

## Rete-netzwerk-red: analyzing and visualizing scholarly networks using the Network Workbench Tool

Katy Börner · Weixia Huang · Micah Linnemeier · Russell J. Duhon ·  
Patrick Phillips · Nianli Ma · Angela M. Zoss · Hanning Guo ·  
Mark A. Price

Received: 8 December 2009 / Published online: 6 January 2010  
© Akadémiai Kiadó, Budapest, Hungary 2010

**Abstract** The enormous increase in digital scholarly data and computing power combined with recent advances in text mining, linguistics, network science, and scientometrics make it possible to scientifically study the structure and evolution of science on a large scale. This paper discusses the challenges of this ‘BIG science of science’—also called ‘computational scientometrics’ research—in terms of data access, algorithm scalability, repeatability, as well as result communication and interpretation. It then introduces two infrastructures: (1) the Scholarly Database (SDB) (<http://sdb.slis.indiana.edu>), which provides free online access to 22 million scholarly records—papers, patents, and funding awards which can be cross-searched and downloaded as dumps, and (2) Scientometrics-relevant plug-ins of the open-source Network Workbench (NWB) Tool (<http://nwb.slis.indiana.edu>). The utility of these infrastructures is then exemplarily demonstrated in three studies: a comparison of the funding portfolios and co-investigator networks of different universities, an examination of paper-citation and co-author networks of major network science researchers, and an analysis of topic bursts in streams of text. The article concludes with a discussion of related work that aims to provide practically useful and theoretically grounded cyberinfrastructure in support of computational scientometrics research, education and practice.

**Keywords** Scientometrics · Science of science · Evolution of science ·  
Computational scientometrics · Data access · Algorithm scalability · Cyberinfrastructure ·  
Scholarly Database · Network Workbench · Related tools · