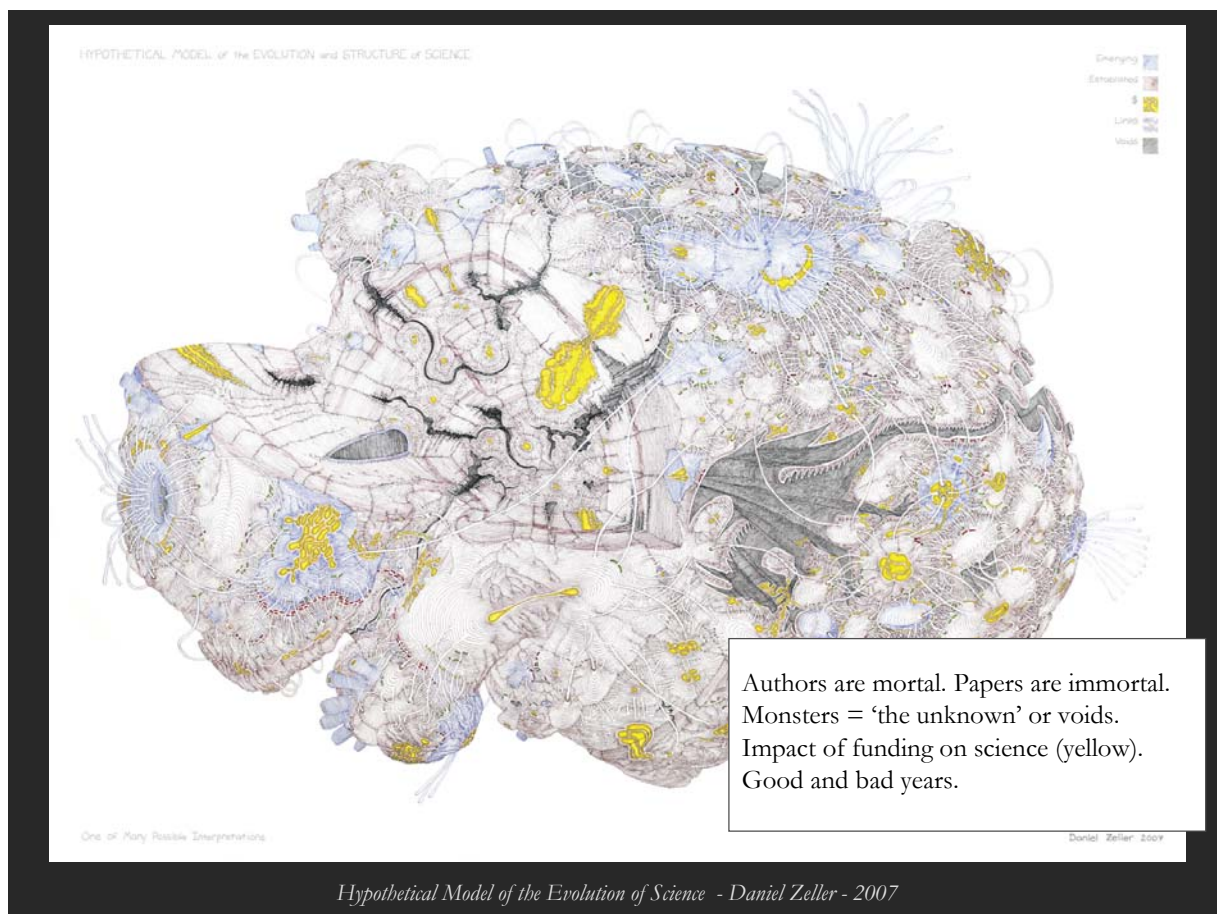
*February 8, 2011 – Workshop on Scholarly Communication and Informetrics
iConference 2011, Seattle, Washington*



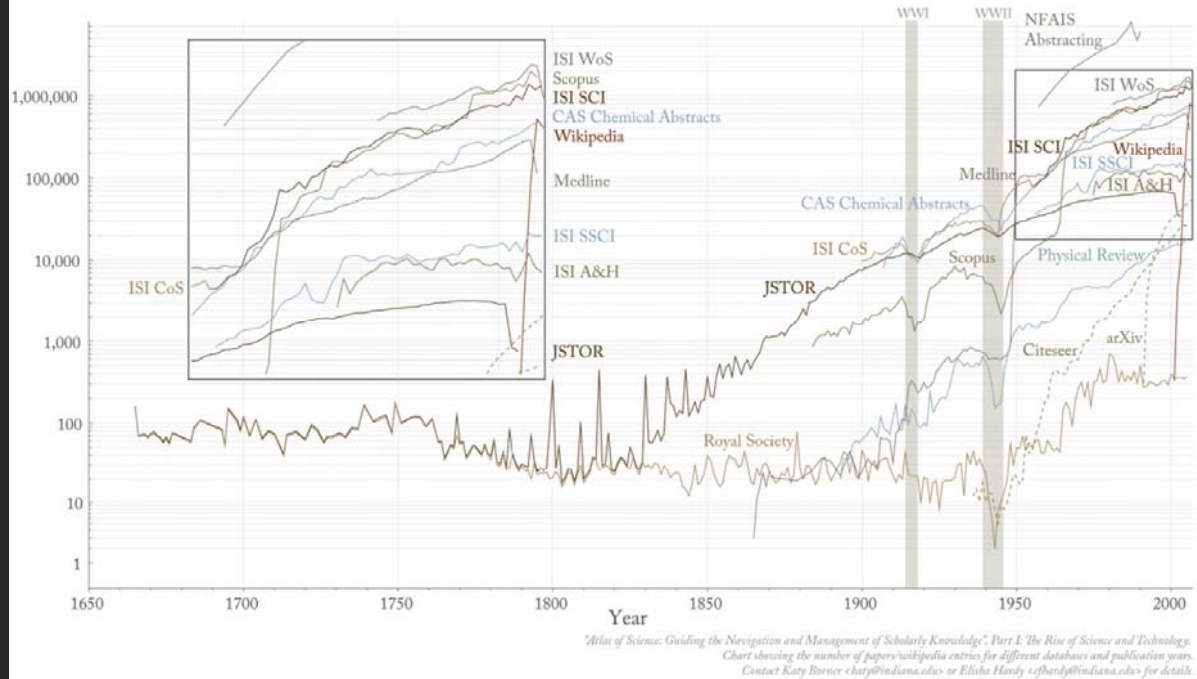
Analysis and Visualization of Science

- **What** is science?
- **Why** do we analyze and visualize science?
- **How** do we analyze and visualize science?

