# Visual Analytics: Mining, Mapping, and Accelerating Science and Technology
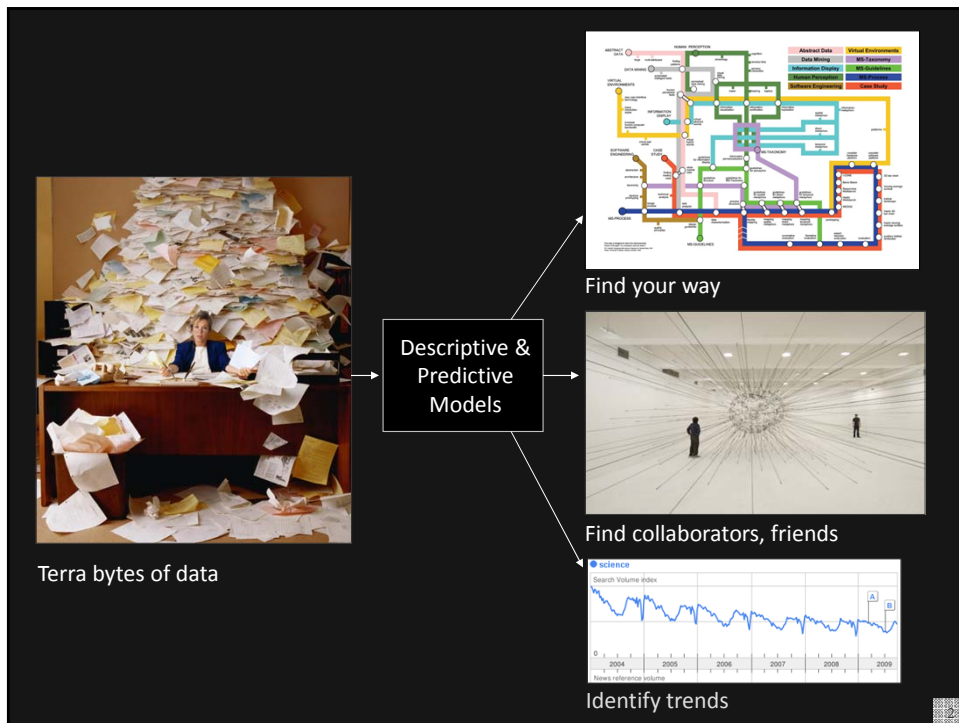
Katy Börner

Director, Cyberinfrastructure for Network Science Center
School of Informatics and Computing, Indiana University, USA

Science, Technology and Innovation Visiting Research Fellow at OECD, France

Presentation at EUROPEAN COMMISSION
DIRECTORATE-GENERAL FOR RESEARCH & INNOVATION
in Brussels, Belgium

June 23, 2014

*Language Communities of Twitter - Eric Fischer - 2012*

---

Find your way

Descriptive & Predictive Models

Find collaborators, friends

Terra bytes of data

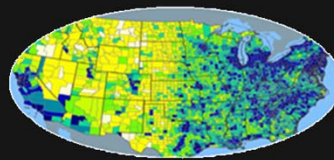Identify trends

*Descriptive Models*

*Multiple levels: Micro … Macro*

*Answering: When? Where? What? With Whom?*



**Different Levels of Abstraction/Analysis**

Macro/Global
Population Level
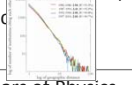
Meso/Local
Group Level

Micro
Individual Level

**Type of Analysis vs. Level of Analysis**

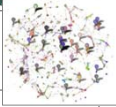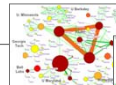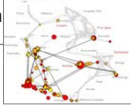|  | Micro/Individual (1-100 records) | Meso/Local (101–100,000 records) | Macro/Global (100,000 < records) |
|---|---|---|---|
| Statistical Analysis/Profiling | Individual person and their expertise profiles | Larger labs, centers, universities, research domains, or states | All of NSF, all of USA, 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 | NIH's core competency |

**Type of Analysis vs. Level of Analysis**



|  | Micro/Individual (1-100 records) | Meso/Local (101–100,000 records) | Macro/Global (100,000 < 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 | topic bursts in PNAS | 113 Years of Physics Research |
| Geospatial Analysis (Where?) | Career trajectory of one individual | intellectual landscape | PNAS publications |
| Topical Analysis (What?) | Base knowledge from which one grant draws. | research | VxOrd/Topic NIH funding |
| Network Analysis (With Whom?) | NSF Co-PI network of one | Co-author network | NIH's core competency |

