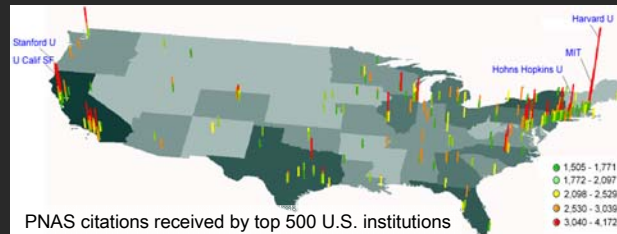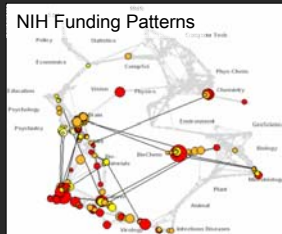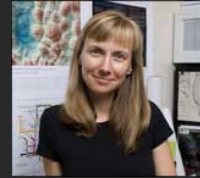Mapping the Structure and Evolution of Scholarly Knowledge:
## Data (Integration) Issues

**Dr. Katy Börner**
Cyberinfrastructure for Network Science Center, Director
Information Visualization Laboratory, Director
School of Library and Information Science
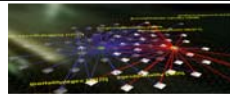Indiana University, Bloomington, IN
*katy@indiana.edu*

*Workshop on Scholarly Databases & Data Integration*
*Indiana University, Bloomington, IN, Aug. 29 & 30, 2006*

NIH Funding Patterns

PNAS citations received by top 500 U.S. institutions

1. Dream Tools for Scholarly Knowledge Management

2. Challenges
3. Opportunities

1. **Dream Tools for Scholarly Knowledge Management**

2. Challenges
3. Opportunities

---

**Dream Tools for Scholarly Knowledge Management**

Tools we developed for ourselves and our clients

**Information Visualization Laboratory Management System**

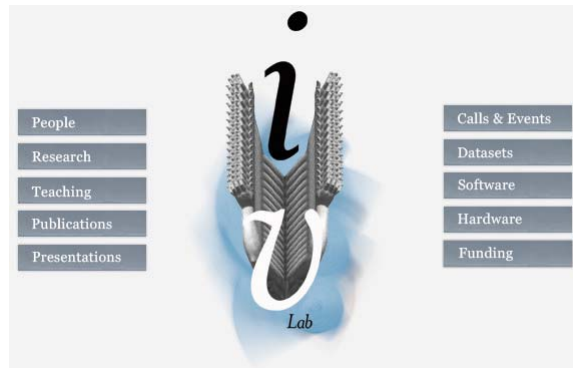**Taxonomy Visualization/ Validation System**