Conceptualizing Science

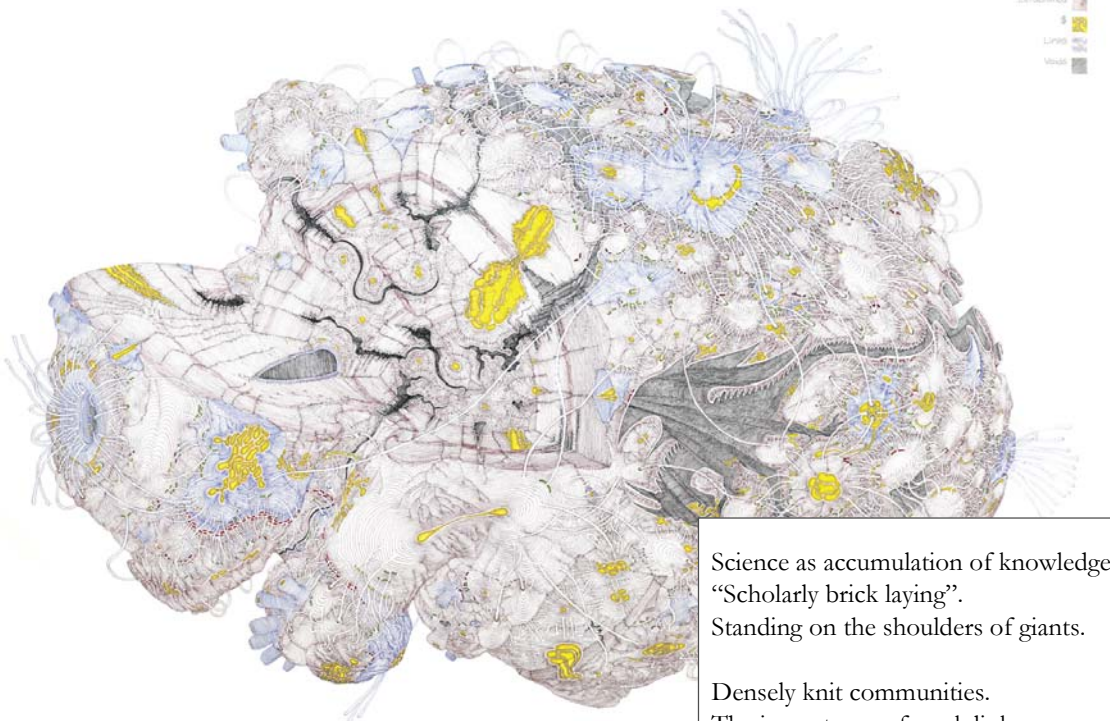


Papers & Wikipedia Entries



Atlas of Science - Katy Borner - 2010

HYPOTHETICAL MODEL of the EVOLUTION and STRUCTURE of SCIENCE



One of Many Possible Interpretations

Hypothetical Model of the Evolution of Science - Daniel Zeller - 2007

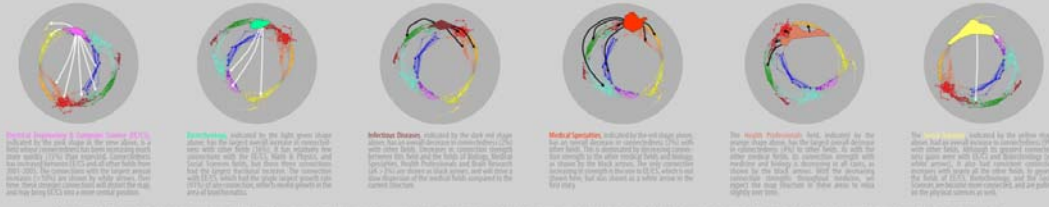
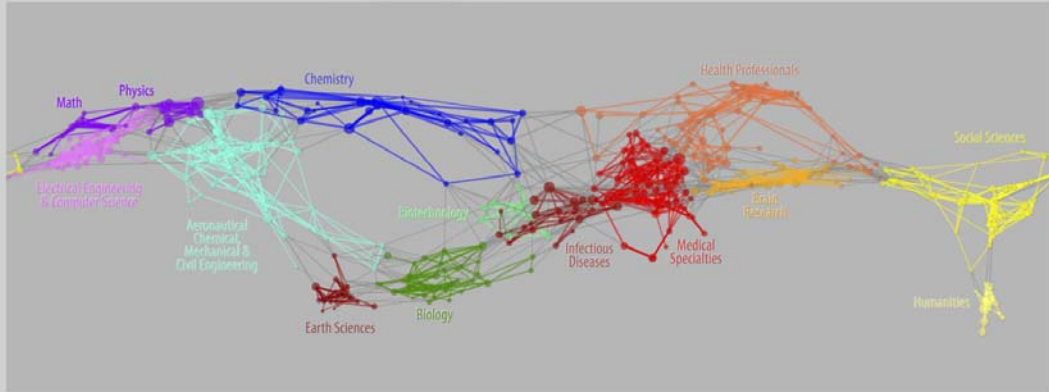
The shape of science was constructed by sorting over 7.2 million scholarly documents, proceedings, and abstracts into clusters. Each cluster represents a common theme. In this visualization, each cluster is represented by a different color and shape. The shape of each cluster is determined by the position of each document on the surface of a sphere based on the keywords it contains. The most frequent keywords in each cluster are shown as a different size of the sphere.

MAPS OF SCIENCE

Forecasting Large Trends in Science

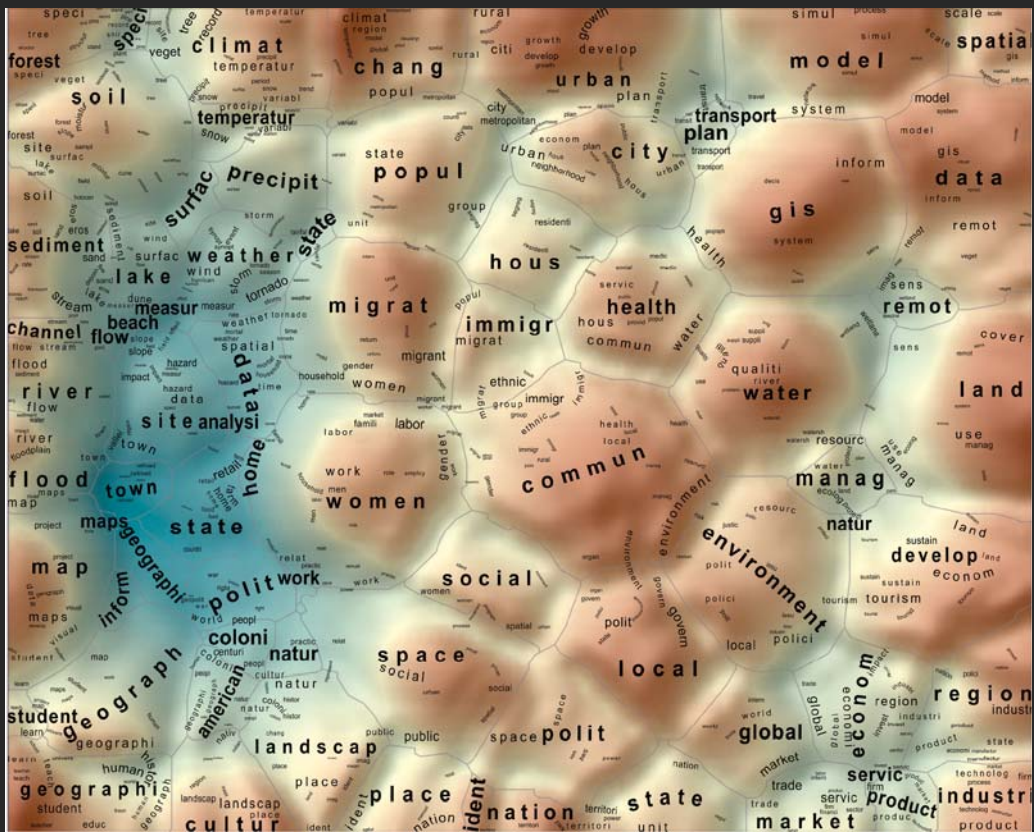
A visualization of 7.2 million scholarly documents
appearing in over 16,000 journals, proceedings or symposia
between Jan, 2001 and Dec, 2005

Calculations were performed using the large network visualization software, Gephi, to generate a map of the relationships between clusters. The structure of science shows below is showing trends in science, and connections between clusters. The structure of science shows below is showing trends in science, and connections between clusters. The structure of science shows below is showing trends in science, and connections between clusters.

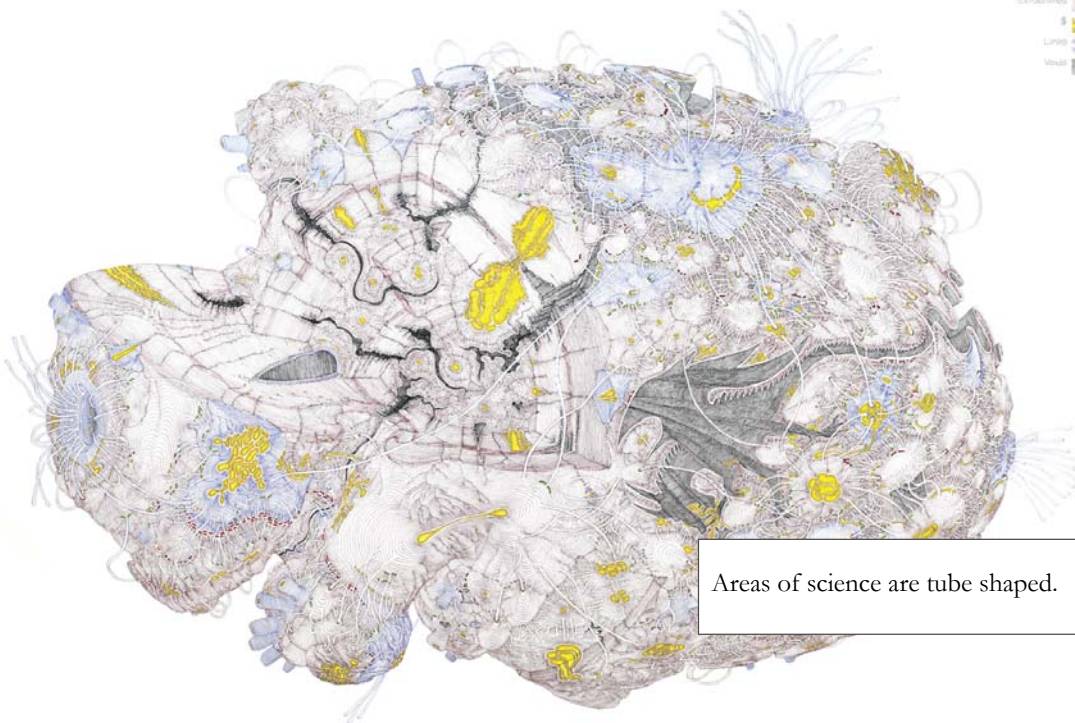


Source: University of California, San Diego. Visualization: Richard Klavans, Kevin Boyack. The visualization is based on the work of Klavans and Boyack. The visualization is based on the work of Klavans and Boyack. The visualization is based on the work of Klavans and Boyack.

Maps of Science: Forecasting Large Trends in Science - Richard Klavans, Kevin Boyack - 2007

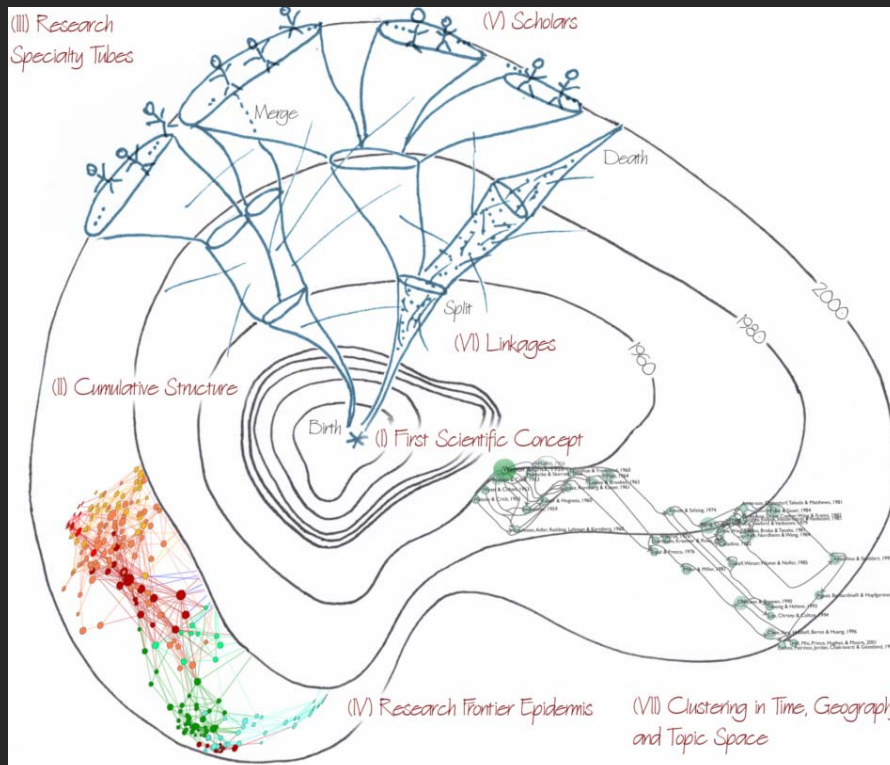


In Terms of Geography - Andre Skupin - 2005

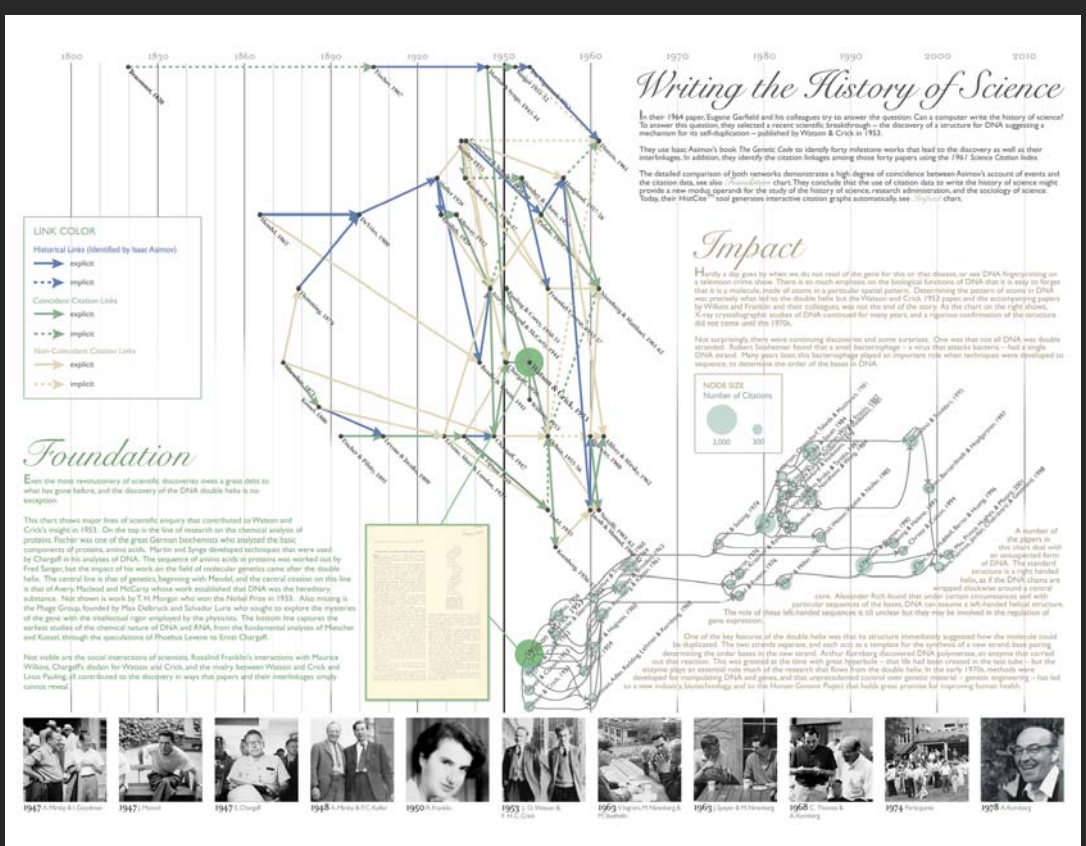


Areas of science are tube shaped.