**Individual Co-PI Network**
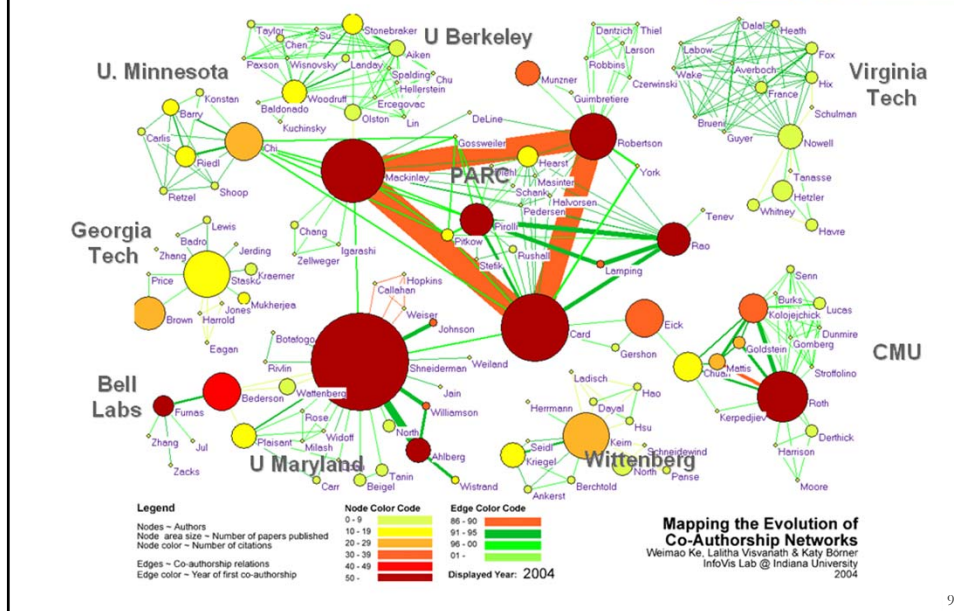*Ke & Börner (2006)*



**Mapping the Evolution of Co-Authorship Networks**
*Ke, Visvanath & Börner (2004) Won 1st price at the IEEE InfoVis Contest.*

**Mapping Transdisciplinary Tobacco Use Research Centers Publications**

Compare R01 investigator based funding with TTURC Center awards in terms of number of publications and evolving co-author networks.

*Stipelman, Hall, Zoss, Okamoto , Stokols & Börner, 2014*

*Supported by NIH/NCI Contract HHSN261200800812*

**Mapping Indiana's Intellectual Space**
*Ke, Light & Börner (2010)*

Identify
➤ Pockets of innovation
➤ Pathways from ideas to products
➤ Interplay of industry and academia

Legend:
● Academic
● Industry
--- Academic vs. Academic
--- Academic vs. Industry
--- Industry vs. Industry

11

---

**Mapping Topic Bursts**

Co-word space of the top 50 highly frequent and bursty words used in the top 10% most highly cited PNAS publications in 1982-2001.

*Mane & Börner. (2004) PNAS, 101(Suppl. 1): 5287-5290.*



**Color Code**
82 - 85
86 - 89
90 - 93
94 - 97
98 - 01

circle size ~ burst weight
cicle color ~ burst onset
ring color ~ year of max word count
years of 2nd and 3rd burst are given in color

12

*Acceleration in The Co-Development of Patented Technologies, 1996-2001 and 2006-11, see OECD Scoreboard 2013, p.171*

**Spatio-Temporal Information Production and Consumption of Major U.S. Research Institutions**

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

**Research questions:**

1. Does space still matter in the Internet age?
2. Does one still have to study and work at major research institutions in order to have access to high quality data and expertise and to produce high quality research?
3. Does the Internet lead to more global citation patterns, i.e., more citation links between papers produced at geographically distant research instructions?

**Contributions:**

➢ Answer to Qs 1 + 2 is YES.
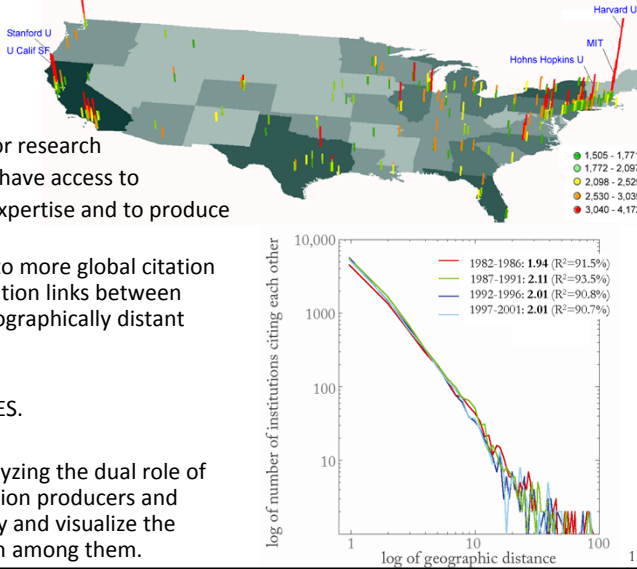➢ Answer to Qs 3 is NO.
➢ Novel approach to analyzing the dual role of institutions as information producers and consumers and to study and visualize the diffusion of information among them.



15

---

**The Global 'Scientific Food Web'**

*Mazloumian, Amin, Dirk Helbing, Sergi Lozano, Robert Light, and Katy Börner. 2013. "Global Multi-Level Analysis of the 'Scientific Food Web'". Scientific Reports 3, 1167.*
*http://cns.iu.edu/docs/publications/2013-mazloumian-food-web.pdf*

**Contributions:**

Comprehensive global analysis of scholarly knowledge production and diffusion on the level of continents, countries, and cities.

Quantifying knowledge flows between 2000 and 2009, we identify global sources and sinks of knowledge production. Our knowledge flow index reveals, where ideas are born and consumed, thereby defining a global 'scientific food web'.

**While Asia is quickly catching up in terms of publications and citation rates, we find that its dependence on knowledge consumption has further increased.**



16

**Figure 2 | World map of the greatest knowledge sources and sinks, based on our scientific fitness index.** Green bars indicate that the number of citations received is over-proportional, red that the number of citations received is lower than expected (according to a homogeneous distribution of citations over all cities that have published more than 500 papers). It can be seen that most scientific activity occurs in the temperate zone. Moreover, areas of high fitness tend to be areas that are performing economically well (but the opposite does not hold).



*Country Mobility Network, 1996-2011, see OECD Scoreboard 2013, p. 62*

**Knowledge networks**

The mobility of researchers contributes to the diffusion of scientific and technological knowledge across institutions, at a national and international level. The trail of affiliation changes left by scientific authors in their scholarly publication records provides a partial means of identifying the international network of researcher flows. As expected, leading research countries tend to attract more scientific authors from abroad than they have authors who leave. Flows within each pair of countries tend to be of a similar order of magnitude in both directions, suggesting the existence of complex patterns of knowledge circulation representing the mobility of individuals at different stages of their careers, from students to established professors. The international mobility network also displays a number of interesting patterns that reveal affinities between different economies based on linguistic, historical as well as political and cultural linkages, such as the link between Spain and Latin America countries.

**57. International mobility network, 1996-2011**
*Counts of bilateral flows, by first and last affiliation*

Source: OECD calculations based on Scopus Custom Data, Elsevier, version 5.2012, June 2013.

18

9

**Type of Analysis vs. Level of Analysis**

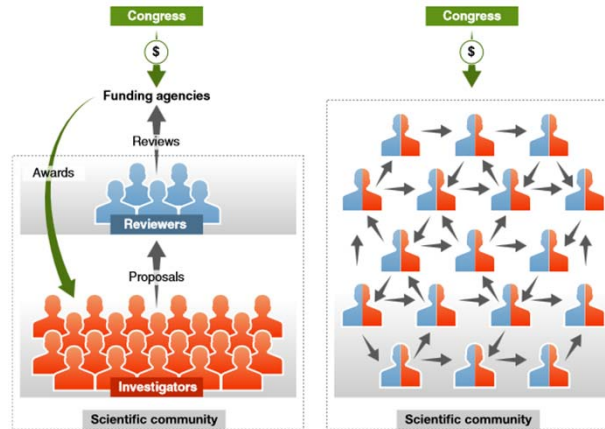| | Micro/Individual (1-100 records) | Meso/Local (101–100,000 records) | Macro/Global (100,000 < records) |
|---|---|---|---|
| Statistical Analysis/Profiling | Individual person and their expertise profiles | Larger labs, centers, universities, research domains, or states | All of NSF, all of USA, all of s... |
| Temporal Analysis (When?) | Funding portfolio of one individual | ...ic bursts in ...PNAS | 113 Years of Physics Research |
| Geospatial Analysis (Where?) | Career trajectory of one individual | ...intellectual la... | PNAS publications |
| Topical Analysis (What?) | Base knowledge from ...s. | ...research | VxOrd/Topic NIH funding |
| Network Analysis (With Whom?) | NSF Co-PI network of one | Co-author network | NIH's core com... |

19

---

*Predictive Models (Why?)*

Example: Collective allocation of science funding as an alternative to peer review

20

**From funding agencies to scientific agency: Collective allocation of science funding as an alternative to peer review**

*Bollen, Crandall, Junk, Ding & Börner. 2014. EMBO Reports 15 (1): 1-121.*



Existing (left) and proposed (right) funding systems. Reviewers in blue; investigators in red.

In the proposed system, all scientists are both investigators and reviewers: every scientist receives a fixed amount of funding from the government and discretionary distributions from other scientists, but each is required in turn to redistribute some fraction of the total they received to other investigators.

21

---

**From funding agencies to scientific agency: Collective allocation of science funding as an alternative to peer review**

*Bollen, Johan, David Crandall, Damion Junk, Ying Ding & Katy Börner. 2014. EMBO Reports 15 (1): 1-121.*
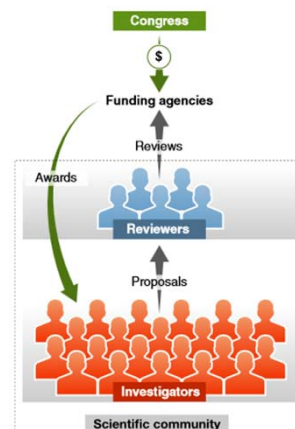
**Current Model is Expensive:**

If four professors work four weeks full-time on a proposal submission, labor costs are about $30k [1]. With typical funding rates below 20%, about five submission-review cycles might be needed resulting in a total expected labor cost of **$150k**. The average NSF grant is **$128k** per year.

U.S. universities charge about 50% overhead (ca. $42k), leaving about **$86k**.

In other words, the four professors lose **$150k-$86k= - $64k** of paid research time by **obtaining** a grant to perform the proposed research.

**To add:** Time spent by researchers to review proposals. In 2012 alone, NSF convened more than 17,000 scientists to review 53,556 proposals.

[1] Taulbee Survey of Salaries Computer Science ,
   *http://cra.org/resources/taulbee*



22

**From funding agencies to scientific agency: Collective allocation of science funding as an alternative to peer review**

*Bollen, Crandall, Junk, Ding & Börner. 2014. EMBO Reports 15 (1): 1-121.*

**Assume**

Total funding budget in year y is $t_y$
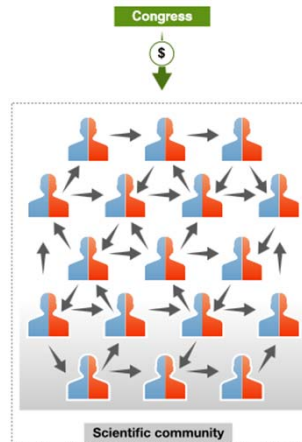
Number of qualified scientists is $n$

**Each year,**

the funding agency deposits a fixed amount into each account, equal to the total funding budget divided by the total number of scientists: $t_y/n$.

Each scientist must distribute a fixed fraction, e.g., 50%, of received funding to other scientists (no self-funding, COIs respected).

**Result**

Scientists collectively assess each others' merit based on different criteria; they "fund-rank" scientists; highly ranked scientists have to distribute more money.

