

# BIG Data for BIG Science of Science Research: The Value of CADRE

Katy Börner @katycns

Victor H. Yngve Distinguished Professor of  
Intelligent Systems Engineering & Information Science  
Director, Cyberinfrastructure for Network Science Center  
School of Informatics, Computing, and Engineering  
Indiana University, Bloomington, IN, USA  
+ 2018 Humboldt Fellow, TU Dresden, Germany



*CADRE Workshop at ISSI 2019*

Rome, Italy | September 2, 2019



# BIG Data



# Datasets used in SoS R&D

2009

2019

Table 1. Data sets and their properties (\* 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

WoS\*

Scopus

Google Scholar

MS Academic Graph

UMETRICS

Social Media Data

News Data

IoT Data

LaRowe, Gavin, Sumeet Adinath Ambre, John W. Burgoon, Weimao Ke, and Katy Börner. 2009. "[The Scholarly Database and Its Utility for Scientometrics Research](#)". *Scientometrics* 79 (2): 219-234.

\* <http://iuni.iu.edu/resources/cadre>

# Web of Science as a Research Dataset

November 14, 2016 - November 15, 2016 | Bloomington, Indiana

Organizers:

<https://cns.iu.edu/workshops/event/161114.html>



## Katy Börner

Victor H. Yngve Distinguished Professor of Information Science  
Director, CI for Network Science Center Curator, Mapping Science Exhibit  
School of Informatics and Computing  
Indiana University, Bloomington  
katy@indiana.edu



## Eamon Duede

Executive Director  
Knowledge Lab  
University of Chicago  
eduede@uchicago.edu



## James Pringle

Head of Industry Development  
Clarivate Analytics





## Session 1: Web of Science “Outside the Box”

Facilitator: Katy Börner

The Web of Science and similar metadata datasets are housed, maintained, and enhanced in local institutional enclaves. This the data (nam

- Matt Hu
- Yadu Ba
- Orion P
- Nicholas
- Vetle To
- Lee Gile

## Session 3: Understanding Web of Science as Research Data

Facilitator: Jason Rollins

For over 50 years, the Web of Science evolved as a dataset in response to changing research contexts and priorities. Today, more researchers are using the Web of Science “at scale” to ask and answer powerful new questions about the shape, dynamics, and veracity of science and scholarship. The Web of Science now appears both an object of inquiry its own right and a vast sensor network for discerning large-scale trends. What is changing in this dataset to support these new uses, and what could change further? Presentations and discussion led by Clarivate Analytics team.

- Jim Pringle: “**WoS Metadata as Research Data**”
- Ted Lawless: “**Web of Science Data Integration**”
- Linge Bai: “**Data Unification and Disambiguation: Institutions and Authors**”
- Sebastien Brien: “Clarivate Analytics in the Cloud: Architecture and Analytics”

*Break*

## Session 4: Hackathon Breakout Sessions

Facilitators: Eamon Duede, Jason Rollins, and Ted Lawless

A mix of sessions determined by 3-4 “big questions” prioritized on Day 1, grouped as:

- Technical Hackathon(s)*: Practical Focus on applying code across research centers in such areas as data disambiguation (names, institutions, geolocations), linking WoS data to other datasets, building models to predict gender, ethnicity, etc.
- Topical Hackathon(s)*: Working across research centers on Authorship & Collaboration; Gender in Science; Topic Modeling and/or other topics defined by attendee interest.
- Community Hackathon*: Focus on establishing an ongoing community (e.g. setting up an enclave, tools & mechanisms for sharing code, citing and acknowledging contributions, and/or what is appropriate for cross-enclave sharing).

# Reproducible Scientometrics Research: Open Data, Code, and Education

## Date

October 17, 2017

## Meeting Place

ISSI 2017, Wuhan University

Wuhan, China

## Session Organizers

Sybille Hinze  
DZHW  
Berlin, Germany

Jesper Schneider  
Aarhus University  
Denmark

Katy Börner  
Indiana University  
Bloomington, IN, USA  
[Slides](#) | [MP4](#)

Jason Rollins  
Clarivate Analytics  
San Francisco, CA, USA

Theresa Velden  
ZTG TU Berlin  
Germany

Andrea Scharnhorst  
KNAW, Amsterdam  
The Netherlands

Jesper Schneider  
Aarhus University  
Denmark

Ludo Waltman  
CWTS, University of Leiden  
The Netherlands

## Workshop Agenda

1. Introduction (Sybille Hinze & Theresa Velden)
2. Reproducibility in Scientometrics: Data Enclaves, Open Code, and Open Education (Katy Börner,
3. Reproducibility in Scientometrics through Quality Assurance (Sybille Hinze)
4. A Vendor's View on Reproducibility — Datasets, Tools, & Partnerships (Jason Rollins)
5. Reproducibility in Scientometrics — A Journal Editor's Perspective (Ludo Waltman)
6. Reproducibility — Principles and Challenges (Jesper Schneider)
7. Reproducibility & the Performativity of Methods (Theresa Velden)

<https://cns.iu.edu/workshops/event/171017.html>



# Web of Science™ as a Research Dataset

Katy Börner,<sup>1</sup> Valentin Pentchev,<sup>2</sup> Matthew Hutchinson,<sup>2</sup> James Pringle,<sup>3</sup> Jason Rollins,<sup>3</sup> Yadu N. Babuji,<sup>4</sup> & Eamon Duede<sup>5</sup>

<sup>1</sup>katy@indiana.edu CNS, SOIC & Network Science Institute, Indiana University, Bloomington, IN, US  
<sup>2</sup>mahutch@iu.edu & vpentche@iu.edu IUNI, Indiana University, Bloomington, USA

<sup>3</sup>jason.rollins@clarivate.com & james.pringle@clarivate.com Clarivate Analytics, USA  
<sup>4</sup>yadunand@uchicago.edu Computation Institute, Knowledge Lab, UofChicago, USA



## Introduction

The Clarivate Analytics Web of Science (WoS) has served as a research dataset for more than 9,000 scholarly articles in the past 15 years alone—across a wide a range of fields and disciplines from toxicology to computer science to economics. Scientists and scholars have been particularly interested in the WoS citation network, a massive graph containing billions of links that can proxy the structure and dynamics of not only scholarly communication, but knowledge diffusion, the evolution of fields, and the career lifecycles of individuals and institutions. To power these investigations, scholars are increasingly employing a number of compute-intensive methodologies, sophisticated big data infrastructures, and so called collaborative “discovery science” tools and techniques. Suddenly, in addition to deep, domain specific expertise, world-class computational knowhow appears to be a new prerequisite for analysis of scholarly data at the scale represented by WoS. While cloud-based computing and tools are more prevalent and accessible than ever before, harnessing these technologies remains both a challenge and opportunity for researchers and data providers (i.e., Clarivate Analytics and similar commercial data vendors and non-commercial aggregators). While the opportunities made possible by scholarly data at the size and scope of WoS for discovery and innovation are limited only by imagination, two general prospects come readily to mind. First, access to these data coupled with the appropriate computational and analytical capabilities opens up a wide range of funding and subsequent publishing opportunities in high impact venues. Second, data providers can pursue new business opportunities, including novel data access models, new types of analytic products, and new kinds of academic/industry partnerships. In this poster paper, we briefly explore 1) the new computational infrastructures that are being developed to enable collaborative research that leverages scholarly datasets such as WoS that are both big and proprietary; 2) some recent findings that have been made possible by these infrastructures; and, 3) new commercial offerings that have been enabled and demanded in response to increasing reliance on the WoS as a research dataset.

## New Computational Infrastructures

Research leveraging big, scholarly datasets like WoS presents researchers with challenges related to the data's size, inherently relational format, and sensitive (proprietary) nature. To overcome these challenges, researchers have developed a new generation of enclave supported, high performance, and cloud-based, collaborative research environments that are both elastic enough to provide substantial computational resources when needed while remaining secure enough to protect data providers'

## IU International Co-Affiliation Network, 2004-2013

CNS @ Indiana University  
2016

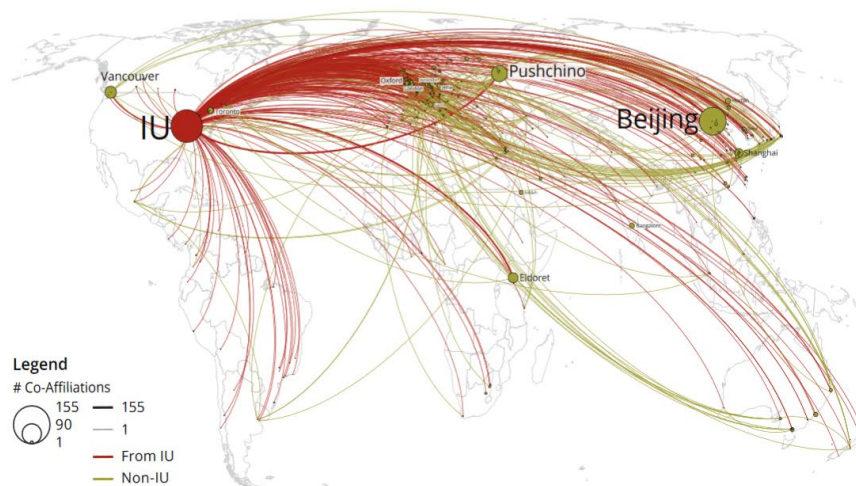


Figure 1. IU Co-Affiliation Network

commercial interests. Moreover, these environments have evolved to allow geographically distributed researchers to collaborate on research projects in the fast paced, iterative style that has come to dominate research in the era of Big Data—namely, “discovery science”.

## IUNI WoS Data Enclave

The Indiana University Network Science Institute (IUNI) acquired the complete set of Clarivate Analytics' Web of Science XML raw data (Web of Knowledge version 5). The data was parsed and stored in a well-documented PostgreSQL database, see entity-relationship diagram, database schema, and data dictionary on <http://iuni.iu.edu/resources/web-of-science>. The code used to parse the WoS XML format and to save data in the PostgreSQL database was made available freely on GitHub, see [https://github.com/iuni/CNS/generic\\_parser](https://github.com/iuni/CNS/generic_parser). All data can be accessed via the IUNI WoS Data Enclave, a secure repository that uses IU's Karst high-throughput computing cluster designed to handle large amounts of processing capacity over long periods of time. Access to the XML data and the PostgreSQL database is granted to a user's Karst account. IU faculty, staff, and qualifying sponsored affiliates can request accounts on Karst to use the data for academic research and without any sharing of data. A simple web browser based query interface to the WoS dataset was implemented to support custom queries for specific terms, journals, or authors. Datasets can be downloaded in CSV data format compatible with data mining and visualization tools such as Gephi or the Sci2 Tool (<http://sci2.cns.iu.edu/>) (Sci2 Team, 2009). More about the IUNI WoS Data Enclave can be found at <http://iuni.iu.edu/resources/web-of-science>.

## Cloud Kotta

One platform specifically developed with WoS in mind is Knowledge Lab's Cloud Kotta (CK). CK is a secure data enclave and analytics platform that serves the research needs of social sciences (Babuji 2016). By hosting CK in the Amazon Web Services cloud, the developers were able to take advantage of virtually limitless compute, cost-effective storage and the ability to implement a fine-grained security model ensuring the authorized collaborators could access both data and compute resources from any where in the world (Babuji 2016). Moreover, CK supports multiuser, rapid ideation and research iteration through a novel Python library that enables specific functions in an analysis code, written in a Jupyter Notebook to be seamlessly and securely submitted to the CK executor (Babuji 2017). By allowing researchers to develop and share analysis code interactively over secure data like WoS, CK has removed the need for deep computational infrastructure expertise. The complete WoS XML dataset was ingested into a relational database housed in CK using a custom parser that has been made freely available on GitHub (see: [https://github.com/alexander-belikov/wos\\_parser](https://github.com/alexander-belikov/wos_parser)). The Cloud Kotta WoS database schema can be found on CK's documentation pages (See: <http://docs.cloudkotta.org/dataguide/wos.html>). More about Cloud Kotta can be found at <http://docs.cloudkotta.org>.

## New Computational Infrastructures

### Fostering Global Collaboration

Among others, IU started to use the IUNI WoS data to understand existing and foster global research collaborations. The world map in Figure 1 shows the co-affiliations of authors that listed “Indiana University” and at least one other non-U.S. institution as affiliation on 1,590 scholarly papers published in 2004-2013. There are 344 affiliation locations (not counting IU) and 641 co-affiliation links. Nodes denote author locations and are area sized coded by degree with the exception of IU, which has 1,592 co-affiliation links. Links denote co-affiliations, e.g., an author with three affiliations IU, X, Y will add three links; the two links that connect IU with X and Y are drawn in red while the link between X and Y is given in green. Links are size coded by the number of co-affiliations with the top-three being Beijing, China (155), Eldoret, Kenya (115), and Pushchino, Russia (90).

### Impact vs. Disruptiveness

Researchers at the University of Chicago's Knowledge Lab and Northwestern University's NICO have used WoS data going back to 1900 to study the relationship between team size and impact and the relationship between team size and disruptiveness. This work, currently under review, finds striking differences between the scientific output of large and small teams. Looking across all fields represented in WoS, small teams are shown to disrupt science, patents, and software with new ideas and opportunities, while large teams contribute to existing ones. Figure 2 shows the relationship between impact and disruptiveness of articles (left panel), patents (middle), and software (right). In all three spaces, there is a strong, inverse relationship between citations and disruptiveness as team size increases.

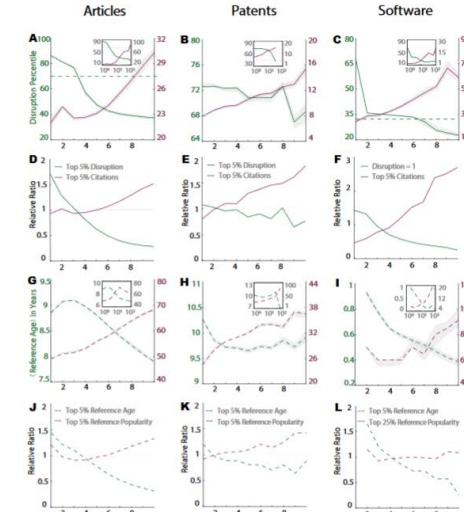


Figure 2. Tracing Inventive Teams

## New Commercial Offerings

The value of the Web of Science as a search and discovery tool is well established at thousands of research institutions worldwide. But the commercial opportunities for the use of its high-quality metadata outside of the platform for big data studies are still emerging. When researchers need to study broad-scale trends in science, technology, and innovation, they very often turn to the Web of Science as the most comprehensive citation source to provide over 100 years of consistent, global publication data. Increasingly, user requests for this data take the form of custom reports, curated data sub-sets, and large-scale raw XML delivery. Clarivate Analytics is actively looking at compelling ways to meet these customer demands with new commercial products and data delivery choices. These options must balance scale and ease of use, with security and control over access to the proprietary WoS data. The lessons learned in the development of Cloud Kotta and IUNI WoS Data Enclave will very likely be instructive here, as they have proven their utility and leverage a mix of custom code built on proven commercial cloud services. Both self-service data access and secure use of analytical tools in a cloud “sandbox” seem like attractive features of these environments that could make commercial sense to meet the evolving expectation of Web of Science customers.