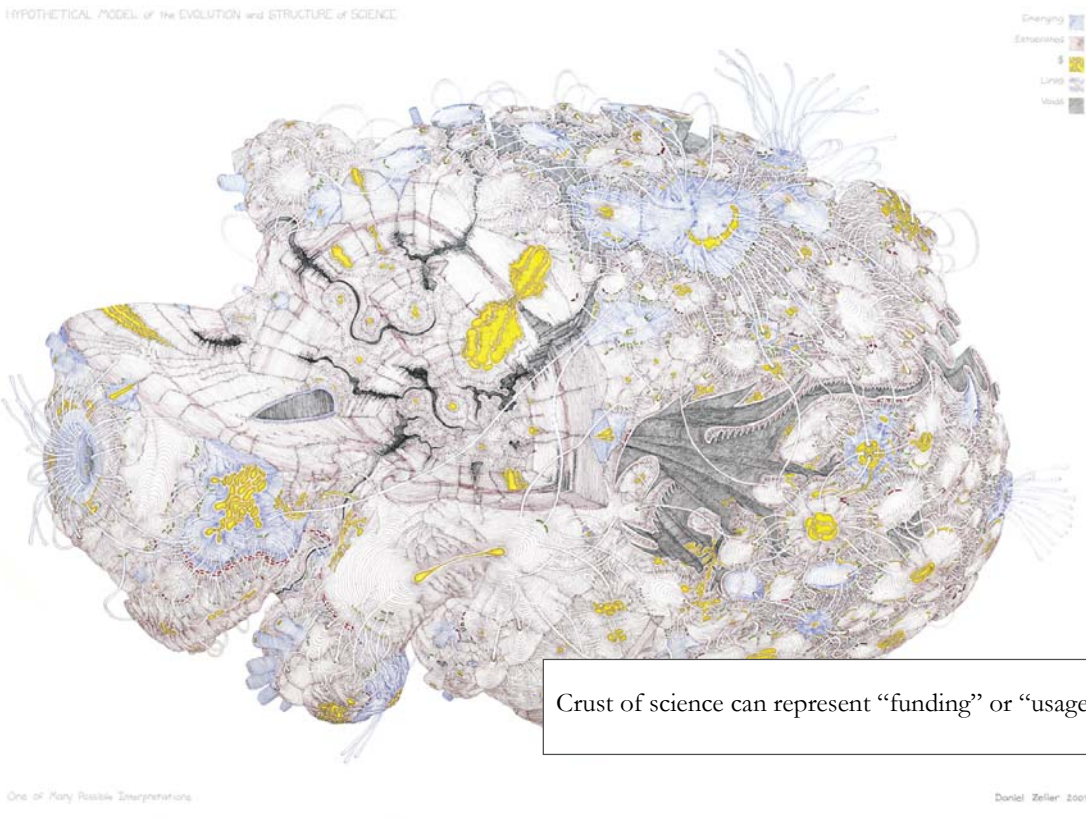
Hypothetical Model of the Evolution of Science - Daniel Zeller - 2007



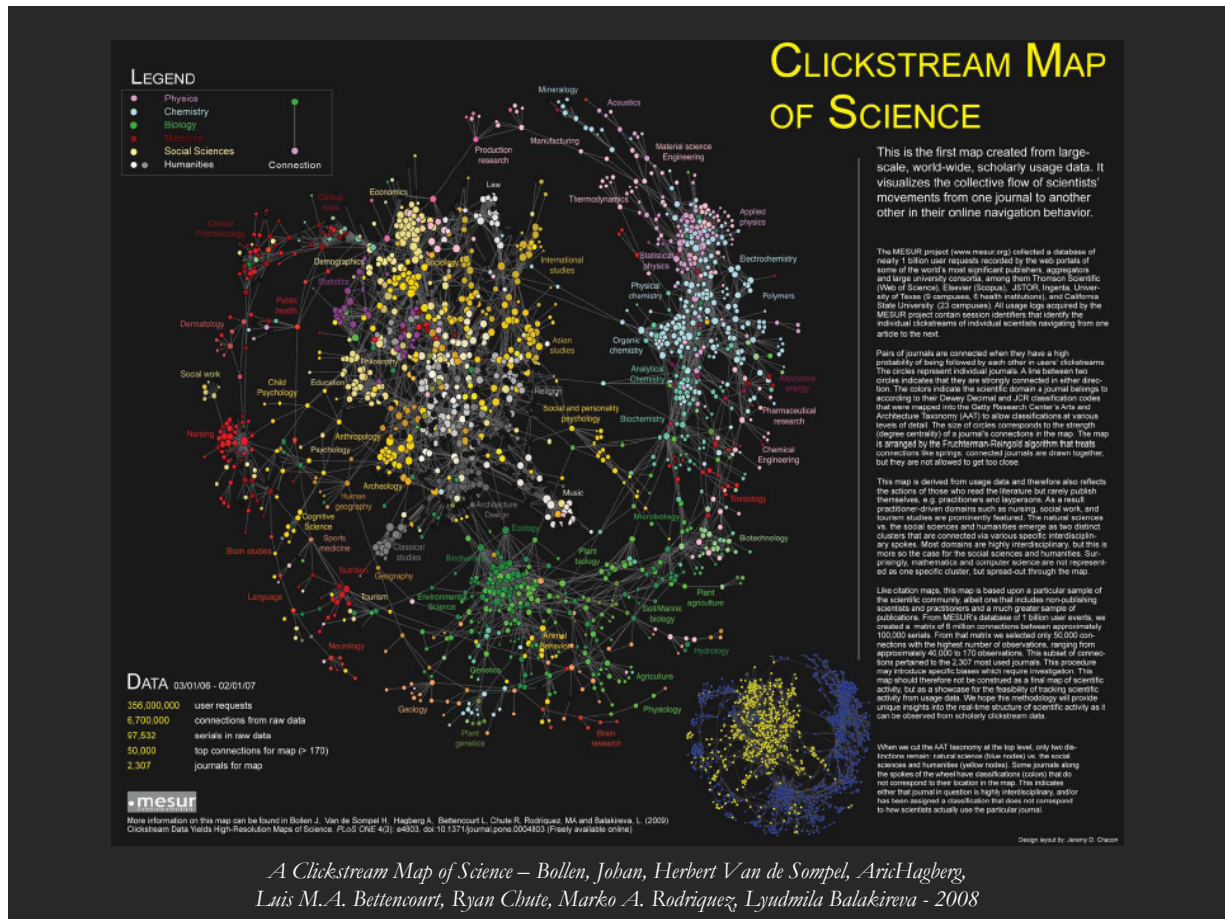
Atlas of Science - Katy Borner - 2010



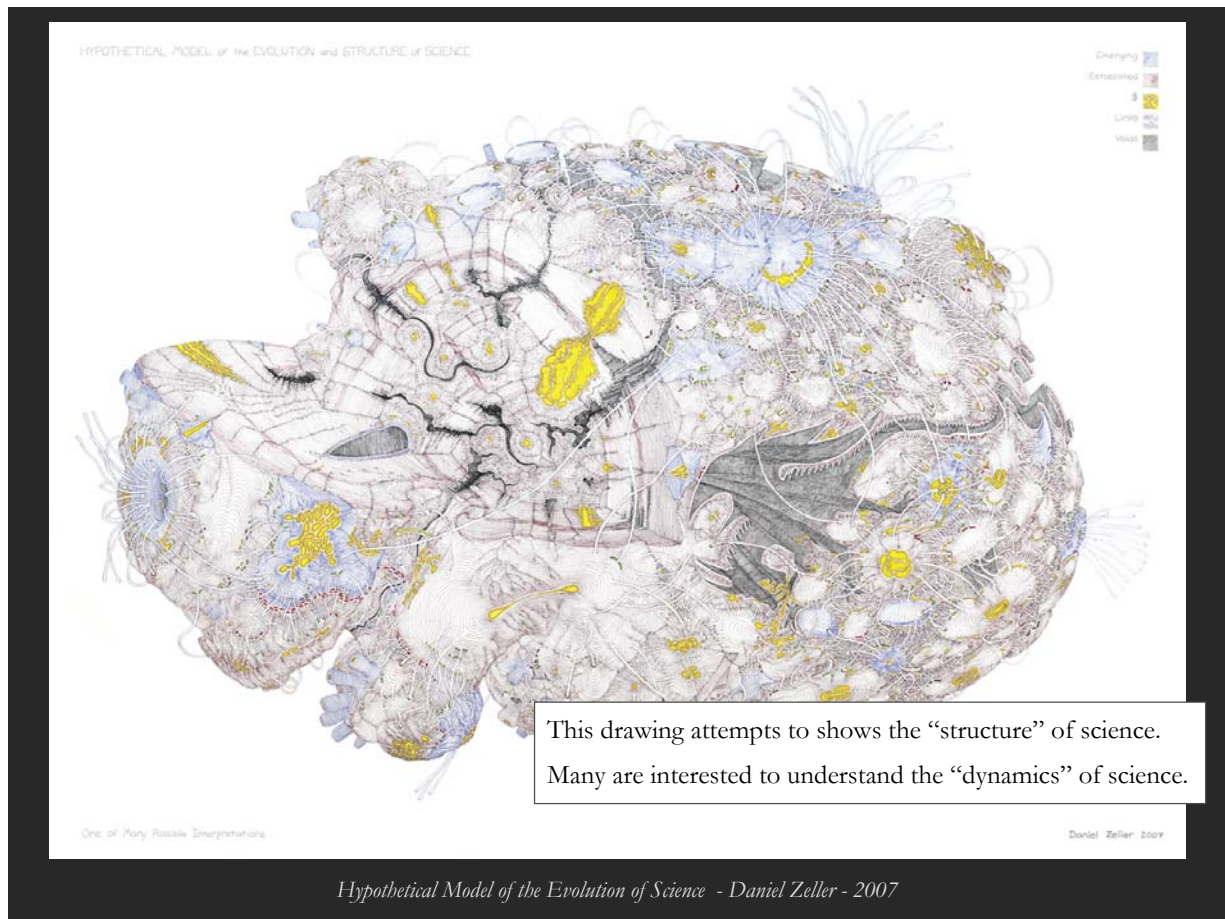
HistCite™ Visualization of DNA Development - Eugene Garfield, Elisha Hardy, Katy Borner, Ludmila Pollock, Jan Witkowski - 2006



Hypothetical Model of the Evolution of Science - Daniel Zeller - 2007



A Clickstream Map of Science – Bollen, Johan, Herbert Van de Sompel, Aric Hagberg, Luis M.A. Bettencourt, Ryan Chute, Marko A. Rodriguez, Lyudmila Balakireva - 2008



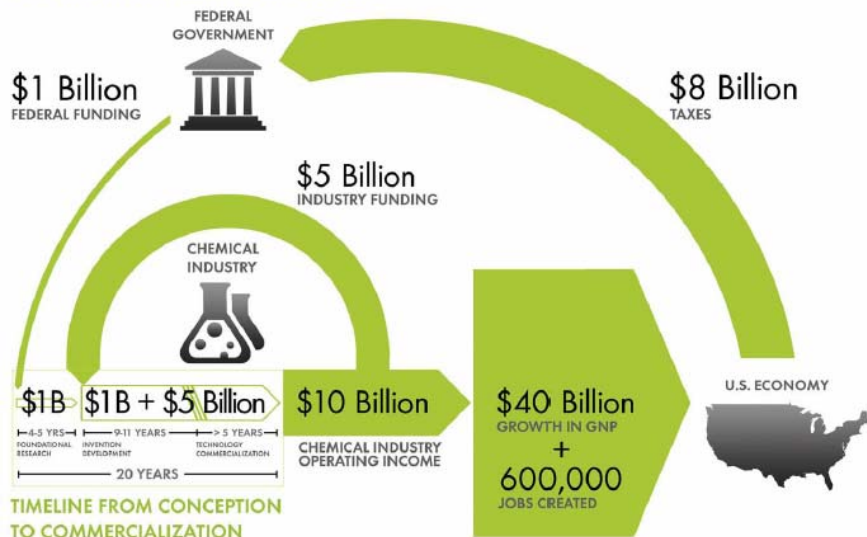
Chemical Research & Development Powers the U.S. Innovation Engine

Macroeconomic Implications of Public and Private R&D Investments in Chemical Sciences

The Council for Chemical Research (CCR)

has provided the U.S. Congress and government policy makers with important results regarding the impact of Federal Research & Development (R&D) investments on U.S. innovation and global competitiveness through its commissioned 5-year two phase study. To take full advantage of typically brief access to policy makers, CCR developed the graphic below as a communication tool that distills the complex data produced by these studies in direct, concise and clear terms.

INVESTMENT IN CHEMICAL SCIENCE R&D



The design shows that an input of \$1B in federal investment, leveraged by \$5B industry investment, brings new technologies to market and results in \$10B of operating income for the chemical industry, \$40B growth in the Gross National Product (GNP) and further impacts the US economy by generating approximately 600,000 jobs, along with a return of \$8B in taxes. Additional details, also reported in the CCR studies, are depicted in the map to the left. This map clearly shows the two R&D investment cycles; the shorter industry investment cycle; the shorter industry investment cycle; and the longer federal investment cycle which begins in basic research and culminates in national economic and job growth along with the increase tax base that in turn is available for investment in basic research.

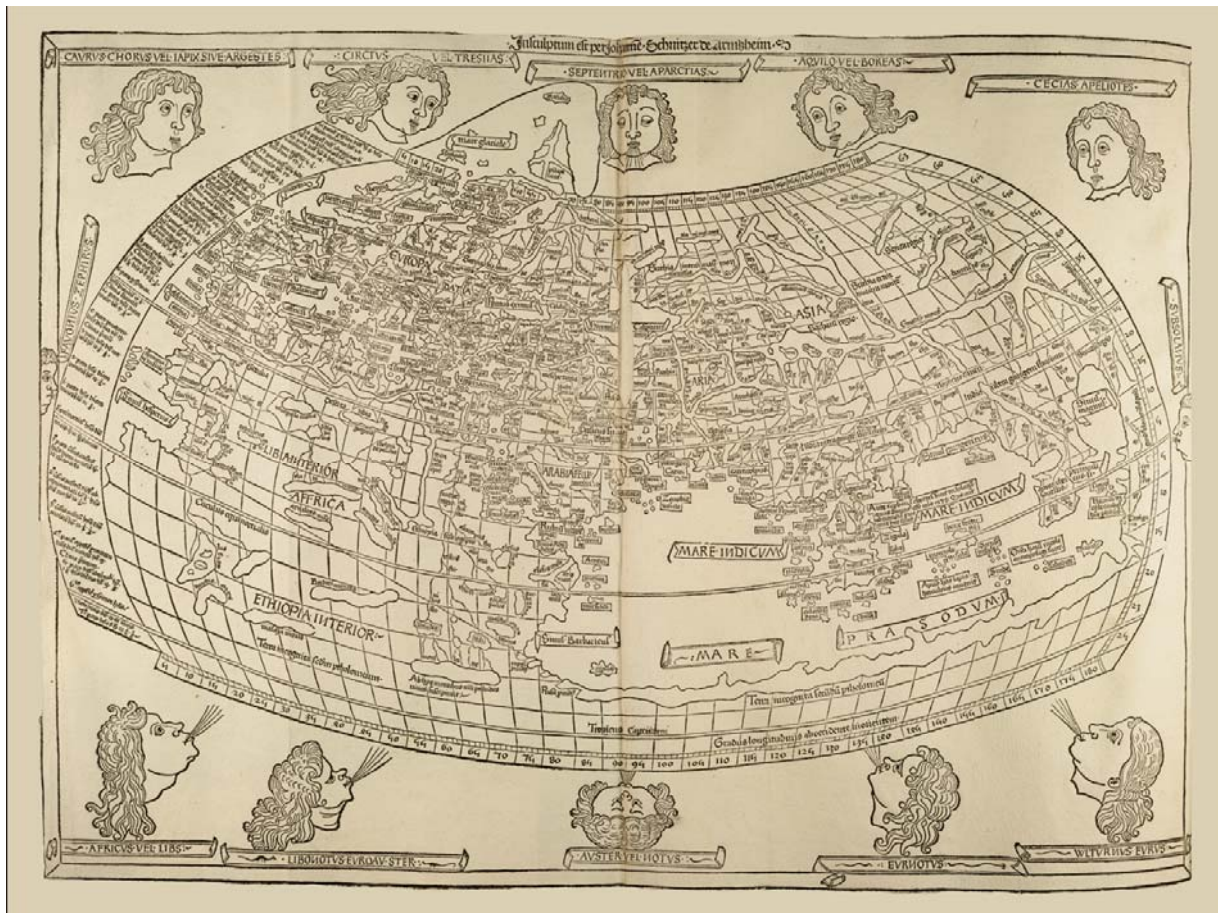
Council for Chemical Research - Chemical R&D Powers the U.S. Innovation Engine. Washington, DC. Courtesy of the Council for Chemical Research - 2009

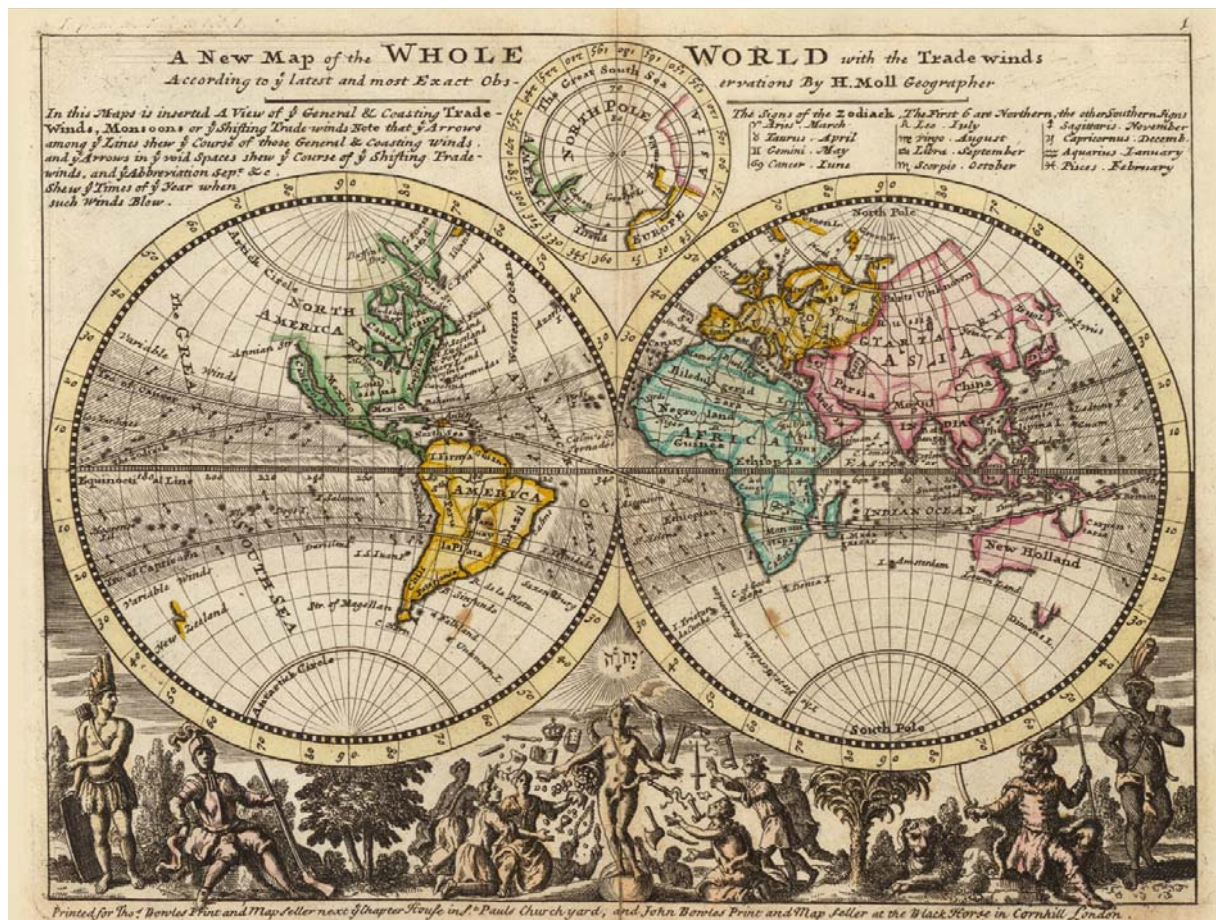
Why Map Science?

Cartographic maps of physical places have guided mankind's explorations for centuries.

They enabled the discovery of new worlds while also marking territories inhabited by the unknown.

Without maps, we would be lost.





Domain maps of abstract semantic spaces aim to serve today's explorers navigating the world of science.

These maps are generated through a scientific analysis of large-scale scholarly datasets in an effort to connect and make sense of the bits and pieces of knowledge they contain.

They can be used to identify objectively major research areas, experts, institutions, collections, grants, papers, journals, and ideas in a domain of interest. Science maps can provide overviews of "all-of-science" or of a specific area.

They can show homogeneity vs. heterogeneity, cause and effect, and relative speed. They allow us to track the emergence, evolution, and disappearance of topics and help to identify the most promising areas of research.

Information Needs for Science Map User Groups

Advantages for Funding Agencies

- Supports **monitoring** of (long-term) money flow and research developments, **evaluation** of funding strategies for different programs, **decisions** on project durations, funding **patterns**.
- Staff resources can be used for **scientific program development**, to identify areas for future development, and the stimulation of new research areas.

Advantages for Researchers

- **Easy access** to research results, relevant funding programs and their success rates, potential collaborators, competitors, related projects/publications (**research push**).
- **More time** for research and teaching.

Advantages for Industry

- Fast and **easy access** to major results, experts, etc.
- Can **influence** the direction of research by entering information on needed technologies (**industry-pull**).

Advantages for Publishers

- **Unique interface** to their data.
- **Publicly funded** development of databases and their interlinkage.

For Society

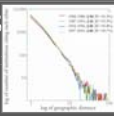








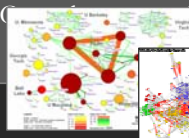

- Dramatically improved **access to scientific knowledge** and expertise.

Analysis and Visualization of Science

Type of Analysis vs. Scale of Level of Analysis

	Micro/Individual (1-100 records)	Meso/Local (101–10,000 records)	Macro/Global (10,000 < records)
Statistical Analysis/Profiling	Individual person and their expertise profiles	Larger labs, centers, universities, research domains, or states	All of NSF, all of USA, all of science.
Temporal Analysis (When)	Funding portfolio of one individual	Mapping topic bursts in 20-years of PNAS	113 Years of physics Research
Geospatial Analysis (Where)	Career trajectory of one individual	Mapping a states intellectual landscape	PNAS publications
Topical Analysis (What)	Base knowledge from which one grant draws.	Knowledge flows in Chemistry research	VxOrd/Topic maps of NIH funding
Network Analysis (With Whom?)	NSF Co-PI network of one individual	Co-author network	NSF's core competency

Type of Analysis vs. Scale of Level of Analysis

	Micro/Individual (1-100 records)	Meso/Local (101–10,000 records)	Macro/Global (10,000 < records)
Statistical Analysis/Profiling	Individual person and their expertise profiles	Larger labs, centers, universities, research domains, or states	All of NSF, all of USA, all of science. 
Temporal Analysis (When)	Funding portfolio of one individual	Mapping topic bursts in 20-years of PNAS 	113 Years of physics Research 
Geospatial Analysis (Where)	Career trajectory of one individual	Mapping a states intellectual landscape 	PNAS publications 
Topical Analysis (What)	Base knowledge from which one grant draws. 	Knowledge flows in Chemistry research 	VxOrd/Topic maps of NIH funding 
Network Analysis (With Whom?)	NSF Co-PI network of one individual 	Co-author network 	NSF's core competency 