23

---

**From funding agencies to scientific agency: Collective allocation of science funding as an alternative to peer review**
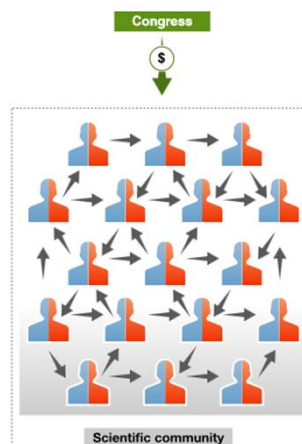
*Bollen, Crandall, Junk, Ding & Börner. 2014. EMBO Reports 15 (1): 1-121.*

**Example:**

Total funding budget per year is 2012 NSF budget

Given the number of NSF funded scientists, each receives a $100,000 basic grant.

Fraction is set to 50%

In 2013, scientist *S* receives a basic grant of $100,000 plus $200,000 from her peers, i.e., a total of $300,000.

In 2013, *S* can spend 50% of that total sum, $150,000, on her own research program, but must donate 50% to other scientists for their 2014 budget.

Rather than submitting and reviewing project proposals, *S* donates directly to other scientists by logging into a centralized website and entering the names of the scientists to donate to and how much each should receive.

24

**From funding agencies to scientific agency: Collective allocation of science funding as an alternative to peer review**
*Bollen, Crandall, Junk, Ding & Börner. 2014. EMBO Reports 15 (1): 1-121.*

**Model Run and Validation:**

Model is presented in http://arxiv.org/abs/1304.1067

It uses **citations as a proxy** for how each scientist might distribute funds in the proposed system.

Dataset: 37M articles from TR 1992 to 2010 Web of Science (WoS) database with **770M citations** and 4,195,734 unique author names. The **867,872 names** who had authored at least one paper per year in any five years of the period 2000–2010 were used in validation.

For each pair of authors we determined the number of times one had cited the other in each year of our citation data (1992–2010).

NIH and NSF funding records from IU's Scholarly Database provided 347,364 grant amounts for 109,919 unique scientists for that time period.

Simulation run begins in year 2000, in which every scientist was given a fixed budget of B = $100k. In subsequent years, scientists distribute their funding in proportion to their citations over the prior 5 years.

The model yields funding patterns similar to existing NIH and NSF distributions.

25

---

NEWSFOCUS

**Making Every Scientist a Research Funder**

When it comes to using peer review to distribute research dollars, Johan Bollen favors radical simplicity.

Over the years, many scientists have suggested that the current system could be improved by changing the composition of the review panels, tweaking the interactions among reviewers, or revising how the proposals are scored. But Bollen, a computer scientist at Indiana University, Bloomington, would simply award all eligible researchers a block grant—and then require them to give some of it away to colleagues they judge most deserving.

That radical step, described in a paper Bollen and four Indiana colleagues recently posted on *EMBO Reports*, retains peer review's core concept of tapping into the views of the most knowledgeable researchers. But it would eliminate the huge investment in time and money required to submit proposals and assemble panels to judge them.

Bollen's process would be almost instantaneous: In a version of expert-directed crowdsourcing, scientists would fill out a form once a year listing their favored researchers, and a predetermined portion of their annual grant money—a total of, say, 50%—would then be transferred to their choices.

"So many scientists spend so much time on peer review, and there's a high level of frustration," Bollen explains. "We already know who the best people are. And if you're doing good work, then you deserve to receive support."

Others are skeptical. "I've known Johan for a long time and have the highest regard for his ability as an out-of-the-box thinker," says Stephen Griffin, a retired National Science Foundation (NSF) program manager who's now a visiting professor of information sciences at the University of Pittsburgh in Pennsylvania. "But there are a number of issues he doesn't address."

Those sticking points include the likely mismatch between what researchers need and what their colleagues give them; the absence of any replacement for the overhead payments in today's grants, which support infrastructure at host institutions; and the dearth of public accountability for the billions of dollars that would flow from public coffers to individuals. "Scientists aren't really equipped to be a funding agency," Griffin notes.

Bollen acknowledges that the process would need safeguards to ensure that scientists don't reward their friends or punish their enemies. But his analysis suggests that the U.S. research landscape would not look all that different if his radical proposal were adopted.

Drawing upon citation data in 37 million papers over 20 years, the Indiana researchers conducted a simulation premised on the idea that scientists would reallocate their federal dollars according to how often they cited their peers. The simulation, he says, yielded a funding pattern "similar in shape to the actual distribution" at NSF and the National Institutes of Health for the past decade—at a fraction of the overhead required by the current system.
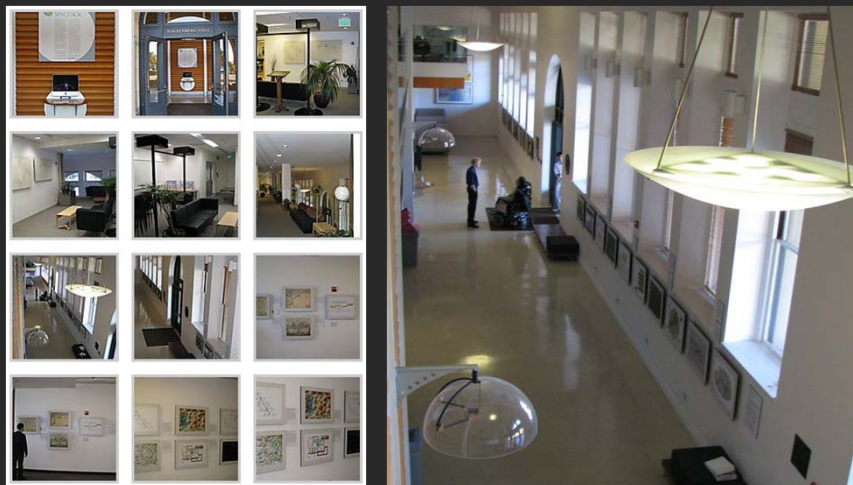
–JDM

February 7, 2014

*Science 7 February 2014: Vol. 343 no. 6171 p. 598*
*DOI: 10.1126/science.343.6171.598*
http://www.sciencemag.org/content/343/6171/598.full?sid=4f40a7f0-6ba2-4ad8-a181-7ab394fe2178
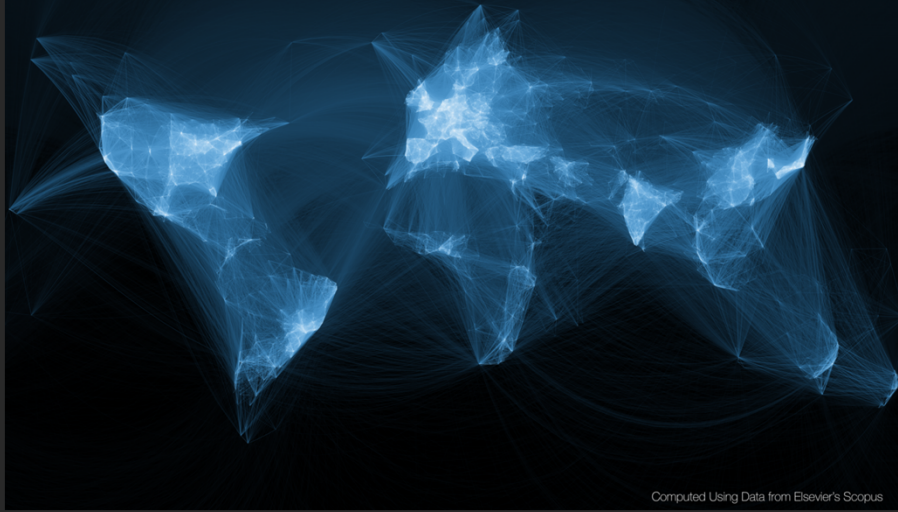
***Visualizing STI Model Results***

Example: Places & Spaces: Mapping Science Exhibit



Mapping Science Exhibit on display at MEDIA X, Stanford University
*http://mediax.stanford.edu*, *http://scaleindependentthought.typepad.com/photos/scimaps*

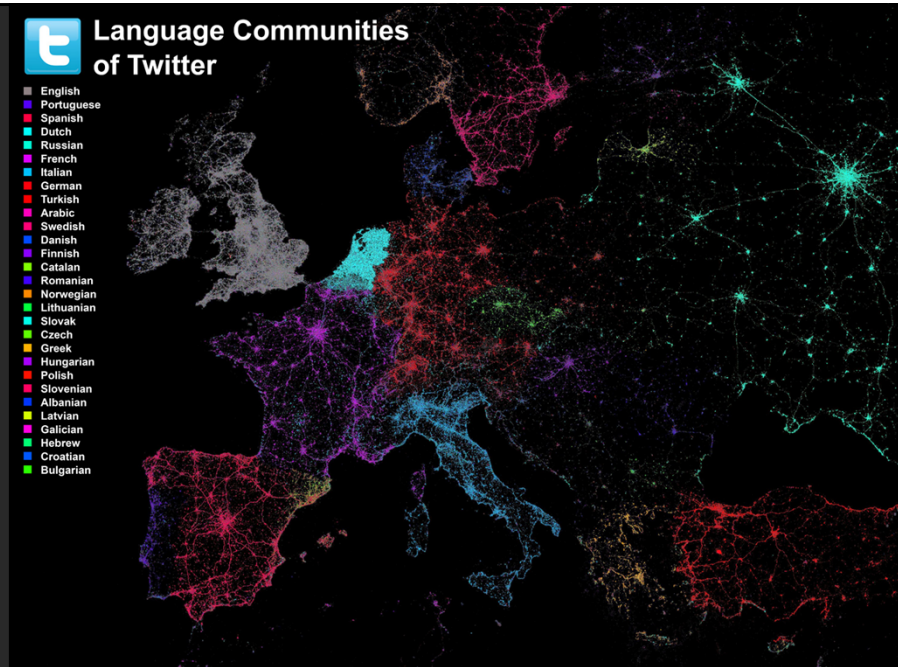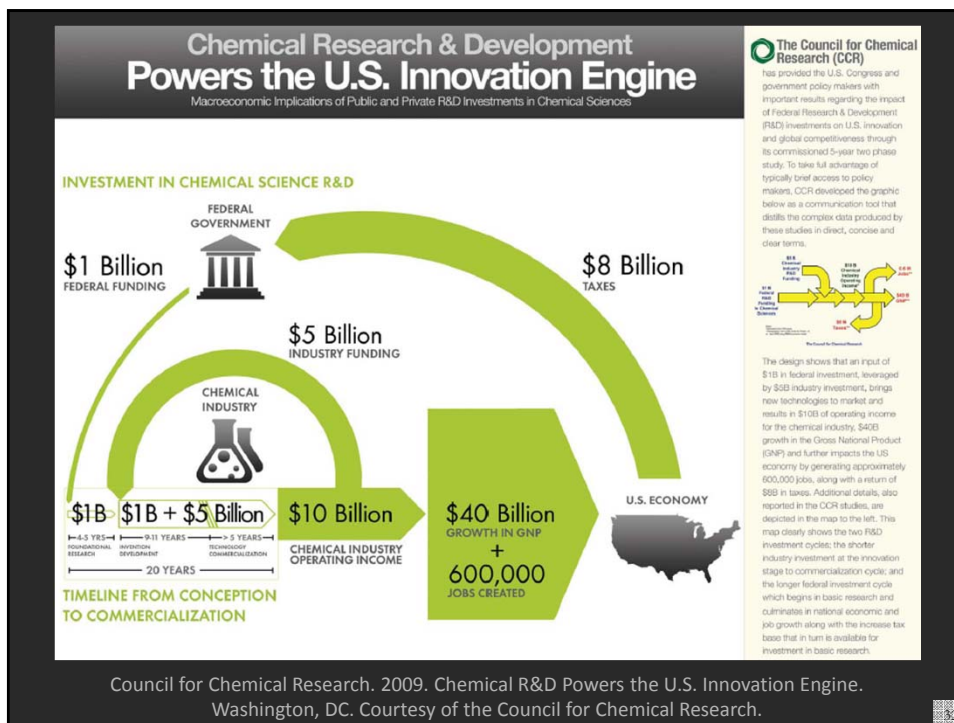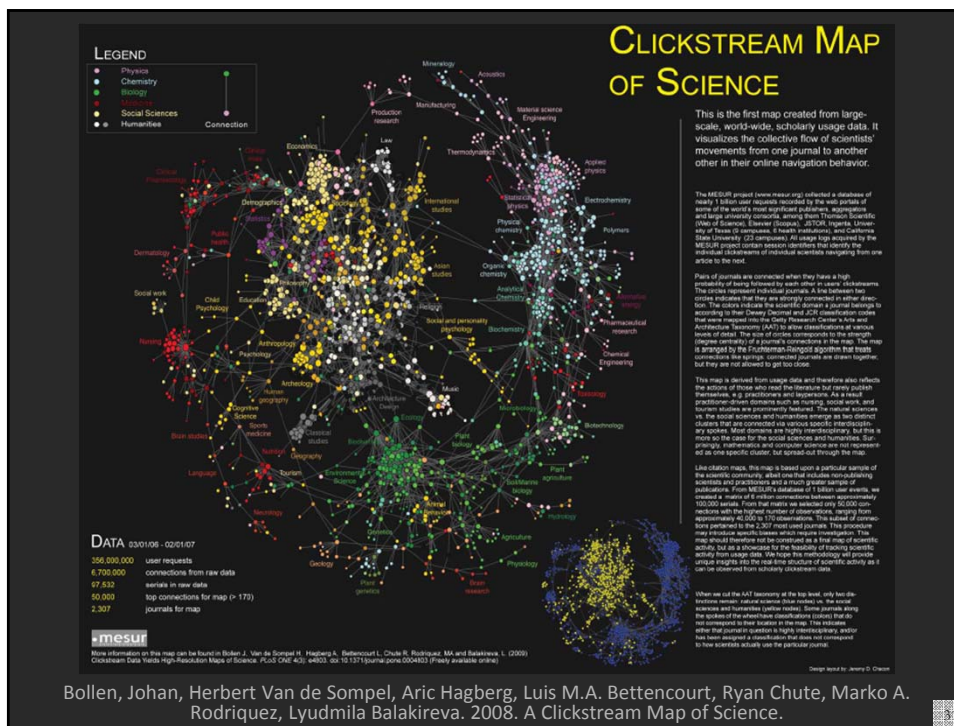Map of Scientific Collaborations from 2005-2009