## Acknowledgements

This work is partially supported by and contributes to research for IBM, Facebook, Jump Trading, AWS, Clarivate Analytics, the National Institutes of Health under awards P01 AG039347 and U01 CA198934 and the National Science Foundation under awards NCSF-1538763, EAGER 1566393, and NCCP Supplement 1553044. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation. This work uses Web of Science data by Clarivate Analytics.

## References

- Babuji, Y. N., Chard, K., Gerow, A., & Duede, E. (2016). Cloud Kotta: Enabling Secure and Scalable Data Analytics in the Cloud. *IEEE Big Data 2016*.
- Babuji, Y. N., Chard, K., Gerow, A., & Duede, E. (2016). A Secure Data Enclave and Analytics Platform for Social Scientists. *IEEE eScience 2016*.
- Babuji, Y. N., Chard, K., & Duede, E. (2017). Enabling Interactive Analytics of Secure Data using Cloud Kotta. *Science Cloud Workshop: ACM International Symposium on High-Performance Parallel and Distributed Computing 2017* (Forthcoming)
- Sci2 Team. (2009). Science of Science (Sci2) Tool. Indiana University and SciTech Strategies, <http://sci2.cns.iu.edu>.

**Börner, Katy, Valentin Pentchev, Matthew Hutchinson, James Pringle, Jason Rollins, Yadu N. Babuji, and Eamon Duede. 2017. "Web of Science™ as a Research Dataset". 16th International Conference on Scientometrics and Informetrics, Wuhan, China.**



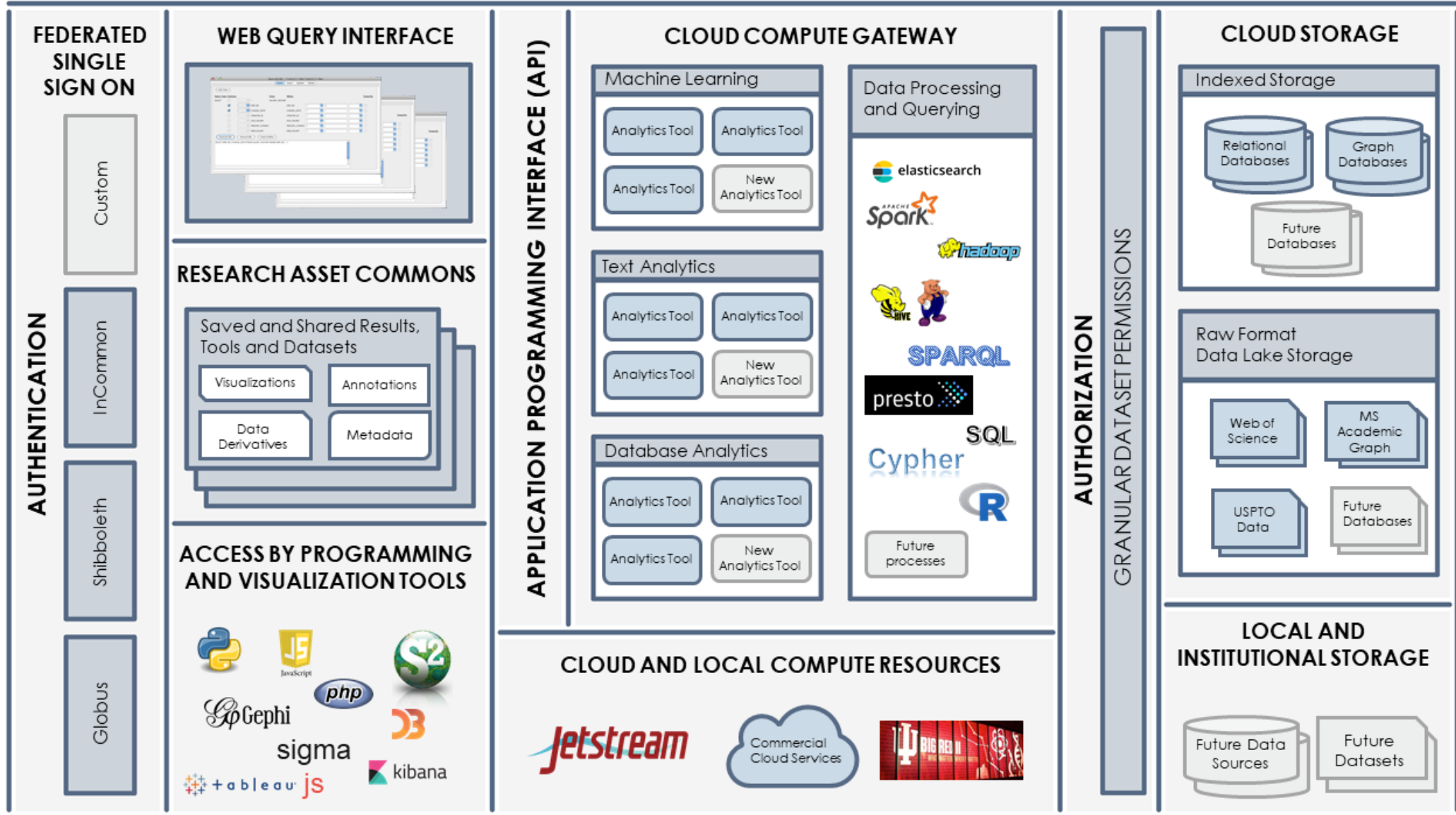
## *IUNI WoS Data Enclave*

The Indiana University Network Science Institute (IUNI) acquired the complete set of Clarivate Analytics' Web of Science XML raw data (Web of Knowledge version 5). The data was parsed and stored in a well-documented Postgresql database, see entity-relationship diagram, database schema, and data dictionary on <http://iuni.iu.edu/resources/web-of-science>. The code used to parse the WoS XML format and to save data in the Postgresql database was made available freely on GitHub, see [https://github.iu.edu/CNS/generic\\_parser](https://github.iu.edu/CNS/generic_parser). All data can be accessed via the IUNI WoS Data Enclave, a secure repository that uses IU's Karst high-throughput computing cluster designed to deliver large amounts of processing capacity over long periods of time. Access to the XML data and the PostgreSQL database is granted to a user's Karst account. IU faculty, staff, and qualifying sponsored affiliates can request accounts on Karst to use the data for academic research and without any sharing of data. A simple web browser based query interface to the WoS dataset was implemented to support custom queries for specific terms, journals, or authors. Datasets can be downloaded in CSV data format compatible with data mining and visualization tools such as Gephi or the Sci2 Tool (<http://sci2.cns.iu.edu>) (Sci2 Team, 2009). More about the IUNI WoS Data Enclave can be found at <http://iuni.iu.edu/resources/web-of-science>.

[https://github.com/lightr/generic\\_parser](https://github.com/lightr/generic_parser)

[Börner, Katy](#), Valentin Pentchev, Matthew Hutchinson, James Pringle, Jason Rollins, Yadu N. Babuji, and Eamon Duede. 2017. "[Web of Science™ as a Research Dataset](#)". *16th International Conference on Scientometrics and Informetrics, Wuhan, China*.

# SHARED BIGDATA-GATEWAY FOR RESEARCH LIBRARIES (SBD-GATEWAY)





# BIG Science of Science R&D





# Maps of Science & Technology

<http://scimaps.org>



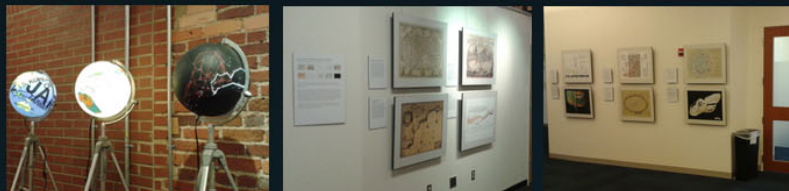
101st Annual Meeting of the Association of American Geographers, Denver, CO.  
April 5th - 9th, 2005 (First showing of Places & Spaces)



University of Miami, Miami, FL.  
September 4 - December 11, 2014.



The David J. Sencer CDC Museum, Atlanta, GA.  
January 25 - June 17, 2016.



Duke University, Durham, NC.  
January 12 - April 10, 2015

100 maps and 20 macrosopes by 250+ experts on display at 350+ venues in 28 countries.

# The Structure of Science

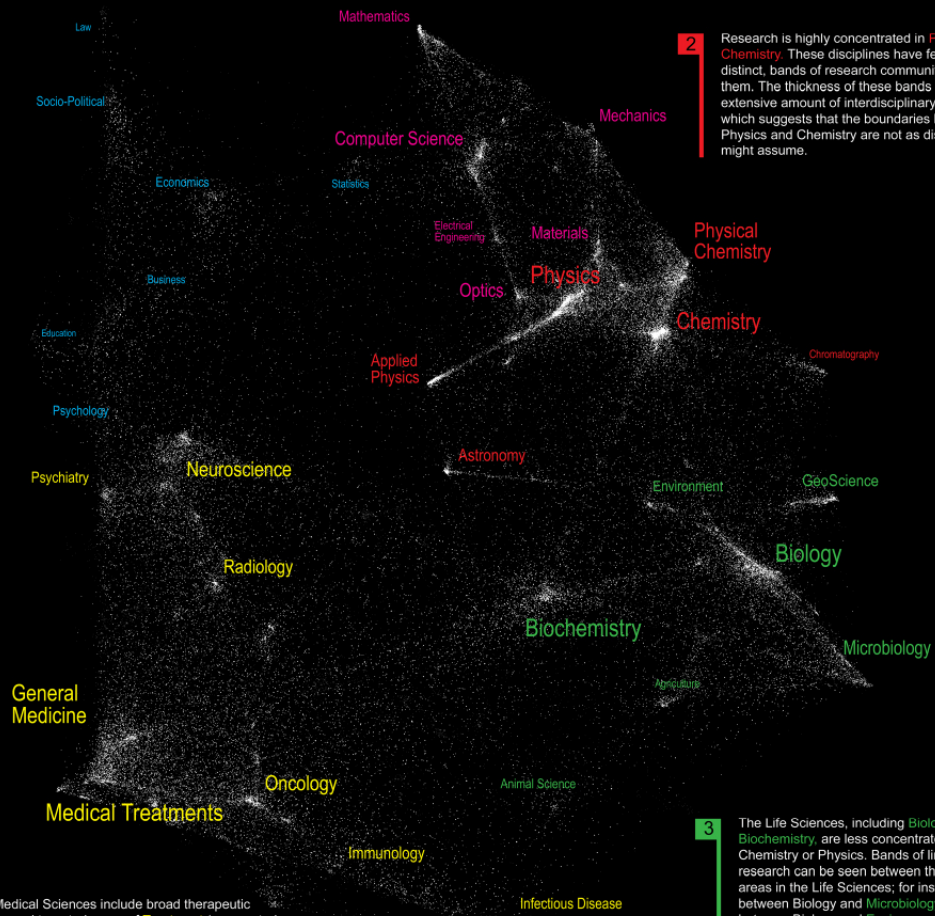
**5** The Social Sciences are the smallest and most diffuse of all the sciences. **Psychology** serves as the link between Medical Sciences (Psychiatry) and the Social Sciences. **Statistics** serves as the link with Computer Science and Mathematics.

**1** **Mathematics** is our starting point, the purest of all sciences. It lies at the outer edge of the map. **Computer Science**, **Electrical Engineering**, and **Optics** are applied sciences that draw upon knowledge in Mathematics and Physics. These three disciplines provide a good example of a linear progression from one pure science (Mathematics) to another (Physics) through multiple disciplines. Although applied, these disciplines are highly concentrated with distinct bands of research communities that link them. Bands indicate interdisciplinary research.

**2** Research is highly concentrated in **Physics** and **Chemistry**. These disciplines have few, but very distinct, bands of research communities that link them. The thickness of these bands indicates an extensive amount of interdisciplinary research, which suggests that the boundaries between Physics and Chemistry are not as distinct as one might assume.

**3** The Life Sciences, including **Biology** and **Biochemistry**, are less concentrated than Chemistry or Physics. Bands of linking research can be seen between the larger areas in the Life Sciences; for instance between Biology and Microbiology, and between Biology and Environmental Science. Biochemistry is very interesting in that it is a large discipline that has visible links to disciplines in many areas of the map, including Biology, Chemistry, Neuroscience, and General Medicine. It is perhaps the most interdisciplinary of the sciences.

**4** The Medical Sciences include broad therapeutic studies and targeted areas of **Treatment** (e.g. central nervous system, cardiology, gastroenterology, etc.) Unlike Physics and Chemistry, the medical disciplines are more spread out, suggesting a more multi-disciplinary approach to research. The transition into Life Sciences (via Animal Science and Biochemistry) is gradual.



We are all familiar with traditional maps that show the relationships between countries, provinces, states, and cities. Similar relationships exist between the various disciplines and research topics in science. This allows us to map the structure of science.

One of the first maps of science was developed at the Institute for Scientific Information over 30 years ago. It identified 41 areas of science from the citation patterns in 17,000 scientific papers. That early map was intriguing, but it didn't cover enough of science to accurately define its structure.

Things are different today. We have enormous computing power and advanced visualization software that make mapping of the structure of science possible. This galaxy-like map of science (left) was generated at Sandia National Laboratories using an advanced graph layout routine (VxOrd) from the citation patterns in 800,000 scientific papers published in 2002. Each dot in the galaxy represents one of the 96,000 research communities active in science in 2002. A research community is a group of papers (9 on average) that are written on the same research topic in a given year. Over time, communities can be born, continue, split, merge, or die.

The map of science can be used as a tool for science strategy. This is the terrain in which organizations and institutions locate their scientific capabilities. Additional information about the scientific and economic impact of each research community allows policy makers to decide which areas to explore, exploit, abandon, or ignore.

We also envision the map as an educational tool. For children, the theoretical relationship between areas of science can be replaced with a concrete map showing how math, physics, chemistry, biology and social studies interact. For advanced students, areas of interest can be located and neighboring areas can be explored.

## Nanotechnology

Most research communities in nanotechnology are concentrated in **Physics**, **Chemistry**, and **Materials Science**. However, many disciplines in the Life and Medical Sciences also have nanotechnology applications.

## Proteomics

Research communities in proteomics are centered in **Biochemistry**. In addition, there is a heavy focus in the tools section of chemistry, such as **Chromatography**. The balance of the proteomics communities are widely dispersed among the Life and Medical Sciences.

## Pharmacogenomics

Pharmacogenomics is a relatively new field with most of its activity in **Medicine**. It also has many communities in **Biochemistry** and two communities in the Social Sciences.



# Science related Wikipedian ACTIVITY

This visualization explores the activity of science, math, and technology (SMT) related articles in the English-language Wikipedia (<http://en.wikipedia.org>). The central image shows 659,388 articles (circles). Overlaid is a 37 x 37 grid of relevant half-inch sized images.

Blue, green, and yellow circles represent the 3,599 math, 6,474 science, and 3,164 technology related articles respectively. The larger the size of a circle the higher the likelihood it is that type of article. The four corners show activity patterns of SMT articles.

**Article Edit Activity**  
Articles are size coded based on how frequently they have been edited from Feb. 6, 2001 to April 6, 2007. More consideration is given to current and major edits. Larger circles have been edited more frequently than smaller circles.

**2007 Major Edits**  
Articles are size coded based on how many major edits they received from January 1st, 2007 to April 6th, 2007. Larger circles have received more edits than smaller circles. The highest number of major edits was 2,627.

For the central image, each article is size coded based on the likelihood that it is math, science, or technology related.

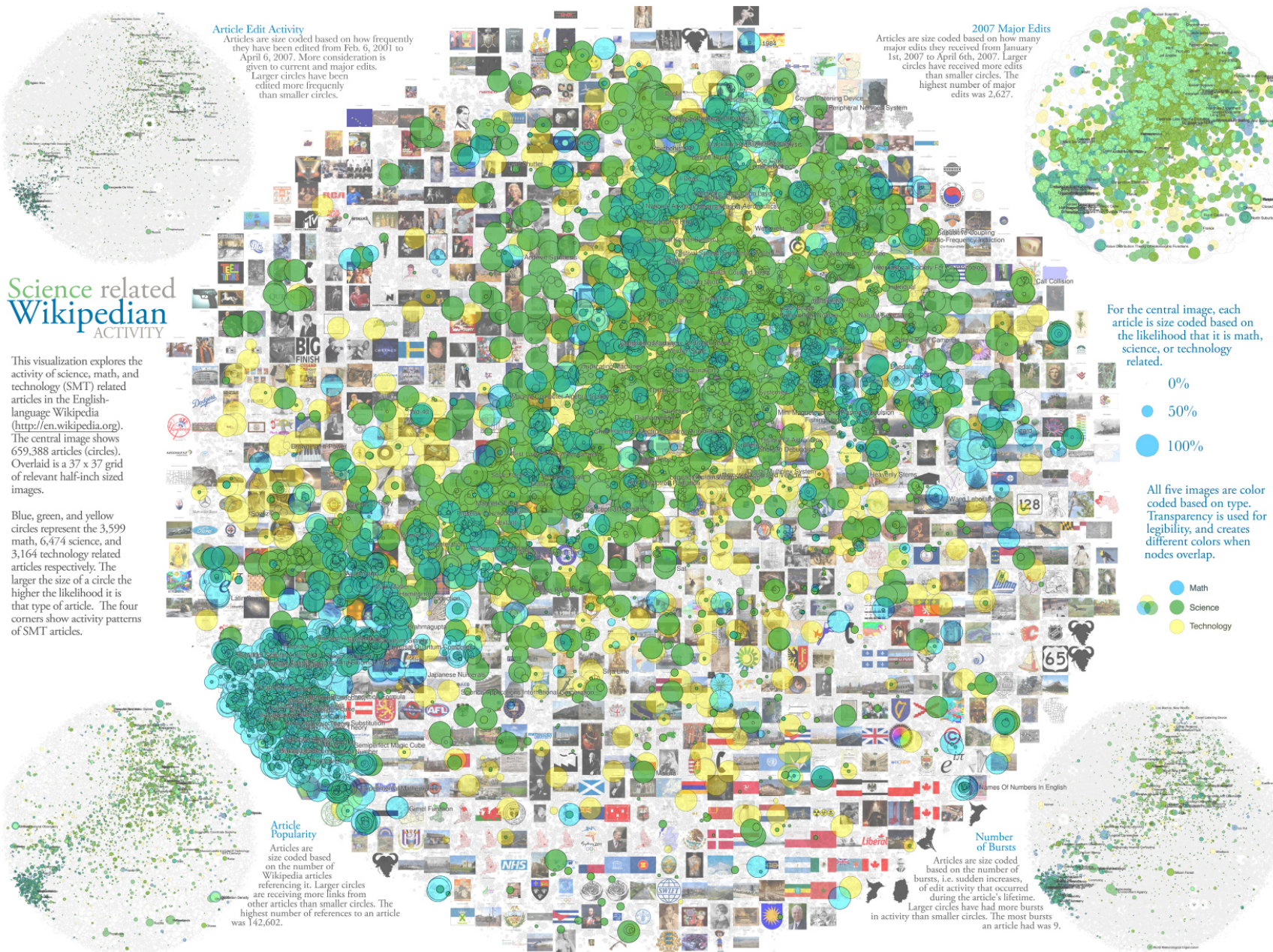
- 0%
- 50%
- 100%

All five images are color coded based on type. Transparency is used for legibility, and creates different colors when nodes overlap.

- Math
- Science
- Technology

**Article Popularity**  
Articles are size coded based on the number of Wikipedia articles referencing it. Larger circles are receiving more links from other articles than smaller circles. The highest number of references to an article was 142,602.

**Number of Bursts**  
Articles are size coded based on the number of bursts, i.e. sudden increases, of edit activity that occurred during the article's lifetime. Larger circles have had more bursts in activity than smaller circles. The most bursts an article had was 9.

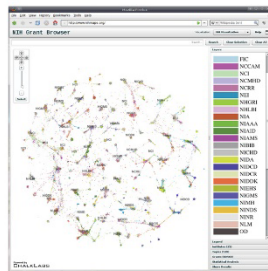




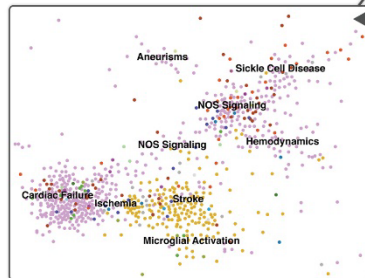
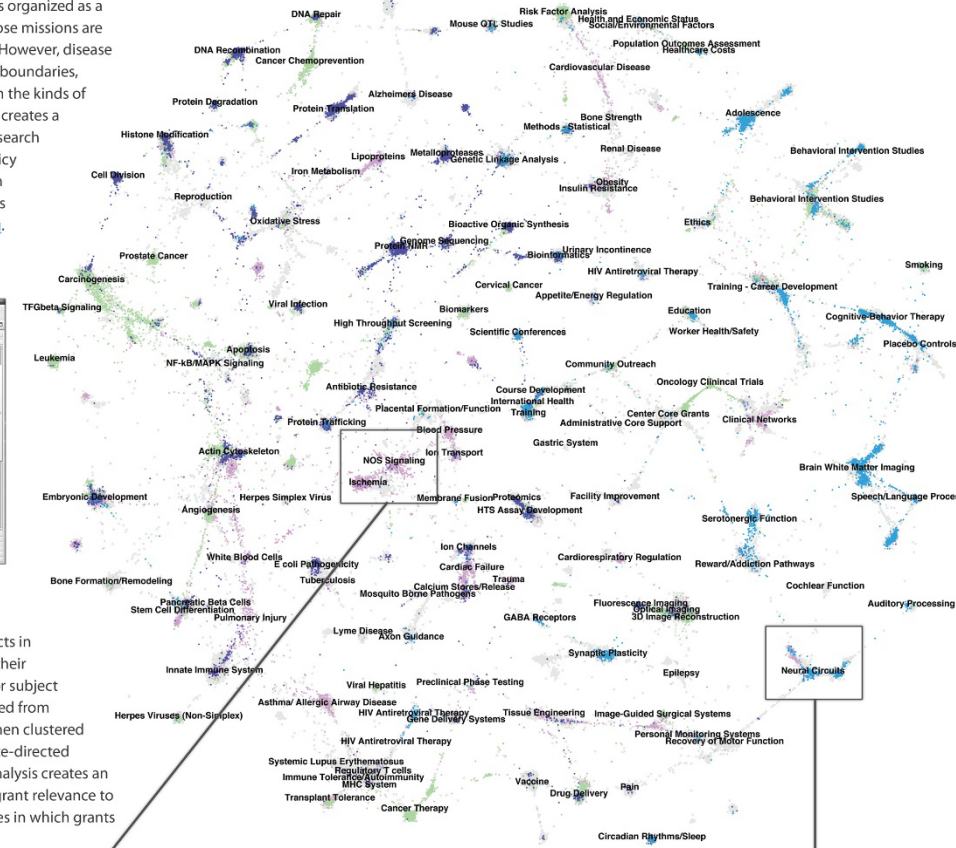
# A Topic Map of NIH Grants 2007

Bruce W. Herr II (Chalklabs & IU), Gully Burns (ISI), David Newman (UCI), Edmund Talley (NIH)

The National Institutes of Health (NIH) is organized as a multitude of Institutes and Centers whose missions are primarily focused on distinct diseases. However, disease etiologies and therapies flout scientific boundaries, and thus there is tremendous overlap in the kinds of research funded by each Institute. This creates a daunting landscape for decisions on research directions, funding allocations, and policy formulations. Shown here is devised an interactive topic map for navigating this landscape, online at [www.nihmaps.org](http://www.nihmaps.org). Institute abbreviations can be found at [www.nih.gov/icd](http://www.nih.gov/icd).

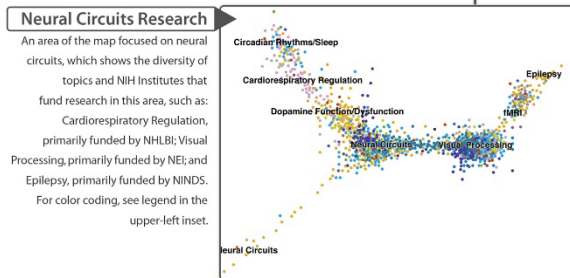


Topic modeling, a statistical technique that automatically learns semantic categories, was applied to assess projects in terms used by researchers to describe their work, without the biases of keywords or subject headings. Grant similarities were derived from their topic mixtures, and grants were then clustered on a two-dimensional map using a force-directed simulated annealing algorithm. This analysis creates an interactive environment for assessing grant relevance to research categories and to NIH Institutes in which grants are localized.



### Cardiac Diseases Research

An area of the map focused on cardiovascular function and dysfunction. Cardiac Failure (primarily funded by NHLBI) is typically clustered next to Stroke (NINDS), since these are the two major medical emergencies associated with ischemia, which results from a restricted blood supply. Also localized in this area are grants focused on Nitric Oxide (NOS) Signaling, a major biochemical pathway for vasodilation, and grants on Hemodynamics, Sickle Cell Disease, and Aneurysms.

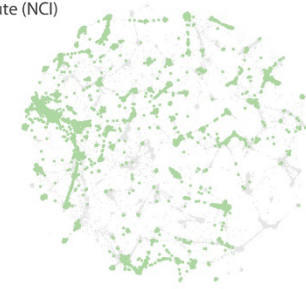


### Neural Circuits Research

An area of the map focused on neural circuits, which shows the diversity of topics and NIH Institutes that fund research in this area, such as Cardiorespiratory Regulation, primarily funded by NHLBI; Visual Processing, primarily funded by NEI and Epilepsy, primarily funded by NINDS. For color coding, see legend in the upper-left inset.

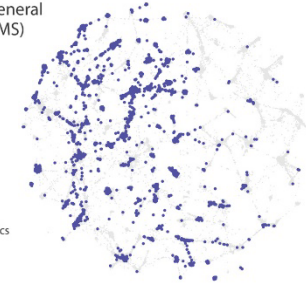
### National Cancer Institute (NCI)

- TOP 10 TOPICS
- 1 Oncology Clinical Trials
  - 2 Cancer Treatment
  - 3 Cancer Therapy
  - 4 Carcinogenesis
  - 5 Risk Factor Analysis
  - 6 Cancer Chemotherapy
  - 7 Metastasis
  - 8 Leukemia
  - 9 Prediction/Prognosis
  - 10 Cancer Chemoprevention



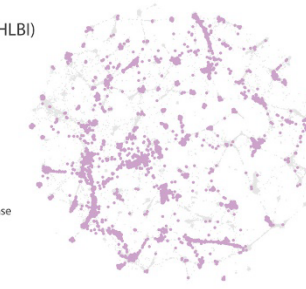
### National Institute of General Medical Sciences (NIGMS)

- TOP 10 TOPICS
- 1 Bioactive Organic Synthesis
  - 2 X-ray Crystallography
  - 3 Protein NMR
  - 4 Computational Models
  - 5 Yeast Biology
  - 6 Metalloproteases
  - 7 Enzymatic Mechanisms
  - 8 Protein Complexes
  - 9 Invertebrate/Zebrafish Genetics
  - 10 Cell Division



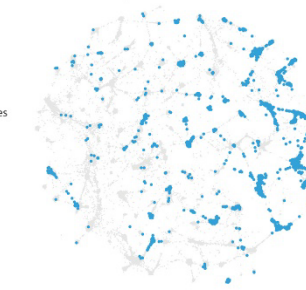
### National Heart, Lung, and Blood Institute (NHLBI)

- TOP 10 TOPICS
- 1 Cardiac Failure
  - 2 Pulmonary Injury
  - 3 Genetic Linkage Analysis
  - 4 Cardiovascular Disease
  - 5 Atherosclerosis
  - 6 Hemostasis
  - 7 Blood Pressure
  - 8 Asthma/ Allergic Airway Disease
  - 9 Gene Association
  - 10 Lipoproteins



### National Institute of Mental Health (NIMH)

- TOP 10 TOPICS
- 1 Mood Disorders
  - 2 Schizophrenia
  - 3 Behavioral Intervention Studies
  - 4 Mental Health
  - 5 Depression
  - 6 Cognitive-Behavior Therapy
  - 7 AIDS Prevention
  - 8 Genetic Linkage Analysis
  - 9 Adolescence
  - 10 Childhood





# Map of Scientific Collaborations from 2005-2009



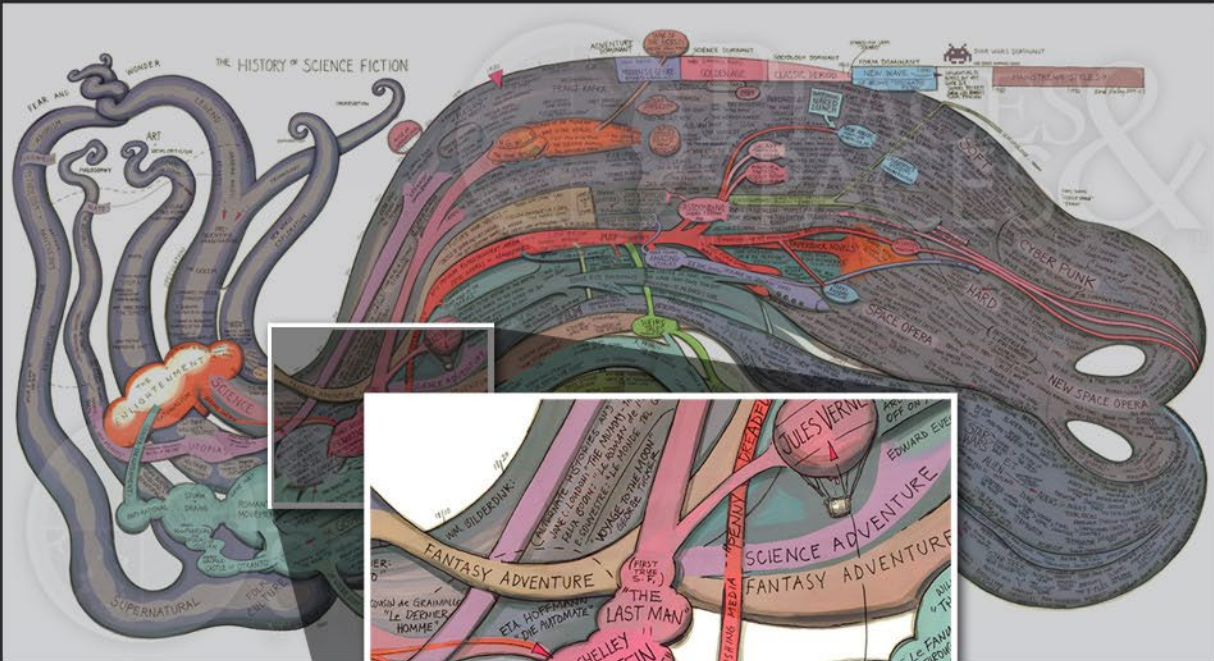
Computed Using Data from Elsevier's Scopus








# Check out our **Zoom Maps** online!



**VII.10**  
History of Science Fiction, by Ward Shelley

BROOKLYN, NY, 2011  
Courtesy of Ward Shelley Studio

Ward Shelley is an artist identified with the Williamsburg scene in Brooklyn, New York, about art and culture. This map plots the science fiction literary genre from its nascence emerging out of the data, here the narrative structure precedes and organizes the data. The monster whose tentacles are like trace roots to pre-historical sources and whose body, Romanticism, which birthed gothic fiction, source not only of Sci-Fi, but also of criticism progressed through a number of distinct periods, which are charted, citing hundreds of



Visit [scimaps.org](http://scimaps.org) and check out all our maps in stunning detail!

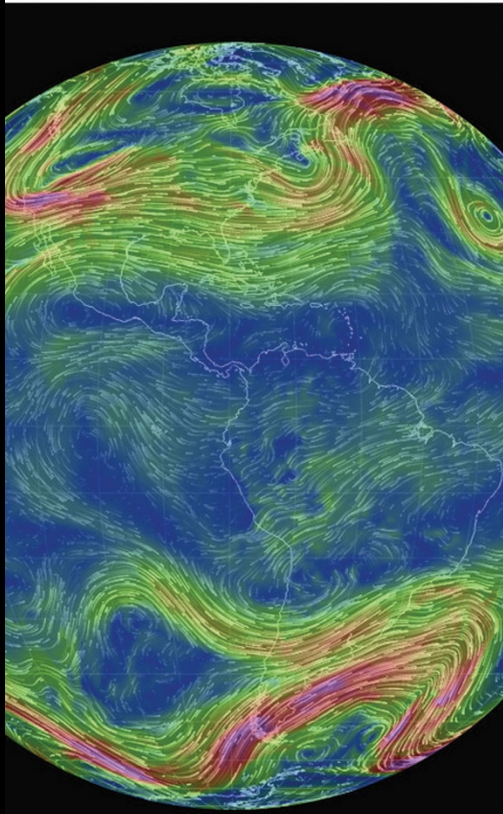




# MACROSCOPES FOR INTERACTING WITH SCIENCE



**PLACES &  
SPACES &**  
MAPPING SCIENCE  
scimaps.org



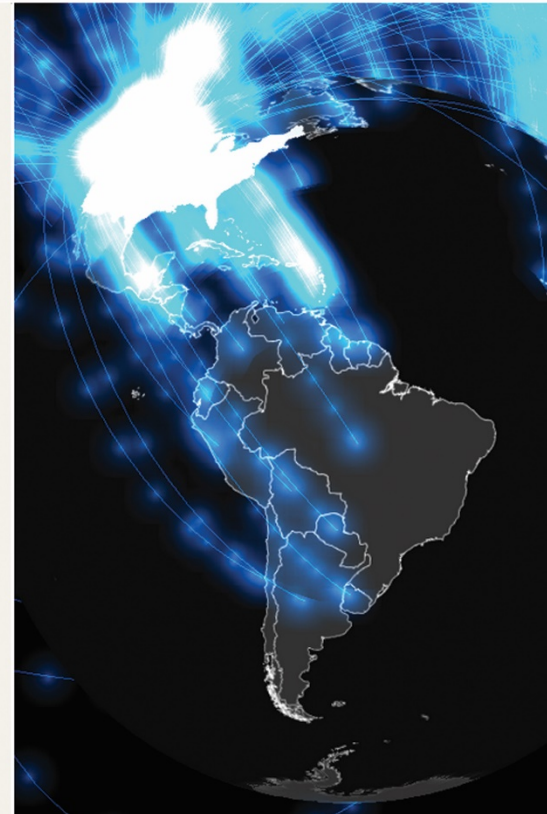
## Earth

*Weather on a worldwide scale*



## AcademyScope

*Exploring the scientific landscape*



## Mapping Global Society

*Local news from a global perspective*



## Charting Culture

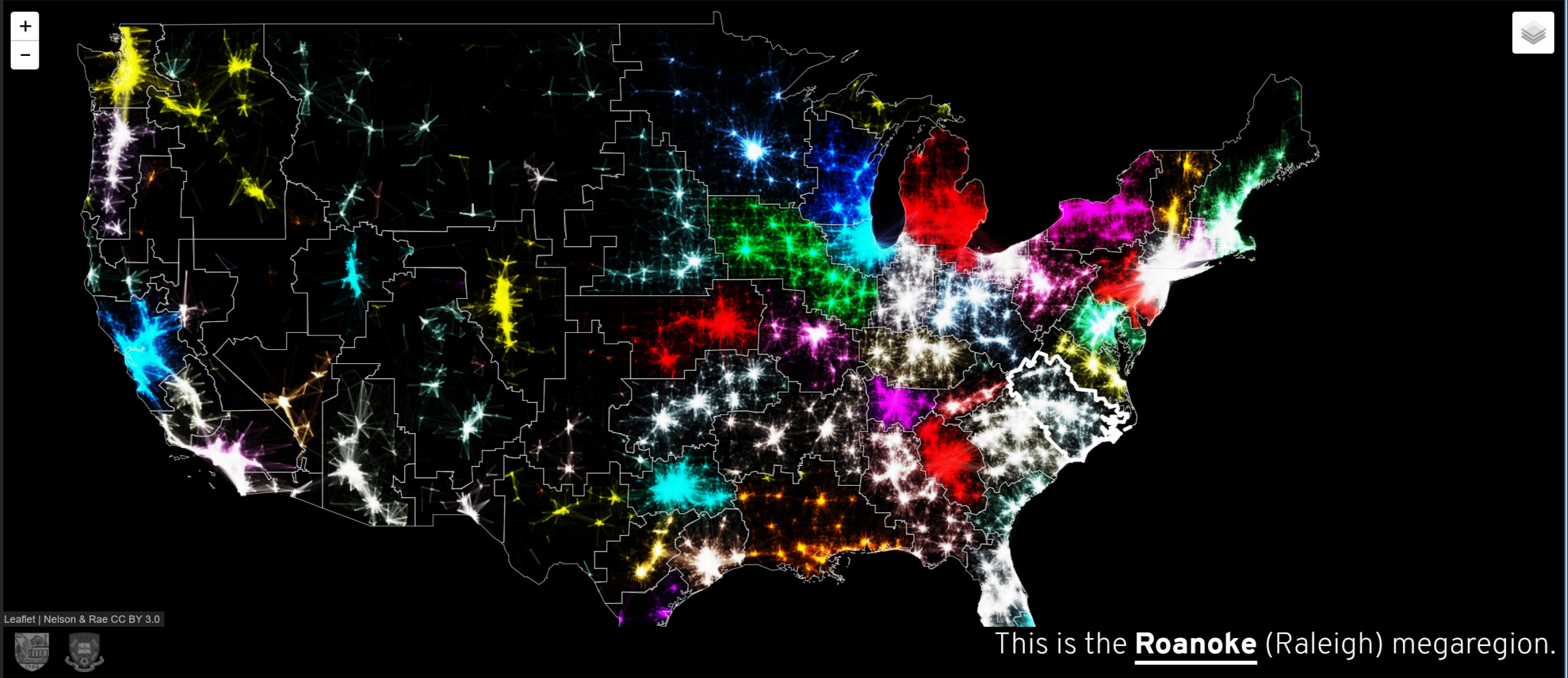
*2,600 years of human history in 5 minutes*

<http://idemo.cns.iu.edu/macroscope-kiosk>



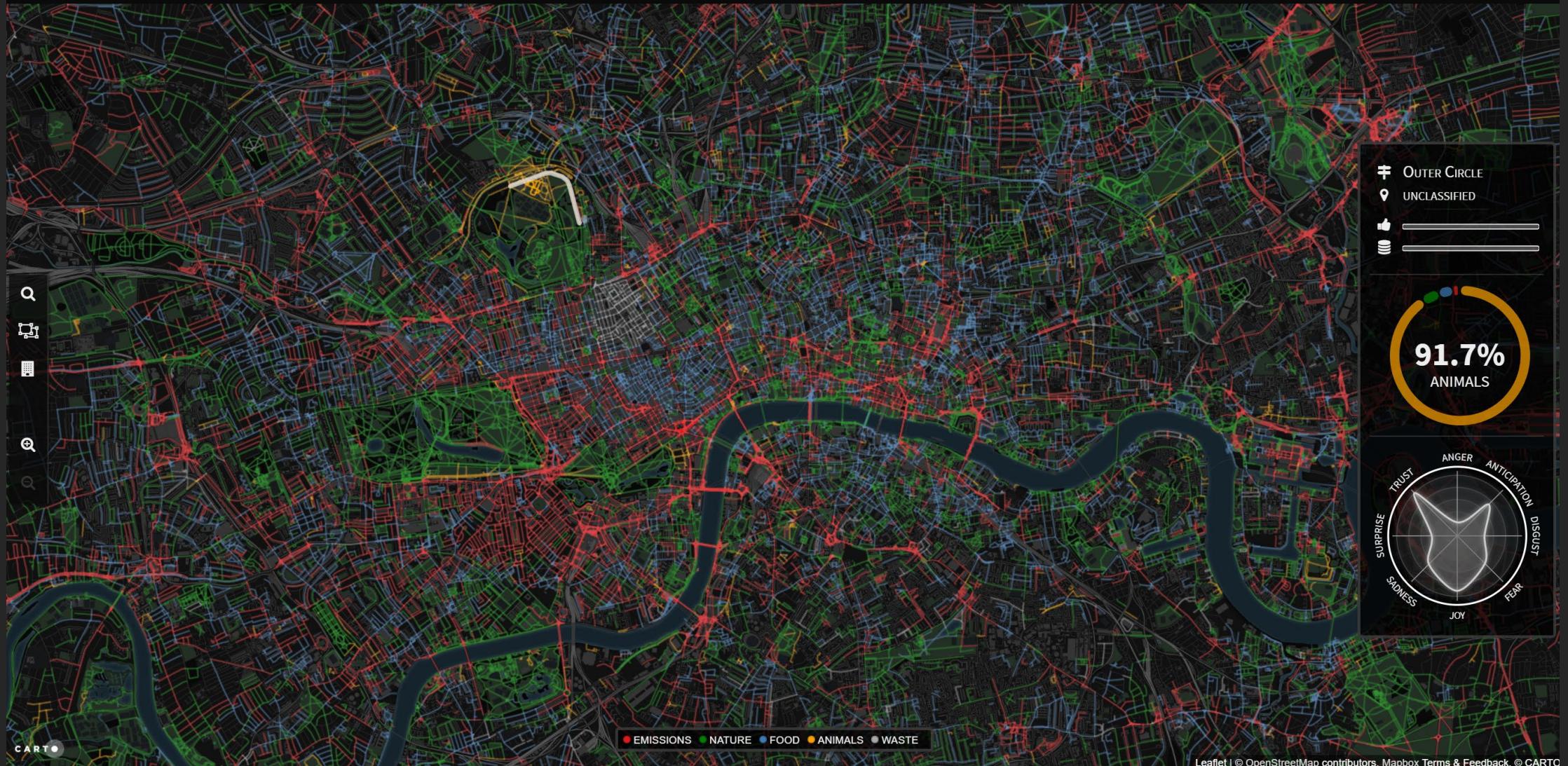
# THE MEGAREGIONS OF THE US

Explore the new geography of commuter connections in the US.  
Tap to identify regions. Tap and hold to see a single location's commuted.





SMELLY  
MAPS



Smelly Maps – Daniele Quercia, Rossano Schifanella, and Luca Maria Aiello – 2015



# The News Co-occurrence Globe

An interactive visualization of how countries are mentioned together in the world's news media

+ - UNITED KINGDOM SEARCH ABOUT

2.92K  
COOCCUR%



**UNITED KINGDOM** cooccurrences in: 2,922%  
cooccurrences out: 80%



COOCCURR

<input checked="" type="checkbox"/>	IN%
<input checked="" type="checkbox"/>	OUT%





# Modeling Science, Technology & Innovation Conference

WASHINGTON D.C. | MAY 17-18, 2016

[View Agenda](#)

Government, academic, and industry leaders discussed challenges and opportunities associated with using big data, visual analytics, and computational models in STI decision-making.

Conference slides, recordings, and report are available via <http://modsti.cns.iu.edu/report>







## Modeling and Visualizing Science and Technology Developments

National Academy of Sciences Sackler Colloquium, December 4-5, 2017, Irvine, CA

### Rankings and the Efficiency of Institutions

H. Eugene Stanley | Albert-László Barabási | Lada Adamic | Marta González | Kaye Husbands Fealing | Brian Uzzi | John V. Lombardi

### Higher Education and the Science & Technology Job Market

Katy Börner | Wendy L. Martinez | Michael Richey | William Rouse | Stasa Milojevic | Rob Rubin | David Krakauer

### Innovation Diffusion and Technology Adoption

William Rouse | Donna Cox | Jeff Alstott | Ben Shneiderman | Rahul C. Basole | Scott Stern | Cesar Hidalgo

### Modeling Needs, Infrastructures, Standards

Paul Trunfio | Sallie Keller | Andrew L. Russell | Guru Madhavan | Azer Bestavros | Jason Owen-Smith



## PROGRAMS

### Sackler Colloquia

- » [About Sackler Colloquia](#)
- » [Upcoming Colloquia](#)
- » [Completed Colloquia](#)
- » [Sackler Lectures](#)
- » [Video Gallery](#)
- » [Connect with Sackler Colloquia](#)
- » [Give to Sackler Colloquia](#)

### Cultural Programs

### Distinctive Voices

### Kavli Frontiers of Science

### Keck Futures Initiative

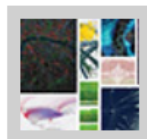
### LabX

### Sackler Forum

### Science & Entertainment Exchange



## Modeling and Visualizing Science and Technology Developments



December 4-5, 2017; Irvine, CA

Organized by Katy Börner, H. Eugene Stanley, William Rouse and Paul Trunfio

### Overview

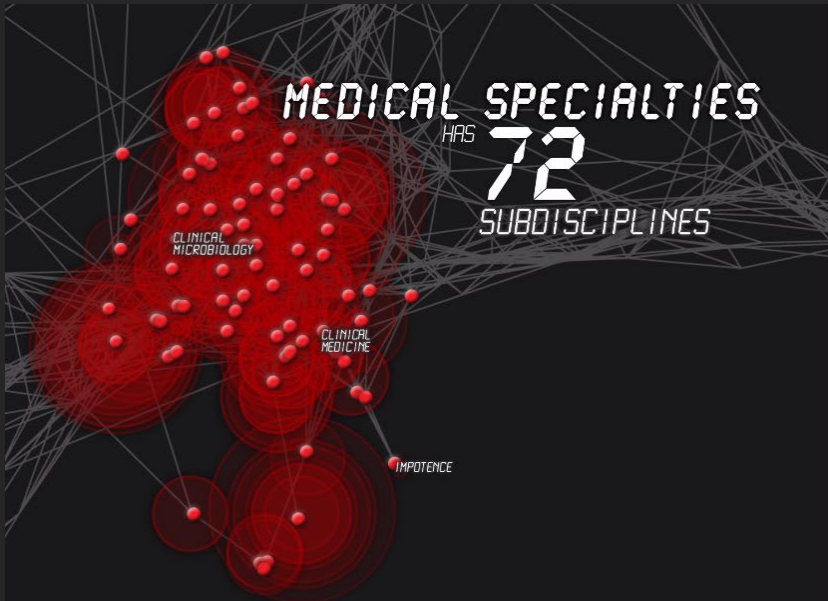
This colloquium was held in Irvine, CA on December 4-5, 2017.

This colloquium brought together researchers and practitioners from multiple disciplines to present, discuss, and advance computational models and visualizations of science and technology (S&T). Existing computational models are being applied by academia, government, and industry to explore questions such as: What jobs will exist in ten years and what career paths lead to success? Which types of institutions will likely be most innovative in the future? How will the higher education cost bubble burst affect these institutions? What funding strategies have the highest return on investment? How will changing demographics, alternative economic growth trajectories, and relationships among nations impact answers to these and other questions? Large-scale datasets (e.g., publications, patents, funding, clinical trials, stock market, social media data) can now be utilized to simulate the structure and evolution of S&T. Advances in computational power have created the possibility of implementing scalable, empirically validated computational models. However, because the databases are massive and multidimensional, both the data and the models tend to exceed human comprehension. How can advances in data visualizations be effectively employed to communicate the data, the models, and the model results to diverse stakeholder groups? Who will be the users of next generation models and visualizations and what decisions will they be addressing.

Videos of the talks are available on the [Sackler YouTube Channel](#).

<https://www.pnas.org/modeling>





Science Forecast S1:E1





[https://www.youtube.com/watch?v=IByX2\\_eb\\_QQ](https://www.youtube.com/watch?v=IByX2_eb_QQ)



# Arthur M. Sackler Colloquium on Modeling and Visualizing Science and Technology Developments

## ✔ **Twin-Win Model: A human-centered approach to research success**

Ben Shneiderman

PNAS December 11, 2018 115 (50) 12590-12594; first published December 10, 2018. <https://doi.org/10.1073/pnas.1802918115>

## ✔ **Forecasting innovations in science, technology, and education**

FROM THE COVER

Katy Börner, William B. Rouse, Paul Trunfio, and H. Eugene Stanley

PNAS December 11, 2018 115 (50) 12573-12581; first published December 10, 2018. <https://doi.org/10.1073/pnas.1818750115>

## ✔ **How science and technology developments impact employment and education**

Wendy Martinez

PNAS December 11, 2018 115 (50) 12624-12629; first published December 10, 2018. <https://doi.org/10.1073/pnas.1803216115>

## ✔ **Scientific prize network predicts who pushes the boundaries of science**

Yifang Ma and Brian Uzzi

PNAS December 11, 2018 115 (50) 12608-12615; first published December 10, 2018. <https://doi.org/10.1073/pnas.1800485115>

## ✔ **The role of industry-specific, occupation-specific, and location-specific knowledge in the growth and survival of new firms**

C. Jara-Figueroa, Bogang Jun, Edward L. Glaeser, and Cesar A. Hidalgo

PNAS December 11, 2018 115 (50) 12646-12653; first published December 10, 2018. <https://doi.org/10.1073/pnas.1800475115>



## Arthur M. Sackler Colloquium on Modeling and Visualizing Science and Technology Developments

### ✔ Skill discrepancies between research, education, and jobs reveal the critical need to supply soft skills for the data economy

Katy Börner, Olga Scrivner, Mike Gallant, Shutian Ma, Xiaozhong Liu, Keith Chewing, Lingfei Wu, and James A. Evans  
PNAS December 11, 2018 115 (50) 12630-12637; first published December 10, 2018. <https://doi.org/10.1073/pnas.1804247115>

### ✔ Changing demographics of scientific careers: The rise of the temporary workforce

Staša Milojević, Filippo Radicchi, and John P. Walsh  
PNAS December 11, 2018 115 (50) 12616-12623; first published December 10, 2018. <https://doi.org/10.1073/pnas.1800478115>

### ✔ The chaperone effect in scientific publishing

Vedran Sekara, Pierre Deville, Sebastian E. Ahnert, Albert-László Barabási, Roberta Sinatra, and Sune Lehmann  
PNAS December 11, 2018 115 (50) 12603-12607; first published December 10, 2018. <https://doi.org/10.1073/pnas.1800471115>

### ✔ Modeling research universities: Predicting probable futures of public vs. private and large vs. small research universities

William B. Rouse, John V. Lombardi, and Diane D. Craig  
PNAS December 11, 2018 115 (50) 12582-12589; first published December 10, 2018. <https://doi.org/10.1073/pnas.1807174115>

and more ...

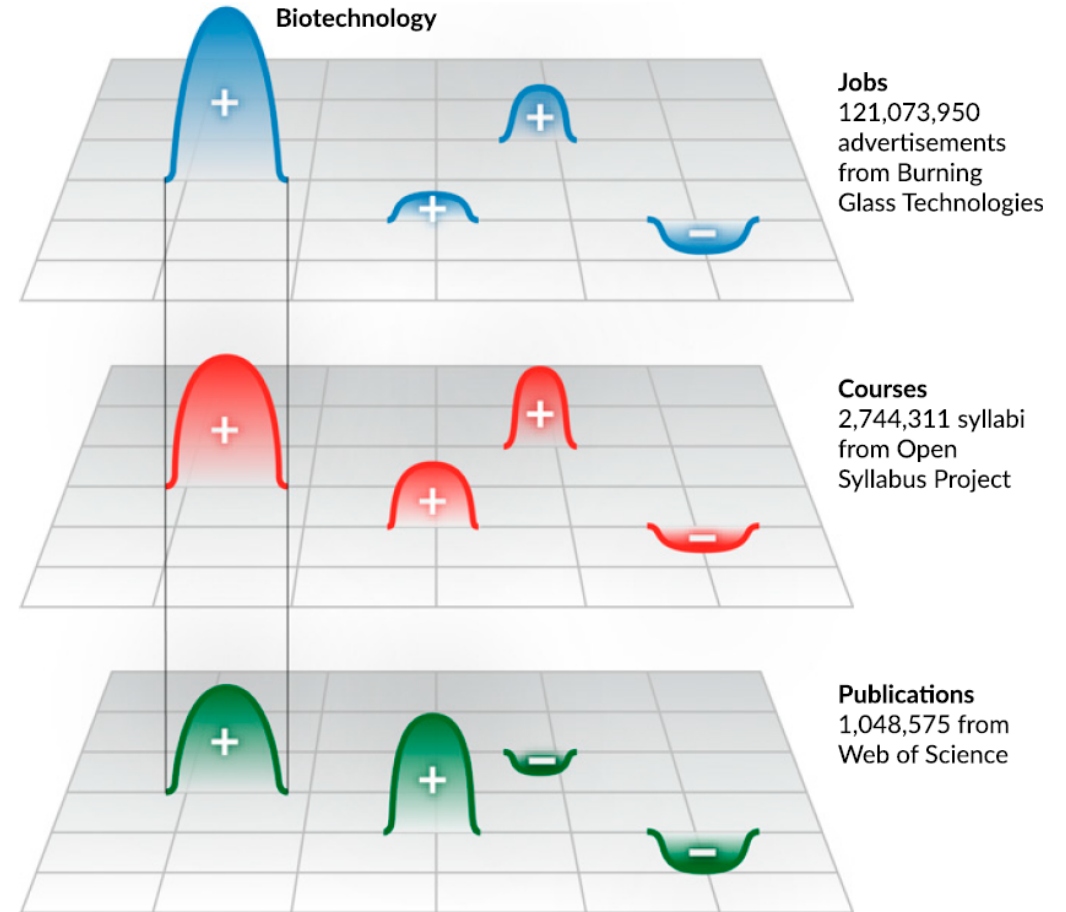




Study the **(mis)match** and **temporal dynamics** of S&T progress, education and workforce development options, and job requirements.

### Challenges:

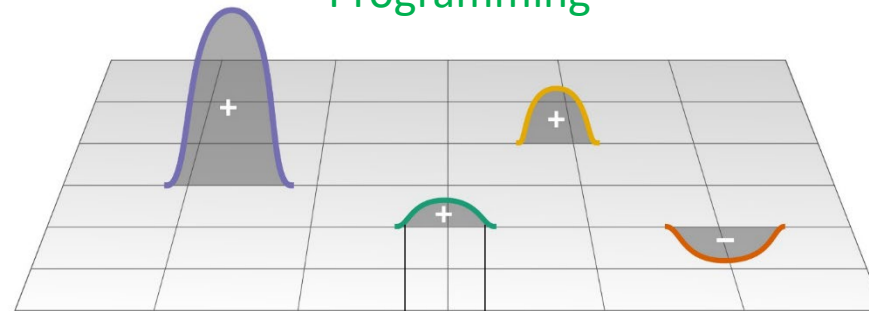
- Rapid change of STEM knowledge
- Increase in tools, AI
- Social skills (project management, team leadership)
- Increasing team size



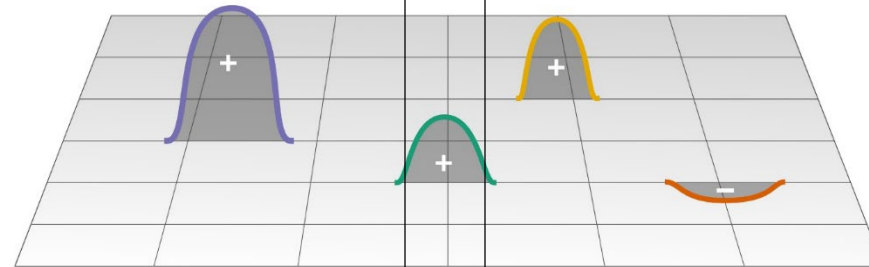
**Fig. 1.** The interplay of job market demands, educational course offerings, and progress in S&T as captured in publications. Color-coded mountains (+) and valleys (-) indicate different skill clusters. For example, skills related to Biotechnology might be mentioned frequently in job descriptions and taught in many courses, but they may not be as prevalent in academic publications. In other words, there are papers that mention these skills, but labor demand and commercial activity might be outstripping publication activity in this area. The numbers of jobs, courses, and publications that have skills associated and are used in this study are given on the right.