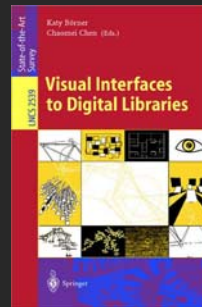
Process of Computational Scientometrics

Data Extraction	Unit of Analysis	Measures	Layout (often one code does both similarity and ordination steps)		Display
			Similarity	Ordination	
Searches •ISI •INSPEC •Eng Index •Medline •ResearchIndex •Patents •etc. Broadening •By citation •By terms	Common Choices •Journal •Document •Author •Term	Counts/Frequencies •Attributes (e.g., terms) •Author citations •Co-citations •By year Thresholds •By counts	Scalar (unit by unit matrix) •Direct citation •Co-citation •Combined linkage •Co-word/co-term •Co-classification Vector (unit by attribute matrix) •Vector space model (words/terms) •Latent Semantic Analysis (words/terms) incl. Singular Value Decomp (SVD) Correlation (if desired) •Pearson's R on any of above	Dimensionality Reduction: •Eigenvector/Eigenvalue solutions •Factor Analysis (FA) and Principal Components Analysis (PCA) •Multi-dimensional scaling (MDS) •LSA •Pathfinder networks (PFNet) •Self-organizing maps (SOM) incl. SOM, ET-maps, etc. Cluster analysis Scalar •Triangulation •Force-directed placement (FDP)	Interaction •Browse •Pan •Zoom •Filter •Query •Detail on demand Analysis

Börner, Katy, Chen, Chaomei, and Boyack, Kevin. (2003) *Visualizing Knowledge Domains*. In Blaise Cronin (Ed.), *Annual Review of Information Science & Technology, Volume 37*, Medford, NJ: Information Today, Inc./ American Society for Information Science and Technology, chapter 5, pp. 179-255.

Computational Scientometrics: Studying Science by Scientific Means

- Börner, Katy, Chen, Chaomei, and Boyack, Kevin. (2003). *Visualizing Knowledge Domains*. In Blaise Cronin (Ed.), *Annual Review of Information Science & Technology*, Medford, NJ: Information Today, Inc./ American Society for Information Science and Technology, Volume 37, Chapter 5, pp. 179-255. <http://ivl.slis.indiana.edu/km/pub/2003-borner-arist.pdf>
- Shiffrin, Richard M. and Börner, Katy (Eds.) (2004). *Mapping Knowledge Domains*. *Proceedings of the National Academy of Sciences of the United States of America*, 101(Suppl_1). http://www.pnas.org/content/vol101/suppl_1/
- Börner, Katy, Sanyal, Soma and Vespignani, Alessandro (2007). *Network Science*. In Blaise Cronin (Ed.), *Annual Review of Information Science & Technology*, Information Today, Inc./ American Society for Information Science and Technology, Medford, NJ, Volume 41, Chapter 12, pp. 537-607. <http://ivl.slis.indiana.edu/km/pub/2007-borner-arist.pdf>
- Places & Spaces: *Mapping Science* exhibit, see also <http://scimaps.org>.
- Börner, Katy. (2010). *Atlas of Science: Visualizing What We Know*. MIT Press. <http://scimaps.org/atlas>
- Börner, Katy. (March 2011). *Plug-and-Play Macroscopes*. *Communications of the ACM*.



Science of Science Cyberinfrastructure

Overview


What cyberinfrastructure will be required to measure, model, analyze, and communicate scholarly data and, ultimately, scientific progress?

This talk presents our efforts to create a science of science cyberinfrastructure that supports:

- Data access and federation via the **Scholarly Database**, <http://sdb.slis.indiana.edu>,
- Data preprocessing, modeling, analysis, and visualization using plug-and-play cyberinfrastructures such as the **Sci² Tool**, <http://sci2.cns.iu.edu>, and
- Communication of science to a general audience via the **Mapping Science Exhibit** at <http://scimaps.org>.

The following demos should be particularly interesting for those interested to

- Map their very own domain of research,
- Test and compare data federation, mining, visualization algorithms on large scale datasets,
- Use advanced network science algorithms in their own research.



Science of Science Cyberinfrastructure
— P O R T A L —

Provided by the [Cyberinfrastructure for Network Science Center](#) at Indiana University.

Introduction
E. O. Wilson writes in *Consilience: The Unity of Knowledge* (1998): "Features that distinguish science from pseudoscience are repeatability, economy, menturation, heuristics, and consilience." Please see Börner's [recent presentation](#) at the *Deeper Look at the Visualization of Scientific Discovery* NSF Workshop for a general introduction of the needs and the resources provided here.

Needs Analysis
As part of the "FIS: [Towards a Macroscopic for Science Policy Decision Making](#)" NSF SBE-0738111 award, interviews with science policy makers are conducted to identify what science of science research results and tools might be most desirable and effective. So far, 30 formal, one-hour interviews have been conducted with science policy makers at university campus level, program officer level, and division director level for governmental, state, and private foundations. Data compilation will start in October 2008 and resulting report can be ordered by sending a request to Mark Price (maaprice@indiana.edu).

Conceptualization of Science
A 'science of science' requires a theoretically grounded and practically useful conceptualization of the structure and evolution of science. A special journal issue entitled "[Science of Science: Conceptualizations and Models of Science](#)" edited by [Katy Börner](#), Indiana University & [Andrea Scharnhorst](#), Royal Netherlands Academy of Arts and Sciences invites contributions on this topic. It will be published in the *Journal of Informetrics* 3(1) in January 2009.

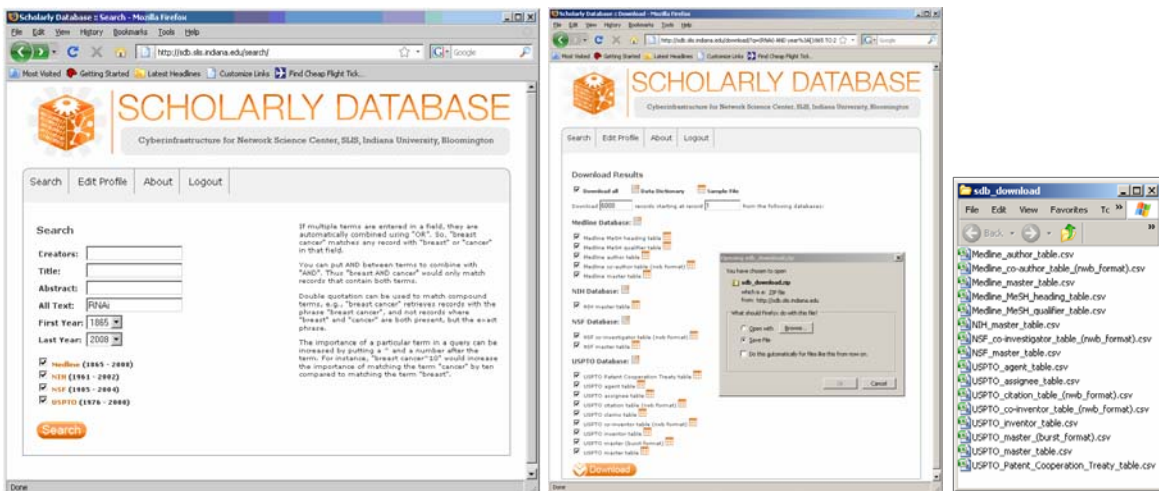
Scholarly Database
The [Scholarly Database \(SDB\)](#) at Indiana University aims to serve researchers and practitioners interested in the analysis, modeling, and visualization of large-scale scholarly datasets. The database currently provides access to over 20 million papers, patents and grants. Resulting datasets can be downloaded in bulk. Register for free access at <https://sdb.slis.indiana.edu/>.

Cyberinfrastructures
The Scientometrics filing of the [Network Workbench \(NWB\) Tool](#) provides a unique distributed, shared resources environment for large-scale network analysis, modeling, and visualization. Thomson Scientific/ISI, Scopus and Google Scholar data, EndNote and BibTeX files, or NSF awards can be read and diverse networks can be extracted and studied. Download [User Manual with focus on Scientometrics](#).