Computed Using Data from Elsevier's Scopus

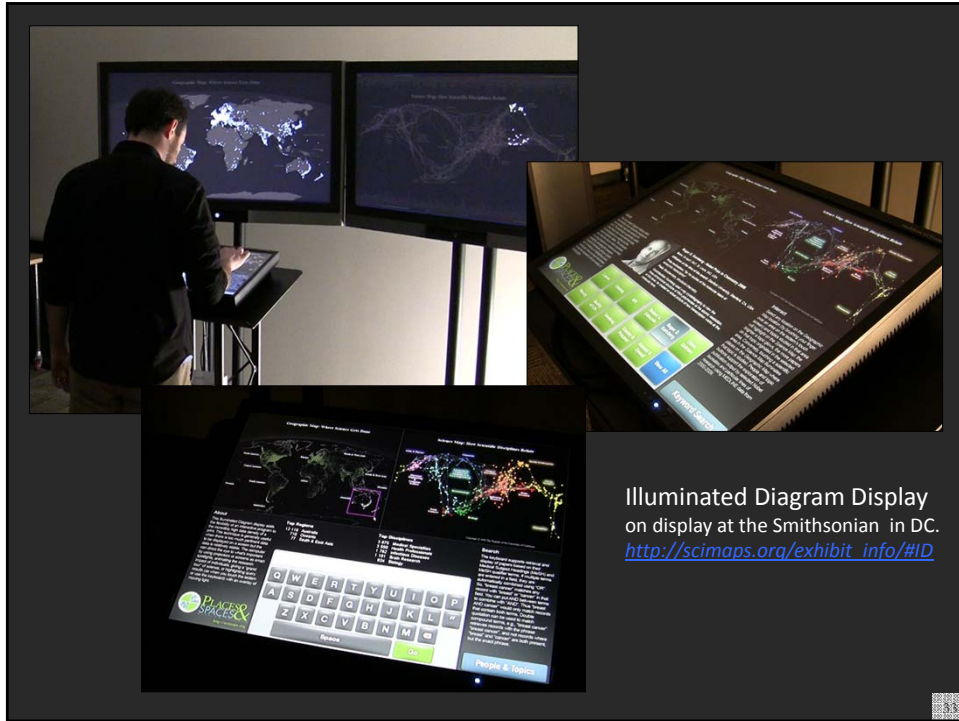Olivier H. Beauchesne, 2011. Map of Scientific Collaborations from 2005-2009.