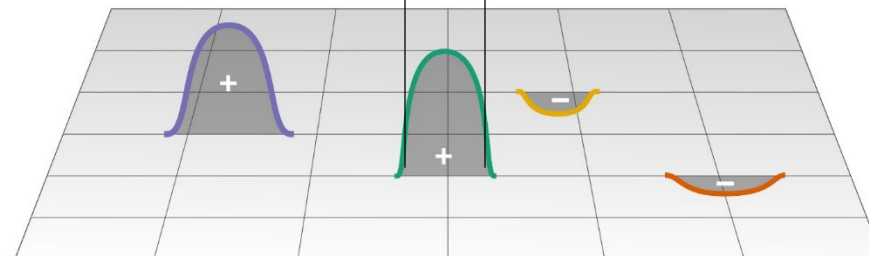
## Programming



**Jobs**

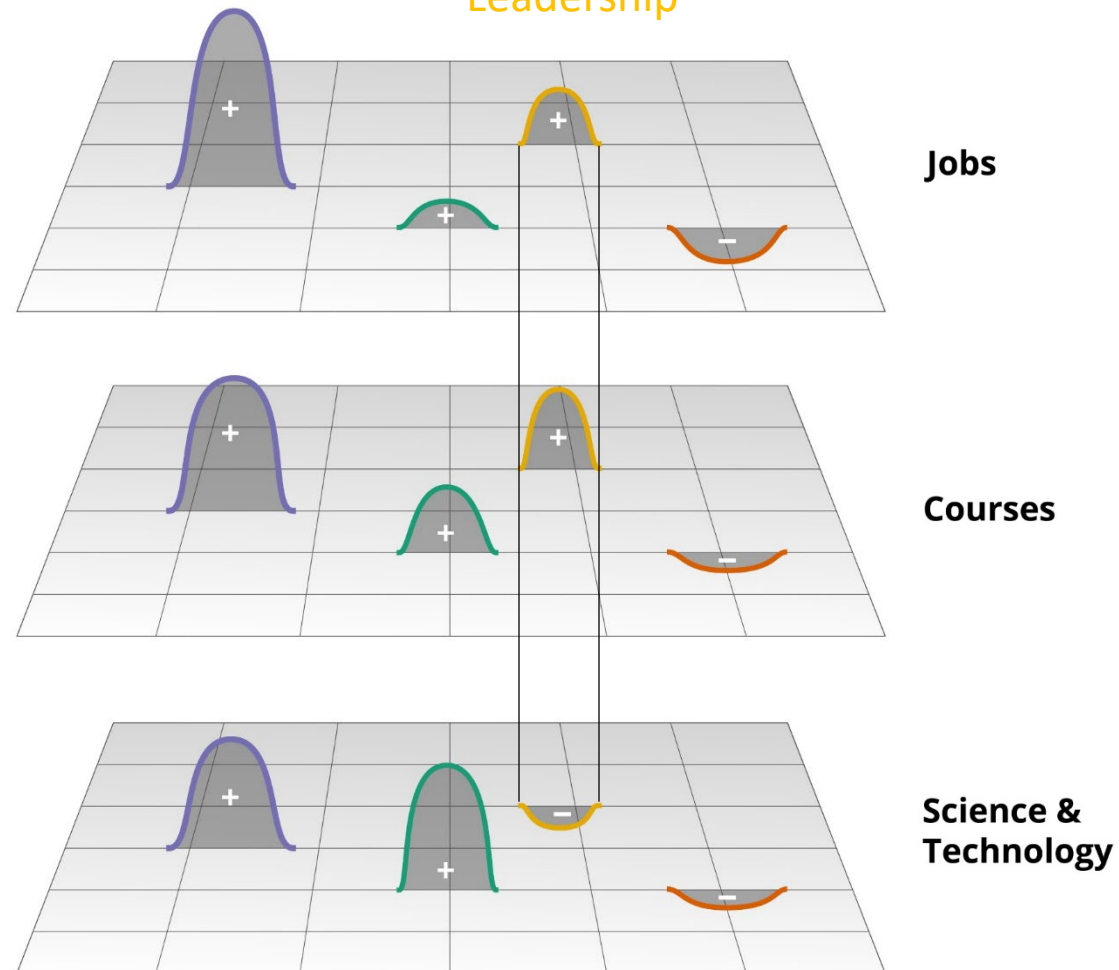


**Courses**



**Science & Technology**

## Leadership



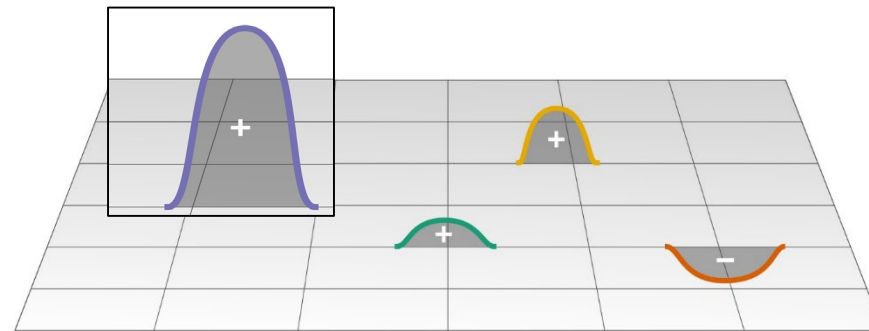
**Jobs**

**Courses**

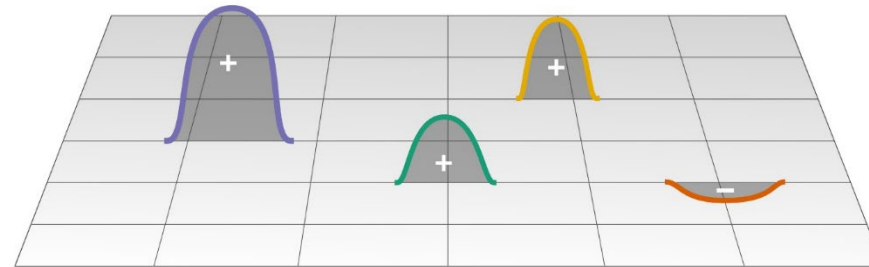
**Science &  
Technology**



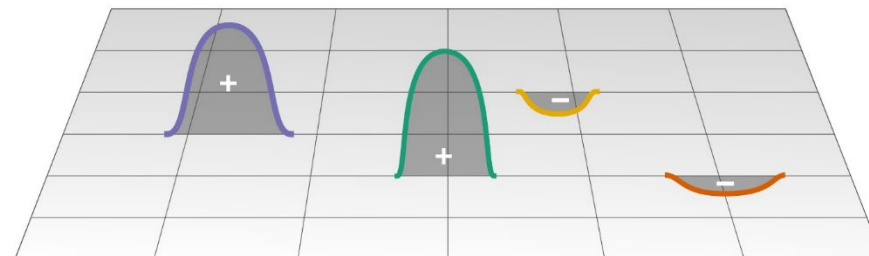
## Biotechnology



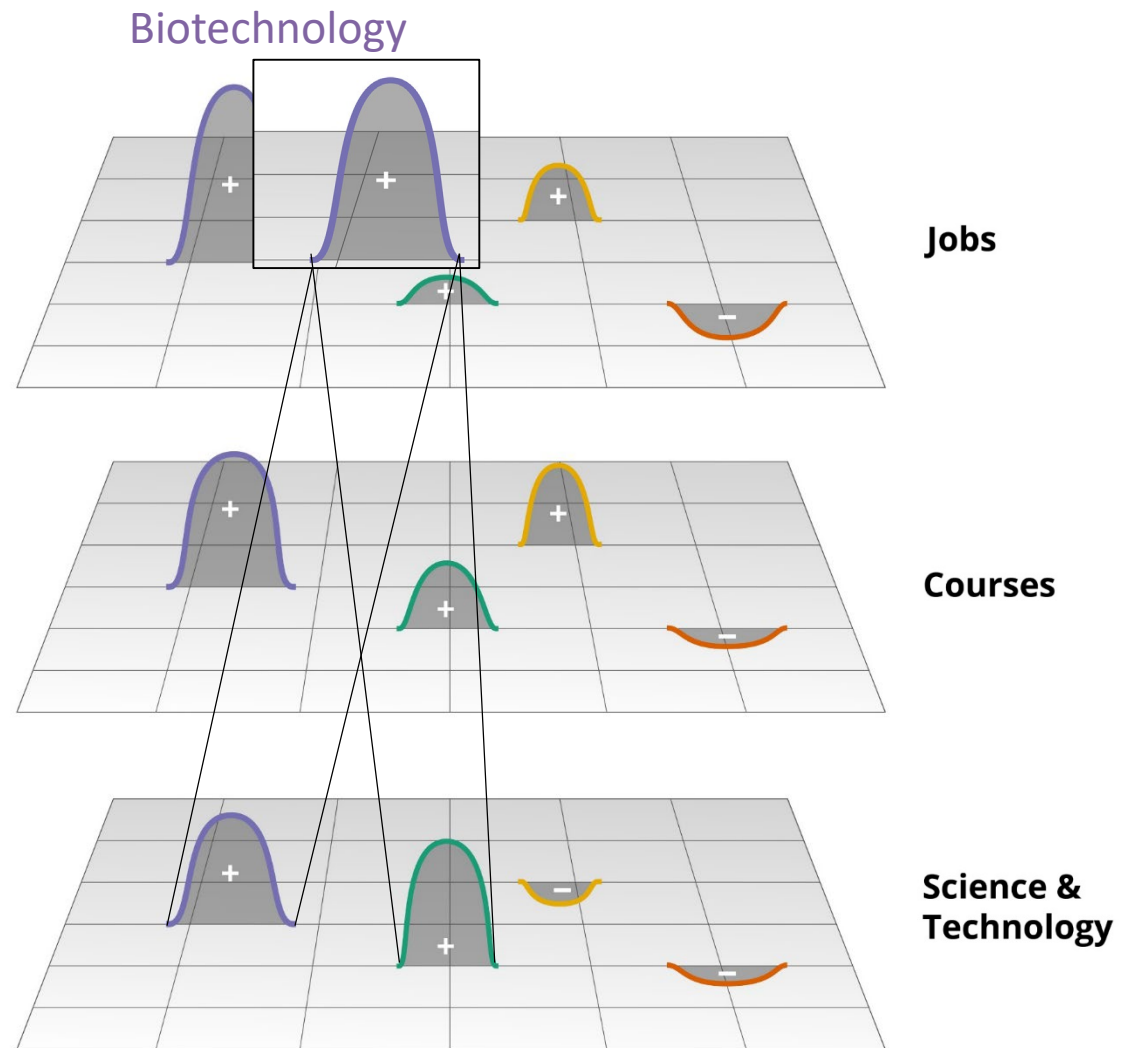
**Jobs**



**Courses**



**Science &  
Technology**





# Stakeholders and Insight Needs

- **Students:** What jobs will exist in 1-4 years? What program/learning trajectory is best to get/keep my dream job?
- **Teachers:** What course updates are needed? What balance of timely and timeless knowledge (to get a job vs. learn how to learn) should I teach? How to innovate in teaching and maintain job security or tenure?
- **Universities:** What programs should be created? What is my competition doing? How do I tailor programs to fit local needs?
- **Science Funders:** How can S&T investments improve short- and long-term prosperity? Where will advances in knowledge also yield advances in skills and technology?
- **Employers:** What skills are needed next year and in 5 and 10 years? Which institutions produce the right talent? What skills does my competition list in job advertisements?
- **Economic Developers:** What critical skills are needed to improve business retention, expansion, and recruitment in a region?

**What is ROI of my time, money, compassion?**

# Urgency

- 35% of UK jobs, and 30% in London, are at high risk from automation over the coming 20 years.  
<https://www2.deloitte.com/content/dam/Deloitte/uk/Documents/uk-futures/london-futures-agiletown.pdf>
- The aerospace industry and NASA have a disproportionately large percentage of workers aged 50 and older compared to the national average, and up to **half of the current workforce** will be eligible for retirement within the coming five years.  
Astronautics AIAA (2012) Recruiting, retaining, and developing a world-class aerospace workforce.  
[https://www.aiaa.org/uploadedFiles/Issues\\_and\\_Advocacy/Education\\_and\\_Workforce/Aerospace%20Workforce-%20030112.pdf](https://www.aiaa.org/uploadedFiles/Issues_and_Advocacy/Education_and_Workforce/Aerospace%20Workforce-%20030112.pdf)
- The rise of artificial intelligence will lead to the displacement of **millions of blue-collar as well as white-collar jobs** in the coming decade. Auerswald PE (2017) The Code Economy: A Forty-thousand-year History; Beyer D (2016) The future of machine intelligence: Perspectives from leading practitioners ; Brynjolfsson E, McAfee A (2014) The second machine age: Work, progress, and prosperity in a time of brilliant technologies; Ford M (2015) Rise of the Robots: Technology and the Threat of a Jobless Future.



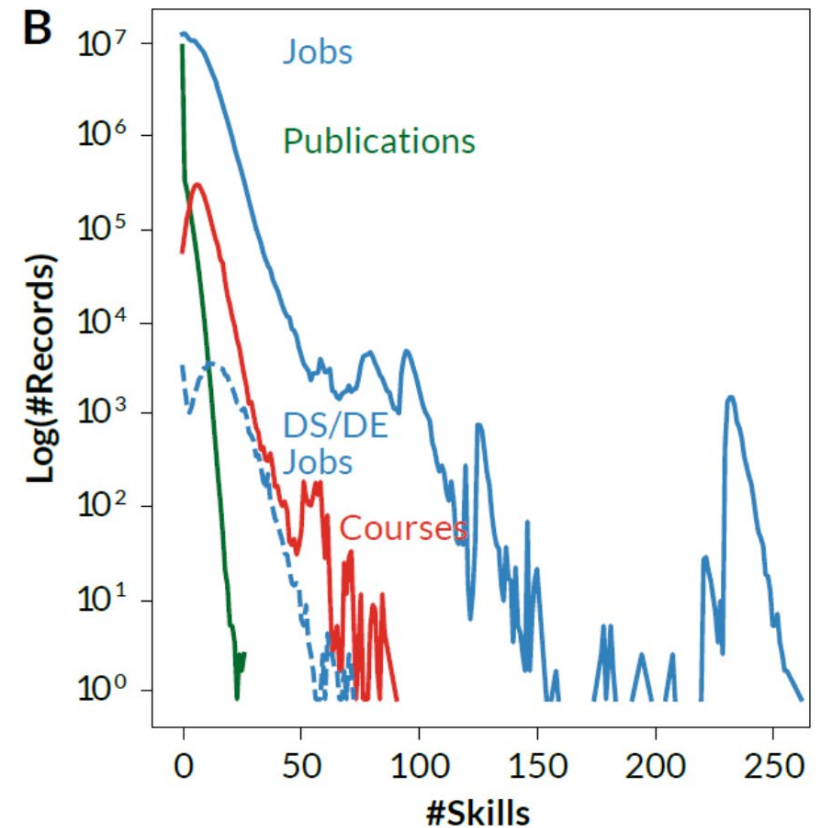
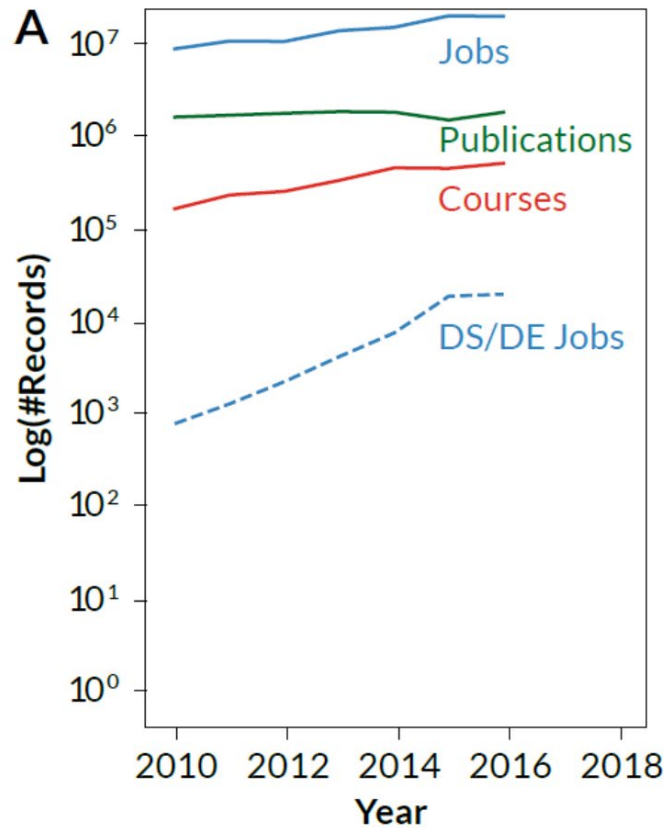


# Datasets Used

Job advertisements by Burning Glass posted between Jan 2010-Dec 2016.

Web of Science publications published Jan 2010-Dec 2016.

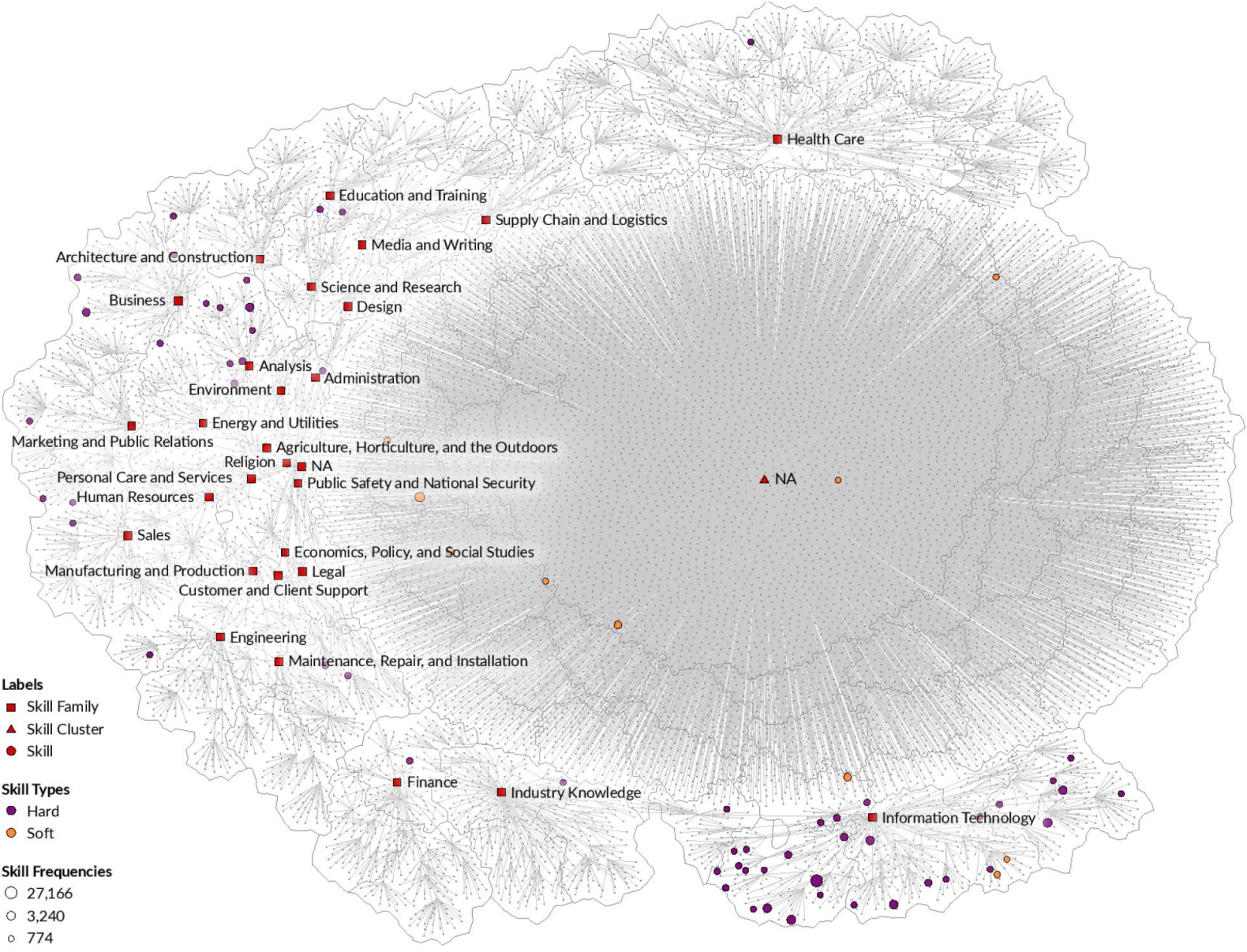
Course descriptions from the Open Syllabus Project acquired in June 2018 for courses offered in 2010-2016.



Data Type	#Records	#Records with skills	#Records without skills
All Courses	3,062,277	2,744,311	54,733
All Jobs	132,011,926	121,073,950	10,937,976
DSDE Jobs	69,405	65,944	3,461
All Publications	15,691,162	1,048,575	14,642,587
DSDE Publications	1,048,575	807,756	240,819

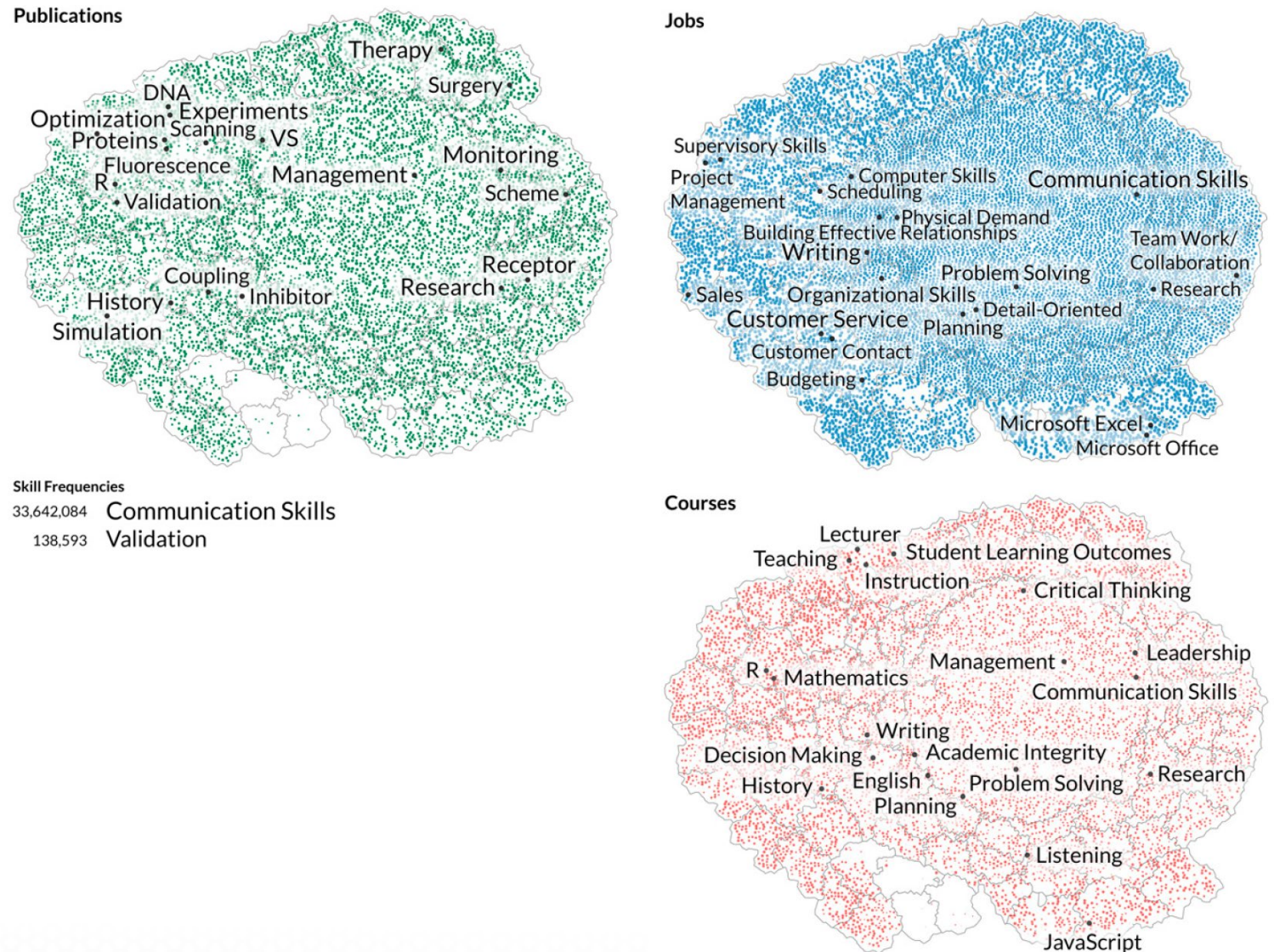


**Fig. 2.** Basemap of 13,218 skills. In this map, each dot is a skill, triangles identify skill clusters, and squares represent skill families from the Burning Glass (BG) taxonomy. Labels are given for all skill family nodes and for the largest skill cluster (NA) to indicate placement of relevant subtrees. Additionally, hard and soft skills are overlaid using purple and orange nodes, respectively; node area size coding indicates base 10 log of skill frequency in DS/DE jobs. Skill area computation uses Voronoi tessellation.

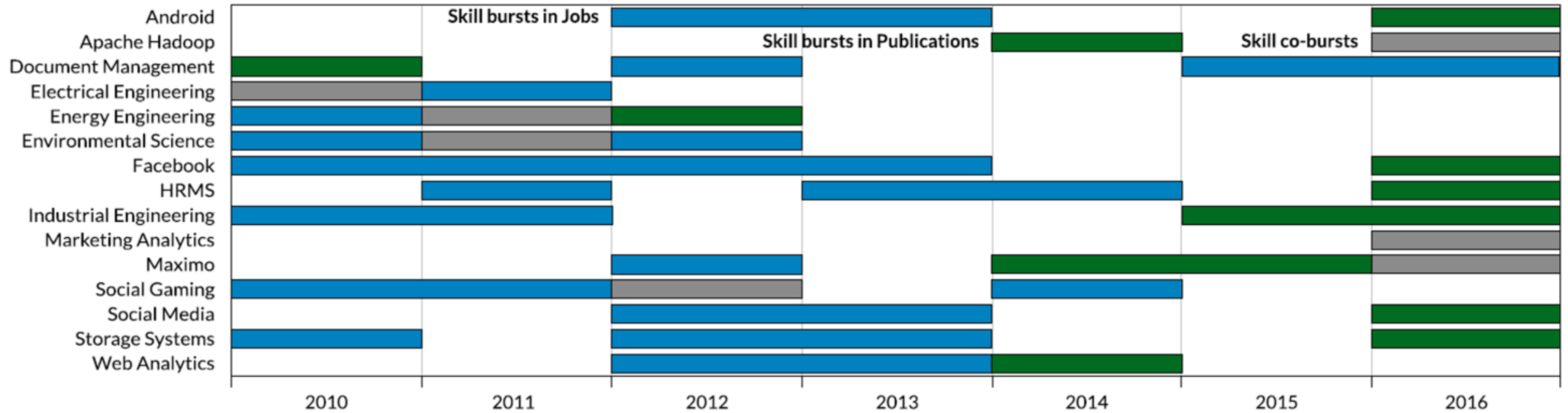




**Fig. 3.** Basemap of 13,218 skills with overlays of skill frequency in jobs, courses, and publications. This figure substantiates the conceptual drawing in Fig. 1 using millions of data records. Jobs skills are plotted in blue, courses are in red, and publications are in green. Node area size coding indicates base 10 log of skills frequency. The top 20 most frequent skills are labeled, and label sizes denote skill frequency.

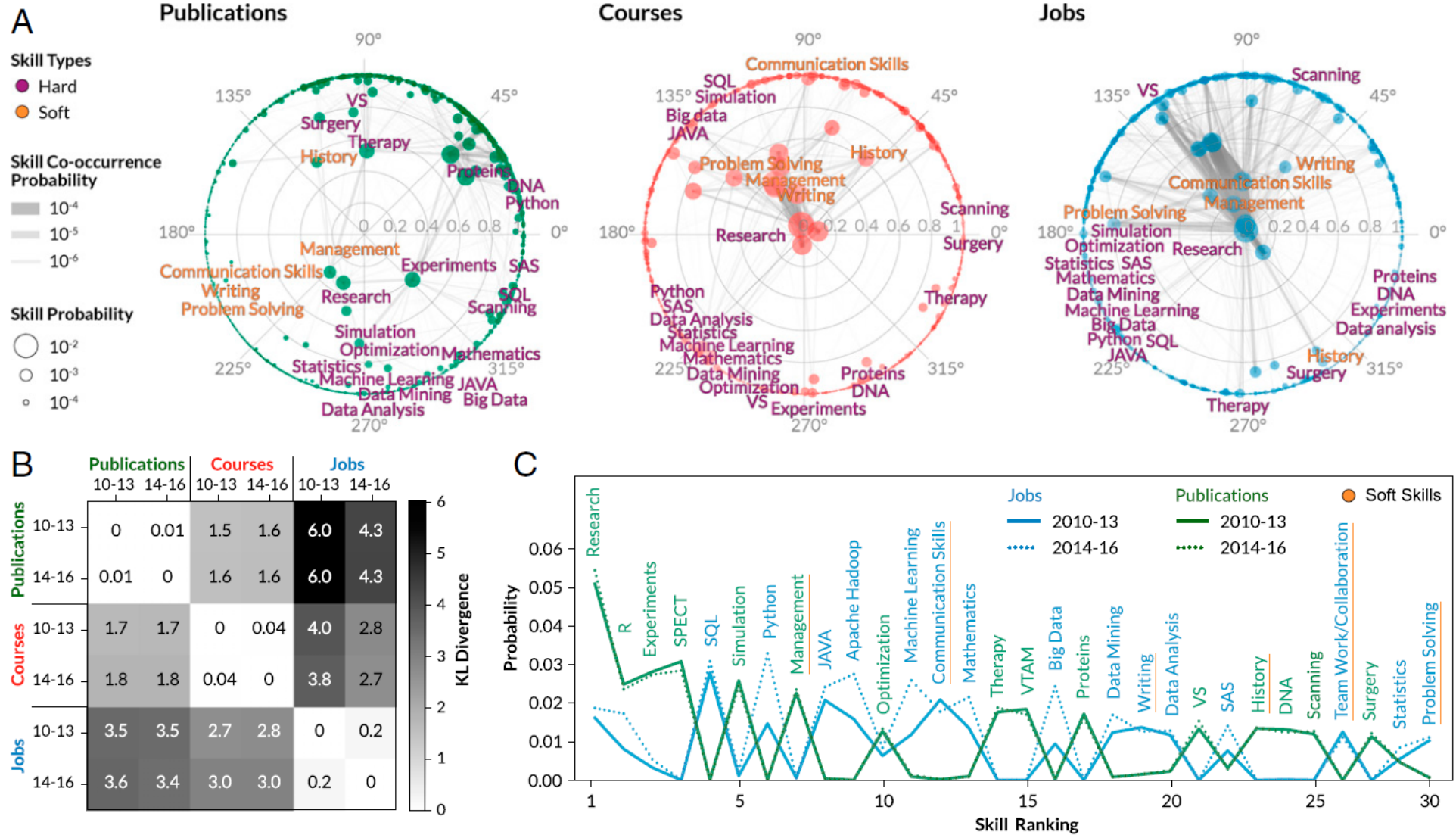






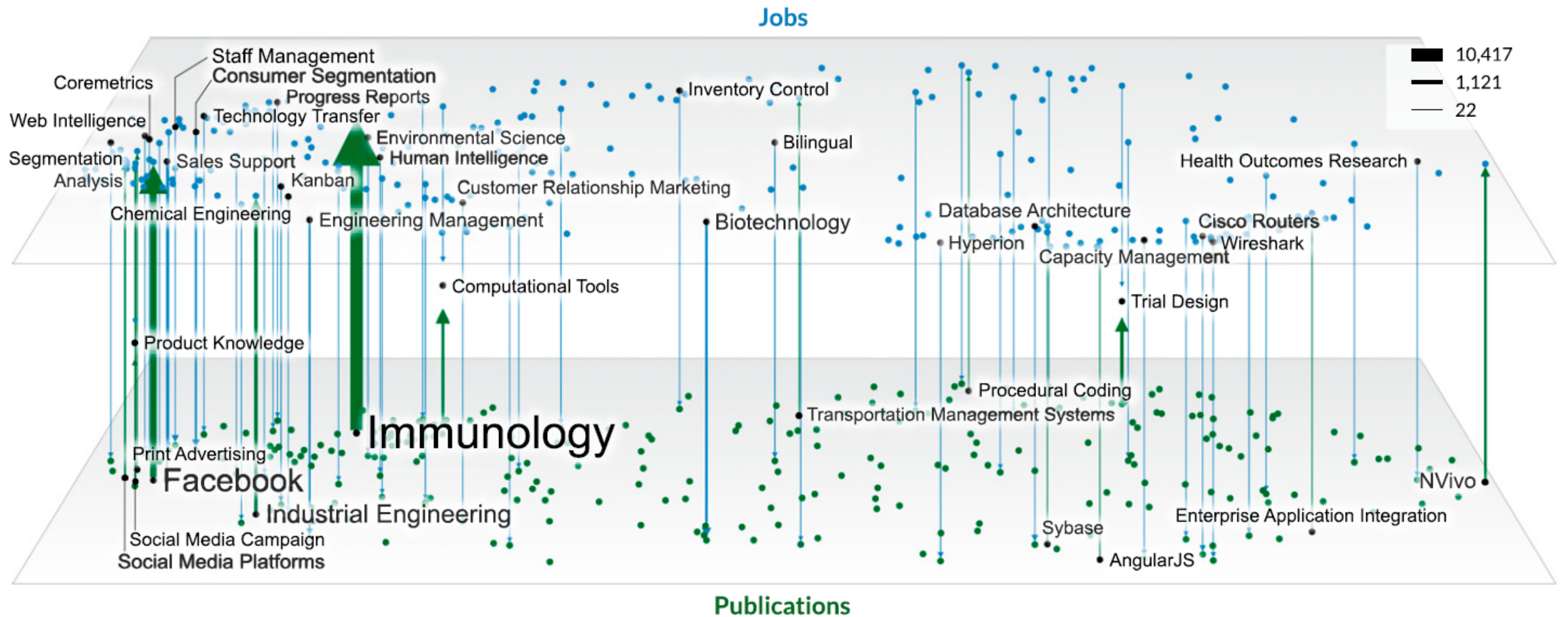
**Fig. 4.** Burst of activity in DS/DE skills in jobs and publications. Each burst is rendered as a horizontal bar with a start and an end date; skill term is shown on the left. Skills that burst in jobs are blue; skills bursting in publications are green. Seven skills burst in both datasets during the same years and are shown in gray. HRMS stands for human resources management system, and Maximo is an IBM system for managing physical assets.

# Kullback-Leibler divergence



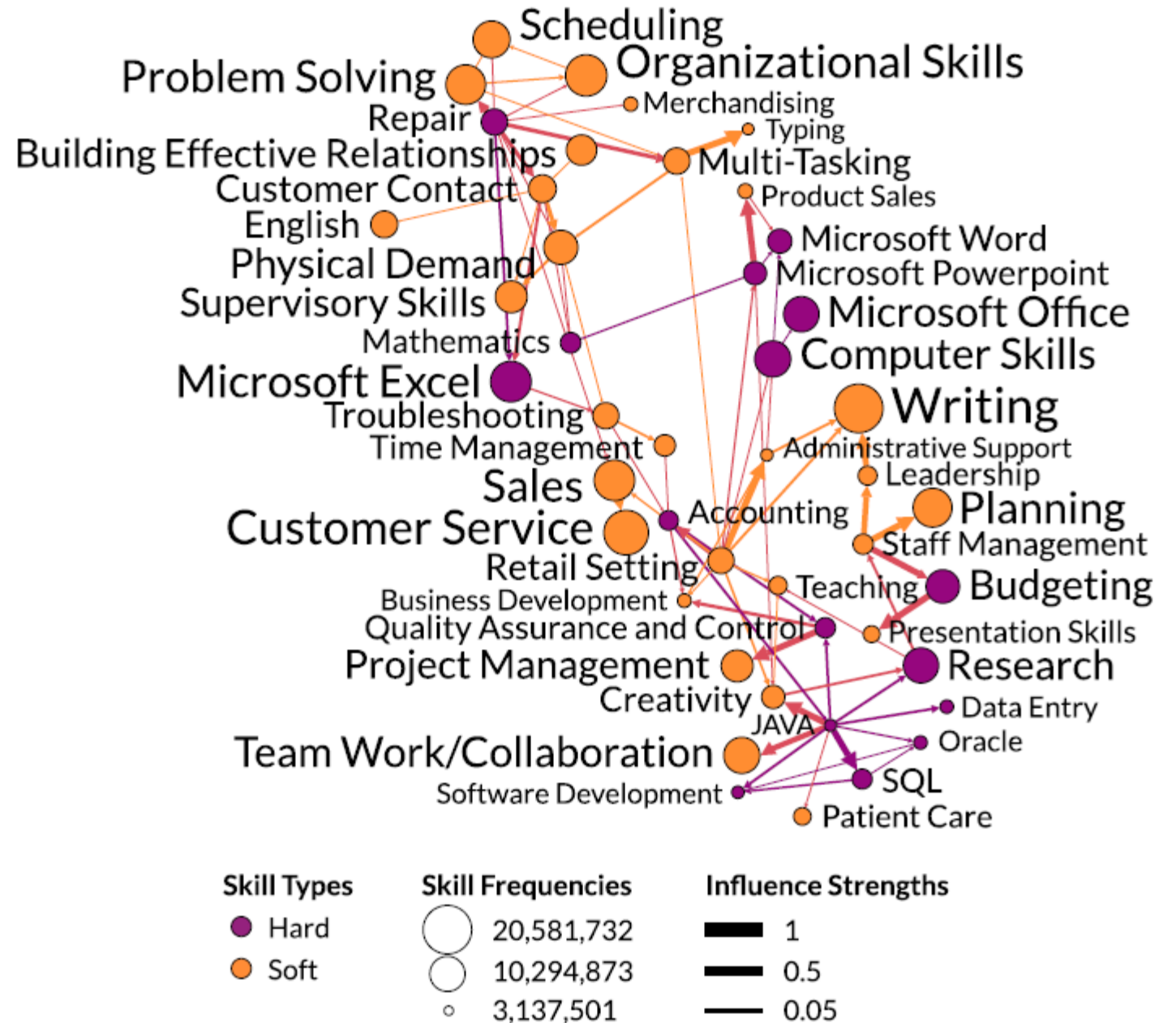
**Fig. 5.** Structural and dynamic differences between skill distributions in jobs, courses, and publications for 2010–2013 and 2014–2016. (A) Poincaré disks comparing the centrality of soft skills (orange) and hard skills (purple) across jobs, courses, and publications. (B) KL divergence matrix for jobs, courses, and publications in 2010–2013 and 2014–2016. (C) The most surprising skills in publications and jobs; *R* is a scripting language, VTAM refers to the IBM Virtual Telecommunication Access Method application, VS is the integrated development environment Visual Studio, and SAS is a data analytics software.





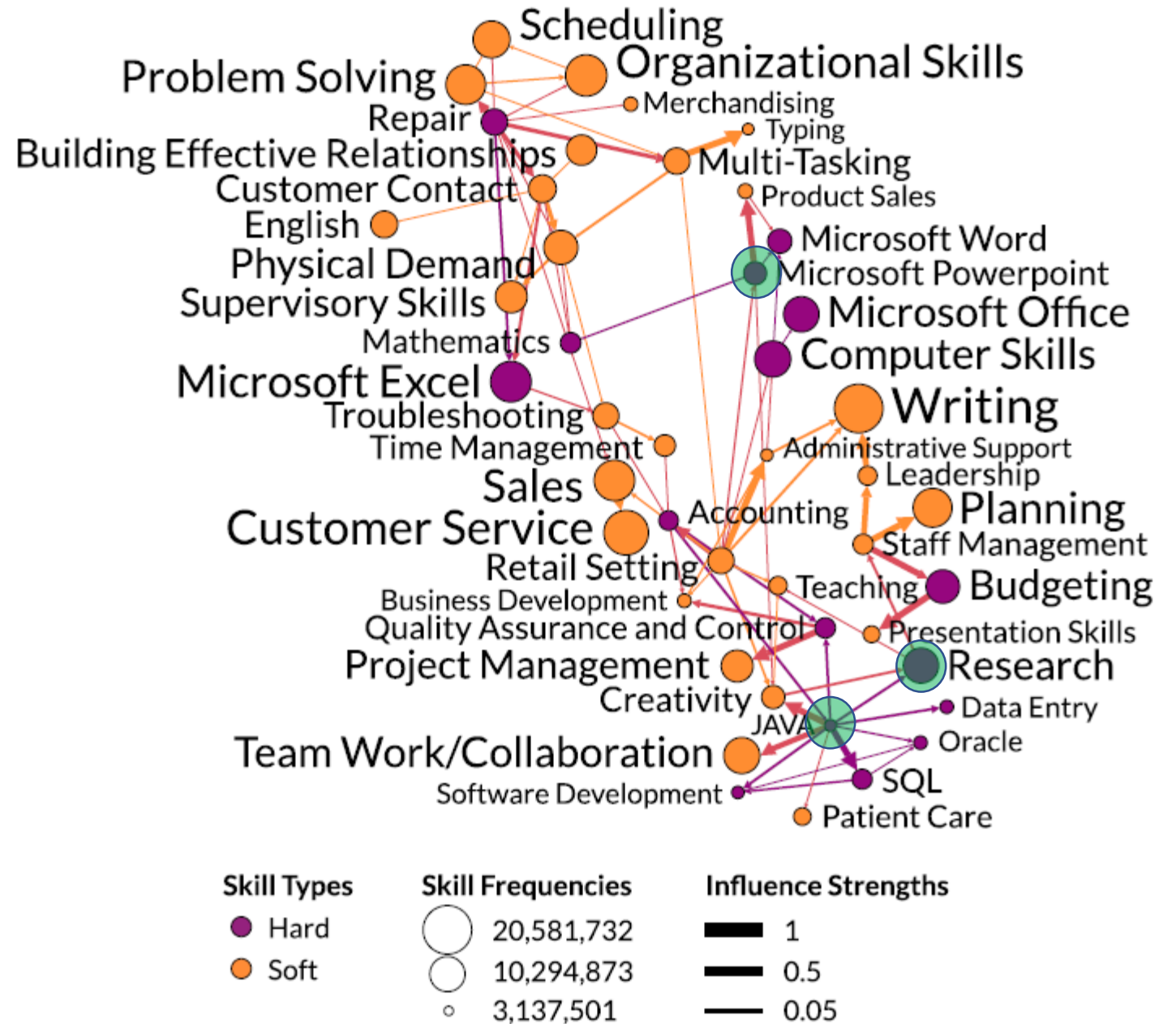
**Fig. 6.** Strength of influence mapping. Top 200 most frequent skills in jobs (blue) and in publications (green) plotted on the skills basemap from Fig. 2. Arrows represent skills with significant Granger causality ( $P$  value  $< 0.05$ ). Line thickness and label size indicate skill frequency. The direction and thickness of each arrow indicate the  $F$ -value strength and direction.

**Fig. 7.** Multivariate Hawkes Process influence network of DS/DE skills within job advertisements 2010–2016. Each of the 45 nodes represents a top-frequency skill (29 soft and 16 hard skills) with a strong influence edge from/to other skill(s) in job advertisements between 2010 and 2016. Node and label size correspond to the number of times that the skill appeared in a job advertisement. Thickness of the 75 directed edges indicates influence strength.





**Fig. 7.** Hawkes influence network of DS/DE skills within job advertisements 2010–2016. Each of the 45 nodes represents a top-frequency skill (29 soft and 16 hard skills) with a strong influence edge from/to other skill(s) in job advertisements between 2010 and 2016. Node and label size correspond to the number of times that the skill appeared in a job advertisement. Thickness of the 75 directed edges indicates influence strength.



# Results

- Novel cross-walk for mapping publications, course offerings, and job via skills.
- Timing and strength of burst of activity for skills (e.g., Oracle, Customer Service) in publications, course offerings, and job advertisements.
- Uniquely human skills such as communication, negotiation, and complex service provision are currently underexamined in research and undersupplied through education for the labor market in an increasingly automated and AI economy.
- The same pattern manifests in the domain of DS/DE where teamwork and communication skills increase in value with greater demand for data analytics skills and tools.
- Skill demands from industry are as likely to drive skill attention in research as the converse.



# 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](http://www.pnas.org/content/vol101/suppl_1)

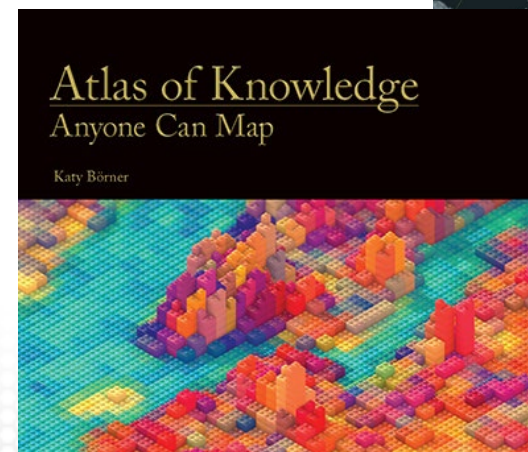
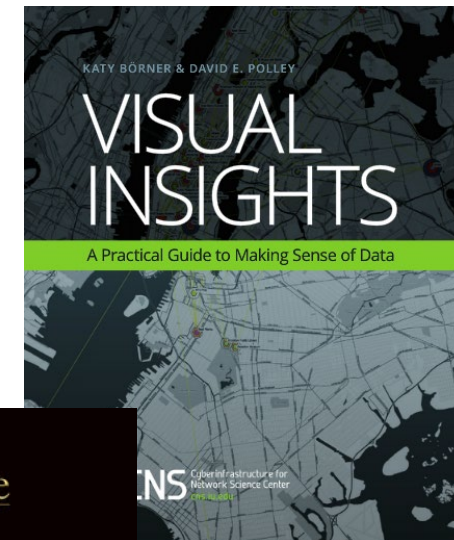
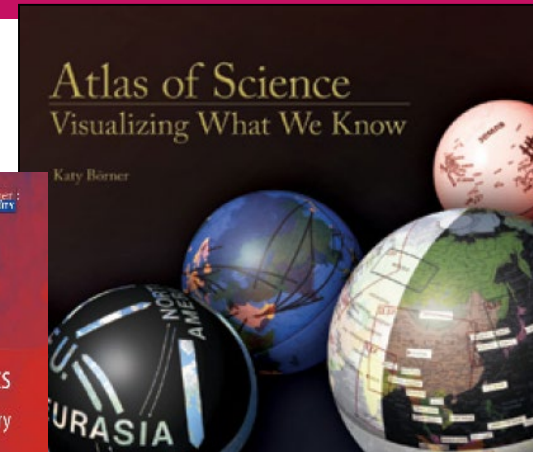
Börner, Katy (2010) **Atlas of Science: Visualizing What We Know**. The MIT Press. <http://scimaps.org/atlas>

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**. The MIT Press.

Börner, Katy (2015) **Atlas of Knowledge: Anyone Can Map**. The MIT Press. <http://scimaps.org/atlas2>





# Visual Analytics Certificate

Advance your skills in one of the most in demand careers through this six-week (6 CEUs) online course focused on understanding and creating data visualizations that translate complex data into actionable insights.

DOWNLOAD FLYER

REGISTER FOR OCT 7-NOV 17, 2019



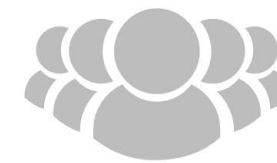
## Learn from Experts

Connect with industry professionals and leading researchers.



## Evolve Yourself

Gain forever knowledge and skill-up in powerful data visualization tools.



## Make a Difference

Embrace data-driven decision-making in your personal and professional life.

<https://visanalytics.cns.iu.edu>