<http://sci.slis.indiana.edu>



Scholarly Database (<http://sdb.cns.iu.edu>)



The image shows three screenshots of the Scholarly Database interface. The first screenshot displays the search page with fields for search, title, abstract, and filters for Medline (1965-2008), ISI (1961-2002), ISI (1965-2004), and USPTO (1976-2008). The second screenshot shows the 'Download Results' page with a list of download options for Medline, ISI, and USPTO databases, including various table formats like 'Medline_author_table.csv' and 'USPTO_patent_cooperation_treaty_table.csv'. The third screenshot shows a file explorer window displaying the downloaded files.

The Scholarly Database at Indiana University provides free access to 25,000,000 papers, patents, and grants. Since March 2009, users can also download networks, e.g., co-author, co-investigator, co-inventor, patent citation, and tables for burst analysis.



Sci² Tool for Science of Science Science (<http://sci2.cns.iu.edu>)

- Explicitly designed for SoS research and practice, well documented, easy to use.
- Empowers many to run common studies while making it easy for exports to perform novel research.
- Advanced algorithms, effective visualizations, and many (standard) workflows.
- Supports micro-level documentation and replication of studies.
- Is open source—anybody can review and extend the code, or use it for commercial purposes.

nature Vol 464|25 March 2010

OPINION

Let's make science metrics more scientific

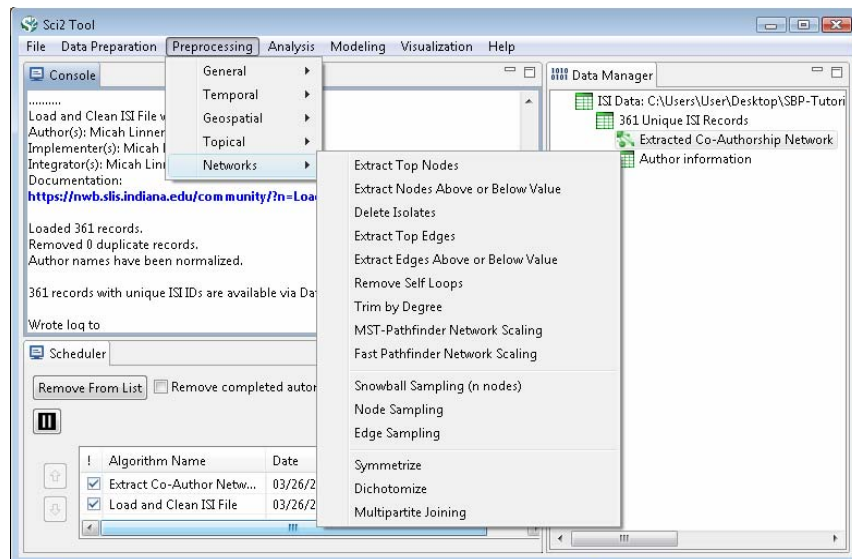
To capture the essence of good science, stakeholders must combine forces to create an open, sound and consistent system for measuring all the activities that make up academic productivity, says **Julia Lane**.

SUMMARY

- Existing metrics have known flaws
- A reliable, open, joined-up data infrastructure is needed
- Data should be collected on the full range of scientists' work
- Social scientists and economists should be involved



Sci² Tool for Science of Science Research and Practice



Acknowledgments

This work is supported in part by the Cyberinfrastructure for Network Science center and the School of Library and Information Science at Indiana University, the National Science Foundation under Grant No. SBE-0738111 and IIS-0513650, and the James S. McDonnell Foundation.





Sci² Tool for Science of Science Research and Practice

Supported Input file formats:

- GraphML (*.xml or *.graphml)
- XGMML (*.xml)
- Pajek .NET (*.net) & Pajek .Matrix (*.mat)
- NWB (*.nwb)
- TreeML (*.xml)
- Edge list (*.edge)
- **CSV (*.csv)**
- **ISI (*.isi)**
- **Scopus (*.scopus)**
- **NSF (*.nsf)**
- **Bibtex (*.bib)**
- **Endnote (*.enw)**

Output file formats:

- GraphML (*.xml or *.graphml)
- Pajek .MAT (*.mat)
- Pajek .NET (*.net)
- NWB (*.nwb)
- XGMML (*.xml)
- CSV (*.csv)

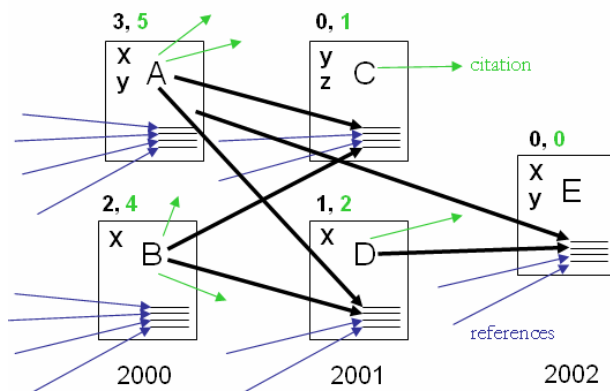
<http://sci2.wiki.cns.iu.edu/2.3+Data+Formats>



Network Extraction

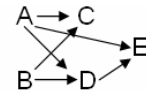
Sample paper network (left) and four different network types derived from it (right). From ISI files, about 30 different networks can be extracted.

Papers A-E written by authors x, y, z over 3 years. Each paper happens to have 4 references.



Paper-Paper Citation Network

Papers are connected via direct citation links. Arrows represent information flow from older papers to younger papers.



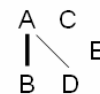
Author-Author (Co-Author) Network

x and y co-author papers A and E together y and z co-author papers A and E



Document Co-Citation (DCA) Network

A and B are co-cited by C and D A and D are co-cited by E



Reference Co-Occurrence (Bibliographic Coupling) Network

C and D are bibliographically coupled as they both cite/reference A and B.

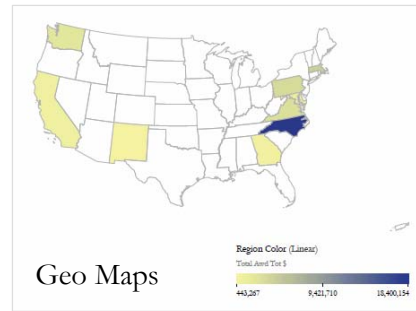
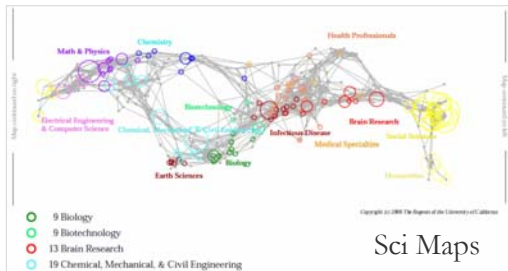


Local citation counts (within this dataset) are given in **black** and global citation counts (ISI times cited) are given in **green** above each paper.



Sci² Tool

Plugins that render into Postscript files:



Börner, Katy, Huang, Weixia (Bonnie), Linnemeier, Micah, Dubon, Russell Jackson, Phillips, Patrick, Ma, Nianli, Zoss, Angela, Guo, Hanning & Price, Mark. (2009). *Rete-Netzwerk-Red: Analyzing and Visualizing Scholarly Networks Using the Scholarly Database and the Network Workbench Tool*. *Proceedings of ISSI 2009: 12th International Conference on Scientometrics and Informetrics, Rio de Janeiro, Brazil, July 14-17*. Vol. 2, pp. 619-630.

Exemplary Analyses and Visualizations

Individual Level

- Loading ISI files of major network science researchers, extracting, analyzing and visualizing paper-citation networks and co-author networks (p. 54-65)
- Loading NSF datasets with currently active NSF funding for 3 researchers at Indiana U (p. 49-53)

Institution Level

- Indiana U, Cornell U, and Michigan U, extracting, and comparing Co-PI networks (p. 65-69)

Scientific Field Level

- Extracting co-author networks, patent-citation networks, and detecting bursts in SDB data (p. 77-85)



cyberinfrastructure for NETWORK SCIENCE CENTER

School of Library and Information Science | Indiana University Bloomington



All papers, maps, cyberinfrastructures, talks, press are linked
from <http://cns.iu.edu>