Language Communities of Twitter - Eric Fischer - 2012

15

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



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

Illuminated Diagram Display
on display at the Smithsonian in DC.
http://scimaps.org/exhibit_info/#ID



Geographic Map: Where Science Gets Done

Science Map: How Scientific Disciplines Relate

**About**
This Illuminated Diagram display adds the flexibility of an interactive program to the incredibly high data density of a print. This technique is generally useful when there is too much pertinent data to be displayed on a screen but the data is relatively stable. The computer can direct the eye to what's important by using projectors or screens as smart spotlights, animating the research impact of individuals, giving a "grand tour" of science, or highlighting query results (as when you touch the lectern or use the keyboard) with an overlay of moving light.

**Top Five Continents**
North America - 4,000 records
South & East Asia - 3,589
Australia - 2,431
Africa - 2,208
South America - 1,562

**Top Five Scientific Disciplines**
Math & Physics - 4,000 records
Health Professionals - 3,589
Social Sciences - 2,431
Aeronautical, Chemical, Mechanical & Civil Engineering - 2,208
Humanities - 1,562

**Search**
The keyboard supports retrieval and display of papers based on their Medical Subject Headings (MeSH) and MeSH qualifier terms. If multiple terms are entered in a field, they are automatically combined using "OR". So, "breast cancer" matches any record with "breast" or "cancer" in that field. You can put AND between terms to combine with "AND". Thus "breast AND cancer" would only match records that contain both terms. Double quotation can be used to match compound terms, e.g., "breast cancer" retrieves records with the phrase "breast cancer", and not records where "breast" and "cancer" are both present, but the exact phrase.

Science Maps in "Expedition Zukunft" science train visiting 62 cities in 7 months 12 coaches, 300 m long Opening was on April 23rd, 2009 by German Chancellor Merkel
*http://www.expedition-zukunft.de*

Places & Spaces *Digital Display in North Carolina State's brand new Immersion Theater*

37

## Places & Spaces: Mapping Science Exhibit
http://scimaps.org



Maps are available for sale and the exhibit can be hosted by anyone.

38

***Visualizing STI Model Results***

Example: The Information Visualization MOOC



Register for free at http://ivmooc.cns.iu.edu. Class will restart in January 2015.

The Information Visualization MOOC
ivmooc.cns.iu.edu

Exterior Color (Linear)
count

1    269    537

Area (Linear)
count

537
269
1

Students from more than 100 countries
350+ faculty members
#ivmooc

---

# Course Schedule

- **Session 1** – Workflow design and visualization framework
- **Session 2** – "When:" Temporal Data
- **Session 3** – "Where:" Geospatial Data
- **Session 4** – "What:" Topical Data

**Mid-Term**

**Students work in teams with clients.**

- **Session 5** – "With Whom:" Trees
- **Session 6** – "With Whom:" Networks
- **Session 7** – Dynamic Visualizations and Deployment

**Final Exam**

Final grade is based on Midterm (**30%**), Final (**40%**), Client Project (**30%**).

42

# Clients



http://ivmooc.cns.iu.edu/clients.html

National Science Foundation Awards to Institutions Conducting "Global Warming" Research from 1979-2009

mjstamper_ivmooc

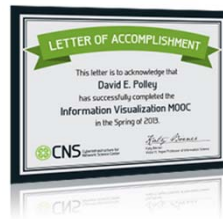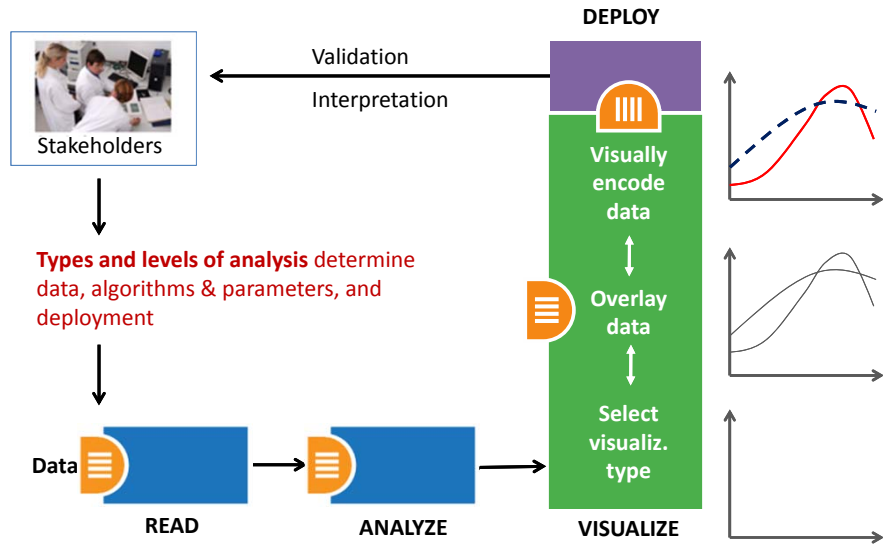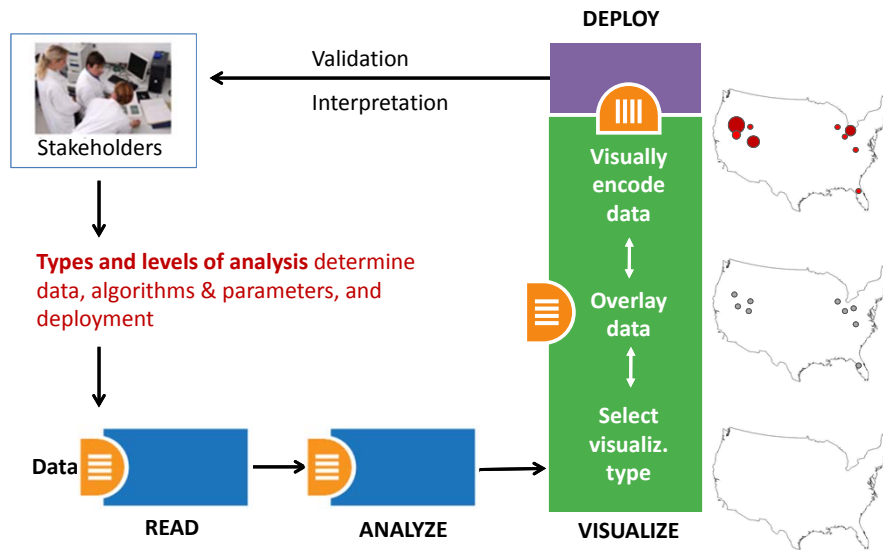# Tutorial at OECD: "Open Source Tools for S&T Data Analysis and Visualization"

*Speaker*:  Katy Börner
*Date*:      25 June, 2014
*Time*:      9:30-13:30
*Venue*:    OECD Conference Centre, Room E, 2 rue André Pascal, Paris 75116

This tutorial is designed for researchers and practitioners interested to use advanced data mining algorithms and visualizations in their research and daily decision making. It gives an overview of open source tools for the analysis and visualization of science and technology (S&T) data. A specific focus is the Science of Science (Sci2) Tool that supports temporal, geospatial, topical, and network analysis and visualization of scholarly datasets at the micro (individual), meso (local), and macro (global) levels. Open data from OECD and other government agencies will be used to demonstrate different analysis and visualization workflows.

The tutorial has two parts: Part 1 provides an overview of diverse international efforts to (1) standardize and federate micro-level datasets of S&T activity, e.g., publication, patent, grant, social media data; (2) design open code tools and online services that are interoperable; (3) develop means to share and teach open datasets and tools. Part 2 is reserved for "hands-on" training. If you plan to attend this part, please bring your laptop and pre-install the Sci2 (v 1.1 beta) tool **prior to the workshop**—the tool is freely available at http://sci2.cns.iu.edu.

For a preview see: http://cns.iu.edu/docs/presentations/2012-borner-sci2tutorial-oecd.pdf
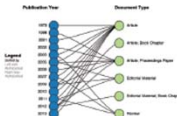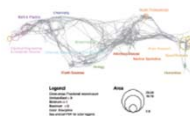
## Slide 49

### Load **One** File and Run **Many** Analyses and Visualizations

| Times Cited | Publication Year | City of Publisher | Country | Journal Title (Full) | Title | Subject Category | Authors |
|---|---|---|---|---|---|---|---|
| 12 | 2011 | NEW YORK | USA | COMMUNICATIONS OF THE ACM | Plug-and-Play Macroscopes | Computer Science | Borner, K |
| 18 | 2010 | MALDEN | USA | CTS-CLINICAL AND TRANSLATIONAL SCIENCE | Advancing the Science of Team Science | Research & Experimental Medicine | Falk-Krzesinski, HJ\|Borner, K\|Contractor, N\|Fiore, SM\|Hall, KL\|Keyton, J\|Spring, B\|Stokols, D\|Trochim, W\|Uzzi, B |
| 13 | 2010 | WASHINGTON | USA | SCIENCE TRANSLATIONAL MEDICINE | A Multi-Level Systems Perspective for the Science of Team Science | Cell Biology \|Research & Experimental Medicine | Borner, K\|Contractor, N\|Falk-Krzesinski, HJ\|Fiore, SM\|Hall, KL\|Keyton, J\|Spring, B\|Stokols, D\|Trochim, W\|Uzzi, B |



Statistical Analysis–p. 44

| Location | Count | # Citations |
|---|---|---|
| Netherlands | 13 | 292 |
| United States | 9 | 318 |
| Germany | 11 | 36 |
| United Kingdom | 1 | 2 |

Temporal Burst Analysis–p. 48

Geospatial Analysis–p. 52

Geospatial Analysis–p. 52

49

## Slide 50

### Load **One** File and Run **Many** Analyses and Visualizations

Topical Analysis–p. 56

Paper Citation Network–p. 60

Bi-Modal Network–p. 60

Co-author and many other bi-modal networks.

50

# References

Börner, Katy, Chen, Chaomei, and Boyack, Kevin. (2003). **Visualizing Knowledge Domains.** In Blaise Cronin (Ed.), *ARIST*, Medford, NJ: Information Today, 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.), *ARIST*, Information Today, Inc., Volume 41, Chapter 12, pp. 537-607.
http://ivl.slis.indiana.edu/km/pub/2007-borner-arist.pdf

Börner, Katy (2010) **Atlas of Science**. MIT Press.
http://scimaps.org/atlas

Scharnhorst, Andrea, Börner, Katy, van den Besselaar, Peter (2012) **Models of Science Dynamics**. Springer Verlag.

Katy Börner, Michael Conlon, Jon Corson-Rikert, Cornell, Ying Ding (2012) **VIVO: A Semantic Approach to Scholarly Networking and Discovery**. Morgan & Claypool.

Katy Börner and David E Polley (2014) **Visual Insights: A Practical Guide to Making Sense of Data**. MIT Press.



51

---



All papers, maps, tools, talks, press are linked from http://cns.iu.edu
These slides will soon be at http://cns.iu.edu/docs/presentations

CNS Facebook: http://www.facebook.com/cnscenter
Mapping Science Exhibit Facebook: http://www.facebook.com/mappingscience

52