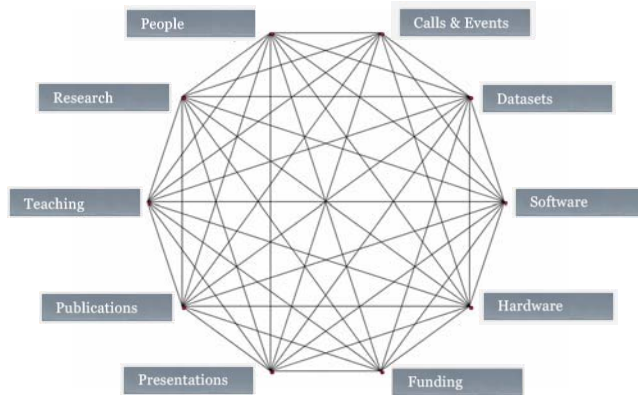**Coupling Geospatial and Topic Space**

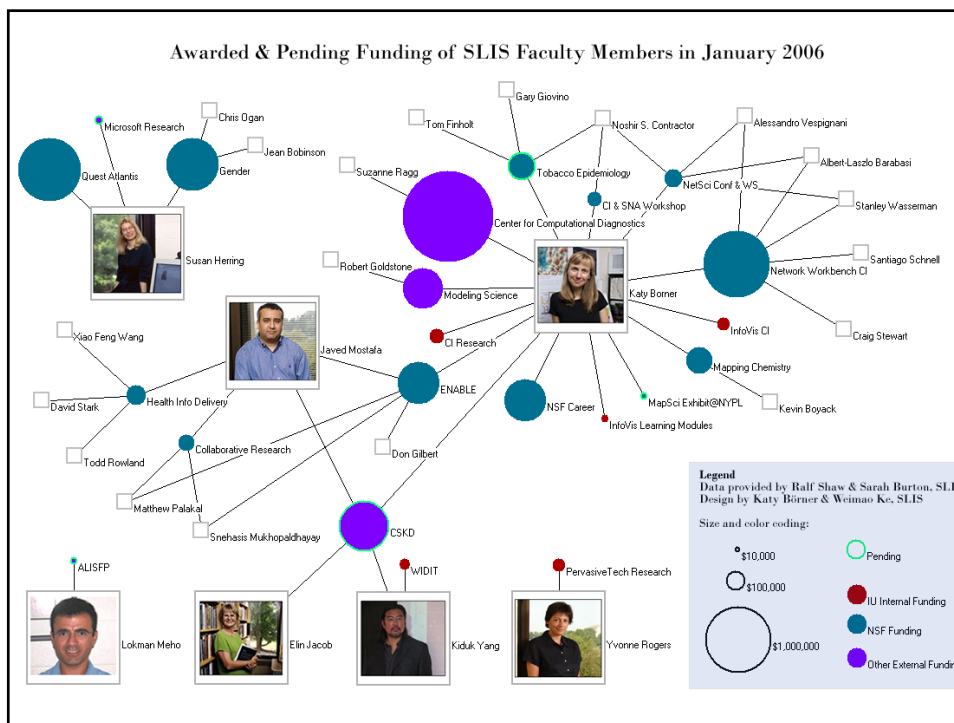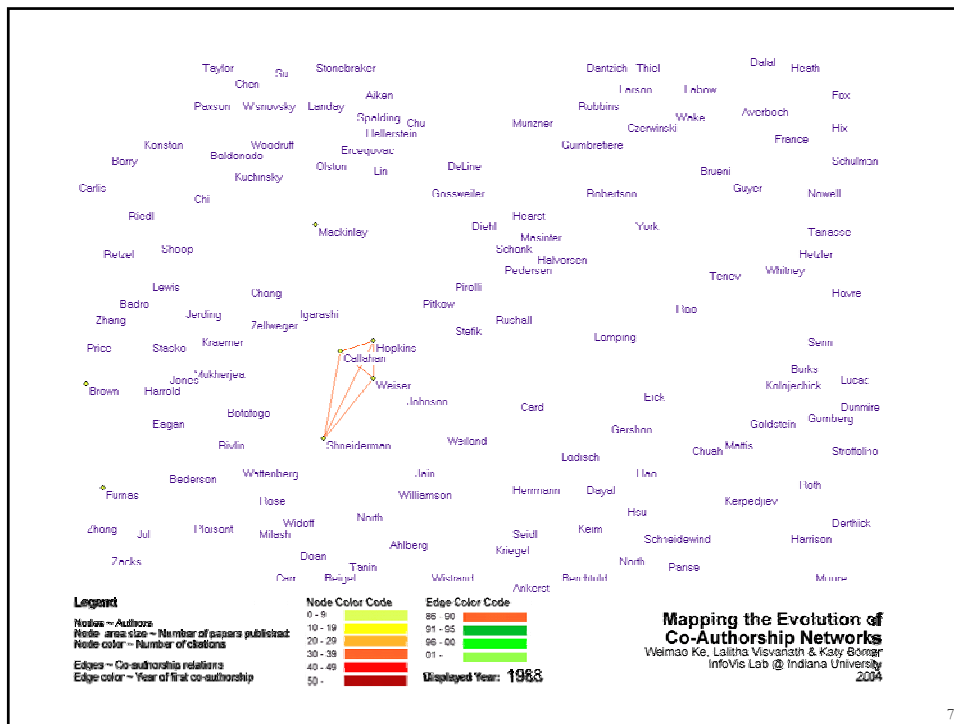**Science Maps for Kids**

4

**Information Visualization Laboratory
Management System**



https://ivl.slis.indiana.edu

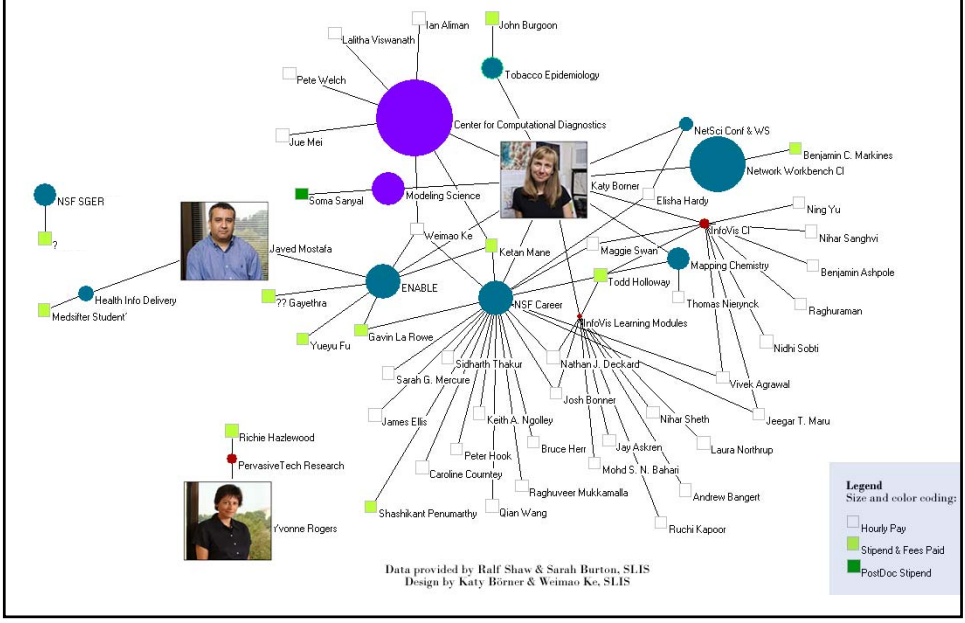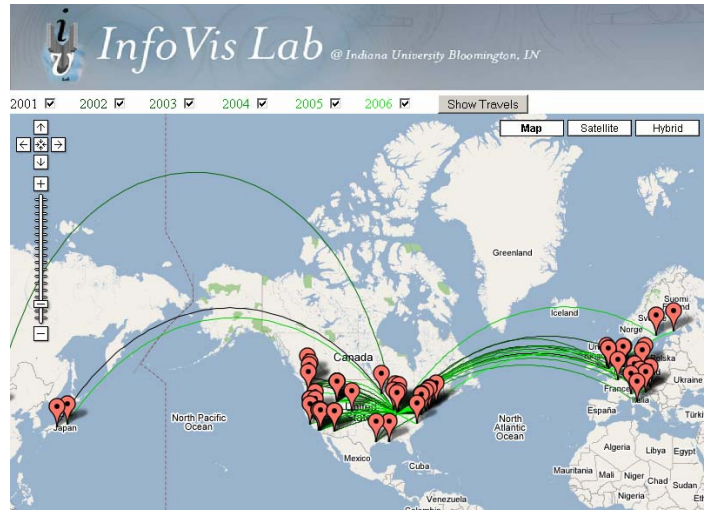Awarded & Pending Funding of SLIS Faculty Members in January 2006

Student Support by SLIS Faculty Members in January 2006



KM: Katy's Travels in 2000-2006

*By Thomas Neirynck, 2007.*



*By Thomas Neirynck, 2007.*

*By Thomas Neirynck, 2007.*



*By Thomas Neirynck, 2007.*

*By Thomas Neirynck, 2007.*



MARC records, Wiki authors

Science in the City

Courses at IU

People

Calls & Events

IVC, NWB data

Research

Datasets

IVC, NWB code

Teaching

Software

Publication DBs

IT Manual

Publications

Hardware

Presentations

Funding

NetSci, IM2 Conference talks

NSF, NIH awards data, IU awards

# US Patent Hierarchy

**Impact**

**Prior Art**

---

Synthetic Resins or Natural Rubbe

Ion-exchange Polymer or Process of Prepari
 Process of Regenerating
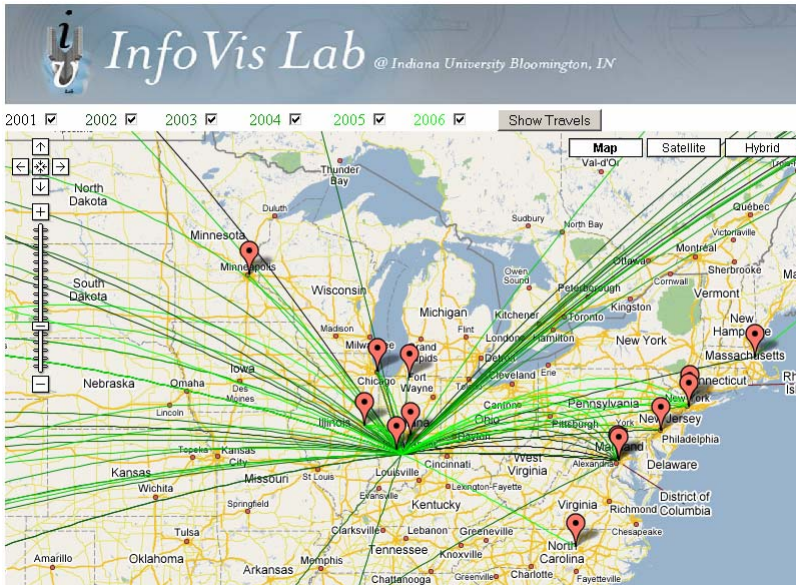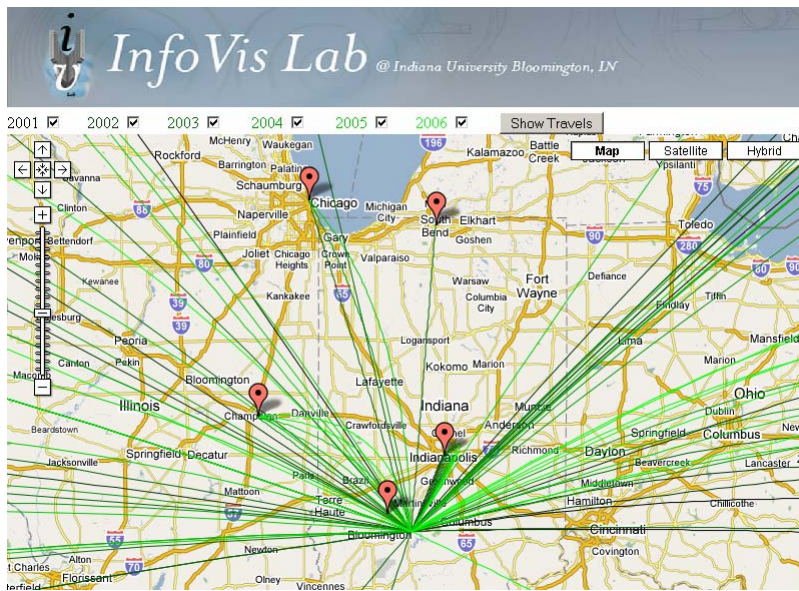 Membrane or Process of Preparing
 Previously Formed Solid Ion-exchange Polymer Admixed With N
 Polymer Characterized By Defined Size or Shape Other than Bea
 Chemically Treated Solid Polymer
  Solid Polymer Derived From Ethylenically Unsaturated Reacta
  Solid Polymer Derived From At Least One 1,2-epoxy Containir
  Solid Polymer Derived From Aldehyde or Derivative
 From Ethylenically Unsaturated Reactant Only
 From Aldehyde or Derivative

Process of Treating Scrap or Waste Product (
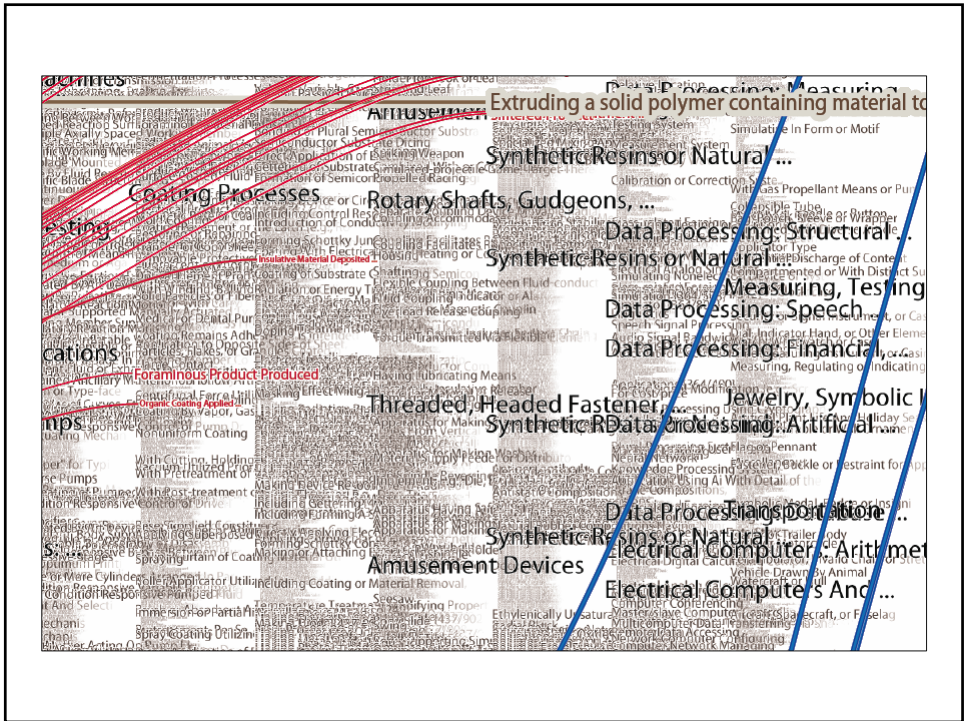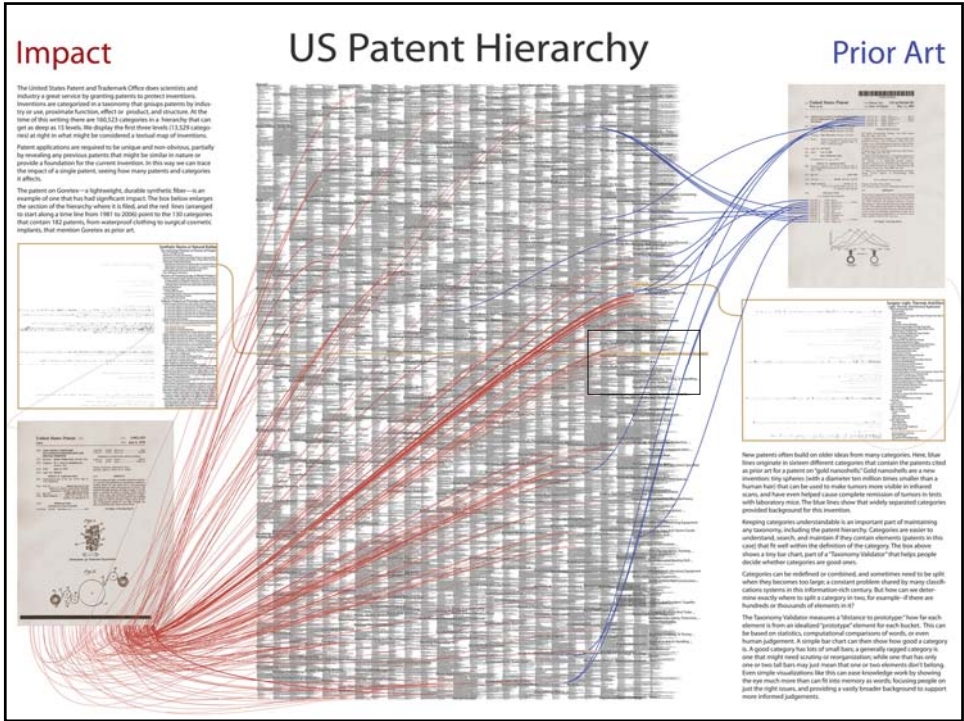 Process of Treating Scrap or Waste Product Containing At Least
  Treating Rubber (or Rubberlike Materials) or Polymer Derived
  Treating Polymer Derived From A Monomer Containing Only (
  Treating Polymer Derived From Hydrocarbon Monomers Only
 Treating Polysiloxane
 Treating Polyester
  Treating With Alcohol
 Treating Polyurethane, Polyurea (excluding Urea-formaldehyde
  Treating With Alcohol or Amine
 Treating Polycarbonamide

Cellular Products or Processes of Preparing /
 Cellular Product Derived From Two or More Solid Polymers or Fr
  At Least One Polymer Is Derived From Reactant Containing Tw
  At Least One Polymer Is Derived From An Aldehyde or Derivat
  At Least One Polymer Is Derived From A -n=c=x Reactant Whe

*Places & Spaces: Mapping Science* exhibit, see also [http://scimaps.org](http://scimaps.org).

"Places & Spaces: Mapping Science"
on display at the NYPL Science, Industry, and Business Library
Madison/34th, New York City
April 3rd - August 31st, 2006.

TOPIC MAP: HOW SCIENTIFIC PARADIGMS RELATE

*Science Puzzle Map for Kids by Fileve Palmer, Julie Smith, Elisha Hardy and Katy Börner, Indiana University, 2006. (Base map taken from Illuminated Diagram display by Kevin Boyack, Richard Klavans, and W. Bradford Paley.)*

Hands-On Science Maps for Kids, by Flávia Palmer (Painting), Julie Smith (Data Acquisition), Elisha Hardy and Katy Börner (Graphic Design), BLOOMINGTON, IN, 2006, Courtesy of Indiana University. Learn more at www.scimaps.org. This map plots the locations of where scientific papers were published each light green dot represents 10 or fewer papers; they are scattered around the exact location for visibility within a labelled green circle whose size is proportional to the number of papers published in that place. The base map is part of an "illuminated diagram" display which used a computer and two projectors, projecting spots of light on the prints to highlight different kinds of scientific research (on a sibling map of scientific paradigms) and the areas in the world where such science was performed. Base map research by Kevin Boyack and Dick Klavans, cartography by John Burgoon, data from Thompson ISI, graphics and typography by W. Bradford Paley. Copyright © 2006 W. Bradford Paley, all rights reserved.



Hands-On Science Maps for Kids, by Flávia Palmer (Painting), Julie Smith (Data Acquisition), Elisha Hardy and Katy Börner (Graphic Design), BLOOMINGTON, IN, 2006, Courtesy of Indiana University. Learn more at www.scimaps.org. This map plots the locations of where scientific papers were published each light green dot represents 10 or fewer papers; they are scattered around the exact location for visibility within a labelled green circle whose size is proportional to the number of papers published in that place. The base map is part of an "illuminated diagram" display which used a computer and two projectors, projecting spots of light on the prints to highlight different kinds of scientific research (on a sibling map of scientific paradigms) and the areas in the world where such science was performed. Base map research by Kevin Boyack and Dick Klavans, cartography by John Burgoon, data from Thompson ISI, graphics and typography by W. Bradford Paley. Copyright © 2006 W. Bradford Paley, all rights reserved.

1. Dream Tools for Scholarly Knowledge Management

**2. Challenges**
3. Opportunities

---

**Challenges - Data Collection & Integration**

Grants  Patents  Papers in Area A  Papers in Area B

Time

Scholarly Knowledge  ○  Citation Links  →

**Figure 1:** The interoperability and cross linkage problem. Many but not all of today's scholarly datasets, e.g., papers, patents, grants, are stored and made available so that 'vertical' citation linkages can be traversed. There are very few instances in which datasets of different origin and/or type are 'horizontally' interlinked.

*Börner, K. (2006) Semantic Association Networks: Using Semantic Web Technology to Improve Scholarly Knowledge and Expertise Management. In Vladimir Geroimenko & Chaomei Chen (eds.) Visualizing the Semantic Web, Springer Verlag, 2nd Edition, chapter 11, pp. 183-198.*

38

**Challenges - Semantic Mining / Integration**

**Works well** if the records are written in similar styles, using similar formatting and conventions, are of similar length, etc.

**Works less well** if applied to interdisciplinary or multi-lingual datasets because
- Words, e.g., 'prototype', have very different meanings in computer science, biology, psychology, architecture, etc.
- Paper titles are frequently used to demonstrate the creativity of authors, e.g., "All you ever wanted to know about x", "A unifying theory of x".
- Author supplied keywords are useful to identify an author but not to find similar papers. Note: Controlled vocabularies/thesauri work well.

Humans might simply be too different and too creative to produce proper raw material that can be analyzed using existing text mining and data mining algorithms.

39

**Challenges - Link Traversal & Link Mining**

**Link Search:** Google

**(Citation) Link Traversal in One Database:** *Thomson Scientific, Google Scholar,* and *CiteSeer* already support citation link traversal. The *Proceedings of the National Academy of Sciences of the United States of* America (PNAS) online interface <http://www.pnas.org>. provides citation maps that shows articles citing or being cited by a selected article.

**(Citation) Link Traversal Across Databases:** The *Library Without Walls* project <http://library.lanl.gov/lww/> at the *Los Alamos National Laboratory* interlinks major publication databases and supporting citation based search across different holdings that have citation linkages.

**(Co-Author) Link Traversal:** Some digital libraries such as the citation indexes published by *Thomson Scientific*, *DBLP Bibliography Server* <http://www.sigmod.org/dblp/db/>, and *ACM Digital Library* <http://portal.acm.org> provide information on co-authorships. Services comprise a listing of all papers by an author, a listing of all co-authors for one author, and co-author link traversal.

40

**Challenges - Interlink $ Input & Publication/Patent Citation Output**



**Need to interlink**

➤ Grants and papers/patents.
➤ Grants/papers/patents and their PIs/authors/inventors, etc.

**Use resulting networks to**

➤ Count #papers, #citations, etc.
➤ Determine strength of co-PI/author/inventor relations, etc.

41

---

**Improved Representation of Scholarly Knowledge**

**Entity and link types:**



**Attributes:**

➤ Records often have a publication date, a publication type (e.g., journal paper, book, patents, grant, etc.), topics (e.g., keywords or classifications assigned by authors and/or publishers).
➤ Authors have an address with information on affiliation and geo-location.

**Derived attributes:**

➤ Because authors and records are associated, the geo-location(s) and affiliation(s) of an author can be attributed to the authors' papers.
➤ Similarly, the publication date, publication type and topic(s) can be associated with a paper's author(s).

42

**Improved Representation of Scholarly Knowledge makes possible**

**Statistics:**
- Number of papers, grants, co-authorships, citation (over time) per author.
- Bursts of activity (#citations, #$, #patents, #collaborators, etc.).
- Changes of topics and geo-locations for authors and their institutions over time.

**Visualizations:**
- Geospatial and topical distribution of funding input & research output.
- Structure and evolution of research topics.
- Evolving research areas (e.g., based on young yet highly cited papers).
- Diffusion of information, people, $s over geospatial and topic space.

43

**Semantic Association Networks**



*Katy Börner. (2006) Semantic Association Networks: Using Semantic Web Technology to Improve Scholarly Knowledge and Expertise Management. In Vladimir Geroimenko & Chaomei Chen (eds.) Visualizing the Semantic Web, Springer Verlag, 2nd Edition, chapter 11, pp. 183-198.*

44

**Open Questions**

➢ Interoperability: How to add more and more databases?

➢ Interlinkage: OAI works. What standard would work best for unique and persistent identifiers for authors/institutions/ countries/journals/geolocations/etc.?

Will 95% automatic and 5% manual data cleaning work?

➢ How to add databases/services while in production?

➢ How to exploit peer-to-peer architectures?

➢ How to resolve proprietary/political issues?

45

1. Dream Tools for Scholarly Knowledge Management

2. Challenges
3. **Opportunities**

46

## Opportunities for Mapping Science

**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.

47

## Opportunities for Modeling Science

- Dynamic science and technology indicators (emerging research frontiers, evolving networks, trends, feedback loops).
- Evolution of scientific communities/fields. Capacity limit to knowledge/ skills knowable by individual researchers.
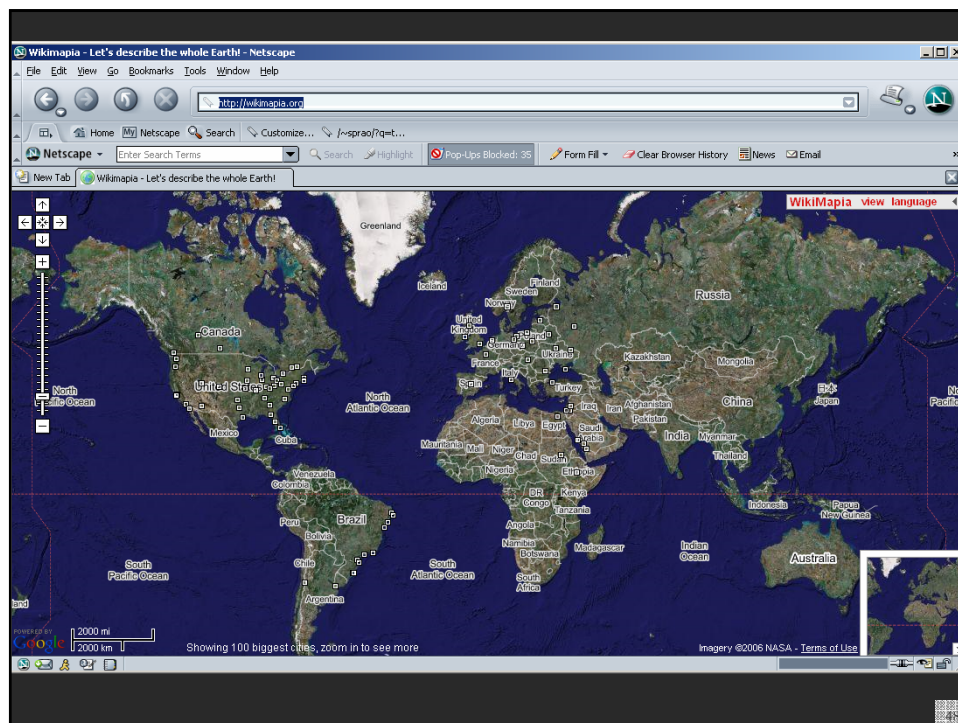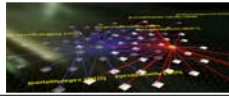- Interplay of competition and collaboration.
- Evolution of fields – birth, growth, mature, decline.
- Interactions among fields. Optimal interdisciplinary collaborations?
- Comparison of different funding models, e.g., few large vs. many small grants, teach the field how to fish or give them fish?
- Impact of publishing/collaboration/funding mechanisms on the dynamics of fields.
- Diffusion of people, ideas, skills, etc.
- How to best communicate modeling results/insights?

**Ultimate goal:**

Learn how to best increase, diffuse, and utilize our collective scholarly knowledge.

49



*CAREER: Visualizing Knowledge Domains. NSF IIS-0238261 award (Katy Börner, $451,000) Sept. 03-Aug. 08. http://iv.slis.indiana.edu/*

*SEI: Network Workbench: A Large-Scale Network Analysis, Modeling and Visualization Toolkit for Biomedical, Social Science and Physics Research. NSF IIS-0513650 award (Katy Börner, Albert-Laszlo Barabasi, Santiago Schnell, Alessandro Vespignani & Stanley Wasserman, Eric Wernert (Senior Personnel), $1,120,926) Sept. 05 - Aug. 08. http://nwb.slis.indiana.edu*

25

**CIShell** Building Market Places not Cathedrals




> Requires the design & implementation of 'software glue' that can interlink datasets and algorithms written in different languages using different data formats.
> The smaller the glue or 'CI Shell', the more likely it can be maintained.
> Dataset and algorithm 'plugins' are provided by application

51

---

**Cyberinfrastructure Shell (CIShell)**
*http://cishell.org*

CIShell is an 'empty shell' that supports

> Easy integration of new datasets and algorithms by <u>algorithm developers</u> and
> Easy usage of algorithms by <u>algorithm users.</u>

Its <u>plug-and-play architecture</u> supports the integration and utilization of diverse

> Datasets, e.g., stored in files, databases, streaming data.
> Algorithms, e.g., data processing, analysis, modeling, visualization.
> Interfaces, e.g., remote services, scripting engines, peer-to-peer clients.
> Services, e.g., workflow support, scheduler.

52

**CIShell – Needs of Algorithm Developers & Users**

Developers

Users

*CIShell Wizards*

**CIShell**

*IVC Interface*

*NWB Interface*

53



**CIShell – Deployment**

**Data-Algorithm Repositories**

**Peer-to-Peer**

**Stand Alone**

**Server-Client**

CIShell applications can be deployed as distributed data and algorithm repositories, stand alone applications, peer-to-peer architectures, and server-client architectures.

54

27

**Scholarly Database: Web Interface**

Search across publications, patents, grants.

Download records and/or (evolving) co-author, paper-citation networks.

Register for free access at https://sdb.slis.indiana.edu.

**SDB** SCHOLARLY DATABASE

**Scholarly Database: # Records & Years Covered**

Datasets available via the Scholarly Database (* future feature)
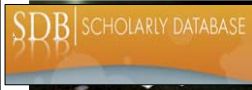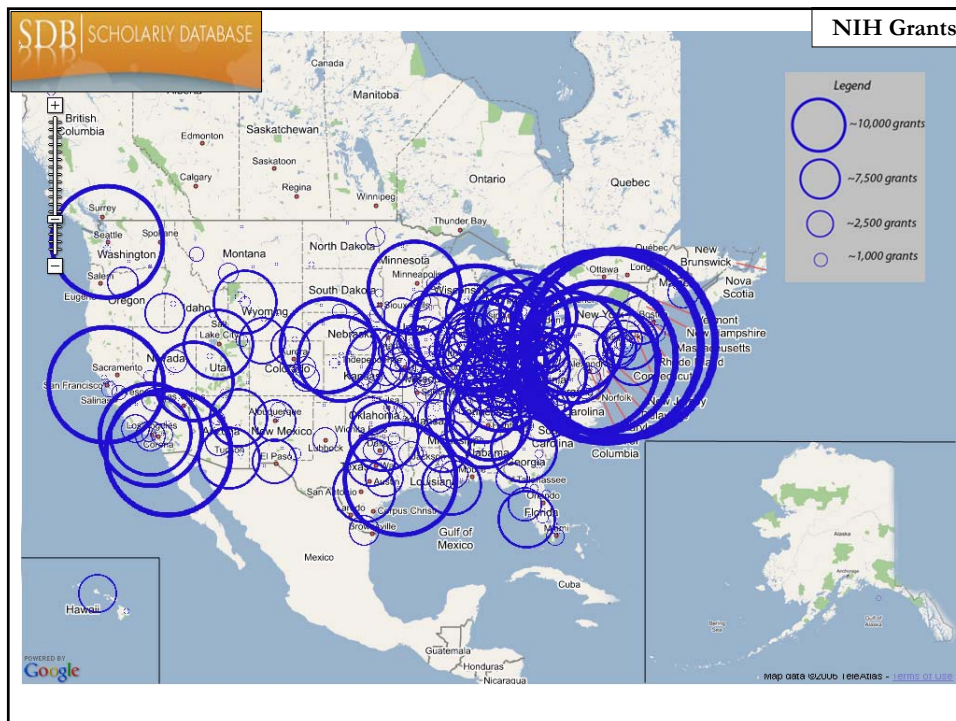
| Dataset | # Records | Years Covered | Updated | Restricted Access |
|---------|-----------|---------------|---------|-------------------|
| Medline | 13,149,741 | 1965-2005 | Yes | |
| PhysRev | 398,005 | 1893-2006 | | Yes |
| PNAS | 16,167 | 1997-2002 | | Yes |
| JCR | 59,078 | 1974, 1979, 1984, 1989 1994-2004 | | Yes |
| USPTO | 3,179,930 | 1976-2004 | Yes* | |
| NSF | 174,835 | 1985-2003 | Yes* | |
| NIH | 1,043,804 | 1972-2002 | Yes* | |
| **Total** | **18,021,560** | **1893-2006** | **4** | **3** |

Aim for comprehensive time, geospatial, and topic coverage.

57

**SDB** SCHOLARLY DATABASE



**NIH Grants**

29

**Medline Publications**



**NSF Grants**

US Patents



Science map applications: Identifying core competency
*Kevin W. Boyack & Richard Klavans, unpublished work.*

Funding patterns of the US Department of Energy (DOE)

**Science map applications: Identifying core competency**
*Kevin W. Boyack & Richard Klavans, unpublished work.*

Funding Patterns of the National Science Foundation (NSF)



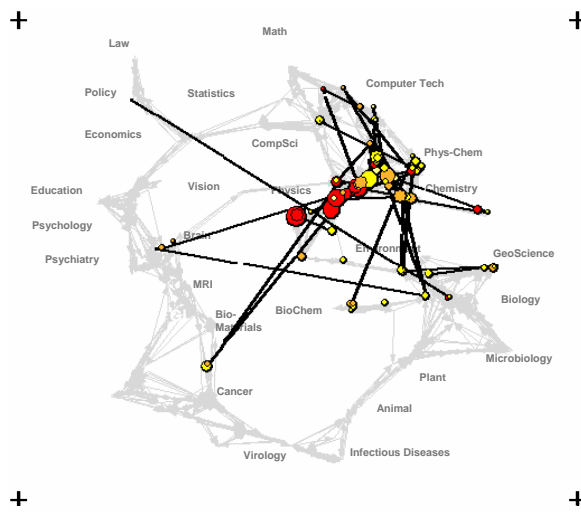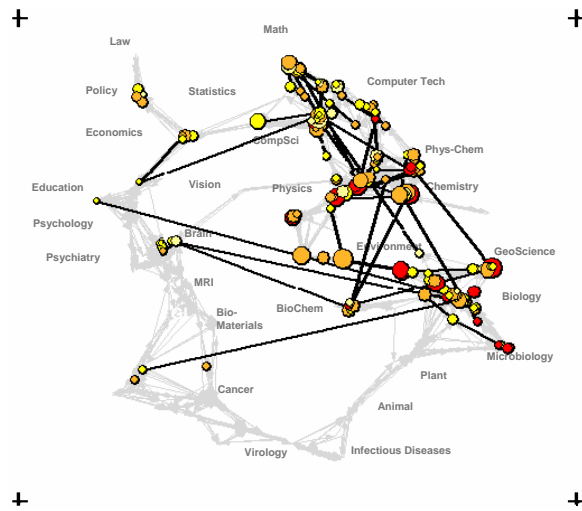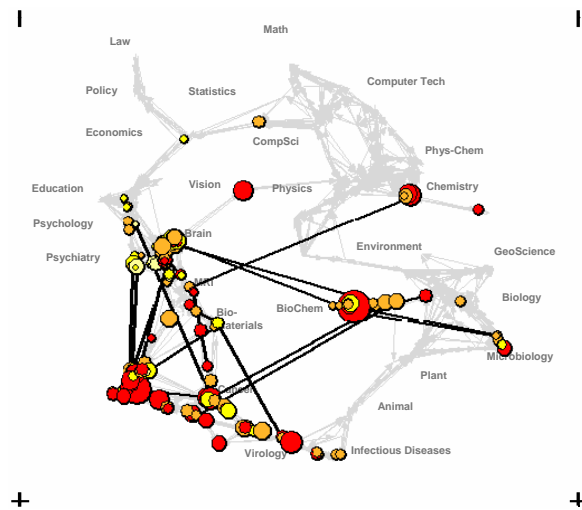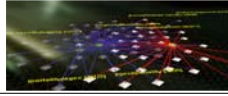**Science map applications: Identifying core competency**
*Kevin W. Boyack & Richard Klavans, unpublished work.*

Funding Patterns of the National Institutes of Health (NIH)

## References

➢ Bruce Herr, Weixia Huang, Shashikant Penumarthy, Katy Börner. Designing Highly Flexible and Usable Cyberinfrastructures for Convergence. Submitted to William S. Bainbridge (Ed.) Progress in Convergence. Annals of the New York Academy of Sciences.

➢ Börner, Katy. Mapping All of Science: How to Collect, Organize and Make Sense of Mankind's Scholarly Knowledge and Expertise. Accepted for *Environment and Planning B*, Special Issue on *Mapping Humanity's Knowledge and Expertise in the Digital Domain.*

➢ Börner, Katy, Penumarthy, Shashikant, Meiss, Mark and Ke, Weimao. (2006) Mapping the Diffusion of Scholarly Knowledge Among Major U.S. Research Institutions. *Scientometrics.* 68(3), pp. 415-426.

➢ Holloway, Todd, Božicevic, Miran and Börner, Katy. Analyzing and Visualizing the Semantic Coverage of Wikipedia and Its Authors. Accepted for *Complexity.* Also available as cs.IR/0512085.

➢ Katy Börner. (2006) Semantic Association Networks: Using Semantic Web Technology to Improve Scholarly Knowledge and Expertise Management. In Vladimir Geroimenko & Chaomei Chen (eds.) *Visualizing the Semantic Web*, Springer Verlag, 2nd Edition, chapter 11, pp. 183-198.

➢ Boyack, Kevin W., Klavans, R. and Börner, Katy. (2005). Mapping the Backbone of Science. *Scientometrics,* 64(3), 351-374.

➢ Hook, Peter A. and Börner, Katy. (2005) Educational Knowledge Domain Visualizations: Tools to Navigate, Understand, and Internalize the Structure of Scholarly Knowledge and Expertise. In Amanda Spink and Charles Cole (eds.) *New Directions in Cognitive Information Retrieval.* Springer-Verlag, Netherlands, chapter 5, pp. 187-208.

➢ Börner, Katy, Dall'Asta, Luca, Ke, Weimao and Vespignani, Alessandro. (April 2005) Studying the Emerging Global Brain: Analyzing and Visualizing the Impact of Co-Authorship Teams. *Complexity*, special issue on *Understanding Complex Systems*, 10(4): pp. 58 - 67. Also available as cond-mat/0502147.

➢ Ord, Terry J., Martins, Emília P., Thakur, Sidharth, Mane, Ketan K., and Börner, Katy. (2005) Trends in animal behaviour research (1968-2002): Ethoinformatics and mining library databases. *Animal Behaviour*, 69, 1399-1413. Supplementary Material.

➢ Mane, Ketan K. and Börner, Katy. (2004). Mapping Topics and Topic Bursts in PNAS. *Proceedings of the National Academy of Sciences of the United States of America*, 101(Suppl. 1):5287-5290. Also available as cond-mat/0402380.

➢ Börner, Katy, Maru, Jeegar and Goldstone, Robert. (2004). The Simultaneous Evolution of Author and Paper Networks. *Proceedings of the National Academy of Sciences of the United States of America*, 101(Suppl_1):5266-5273. Also available as cond-mat/0311459.

65

33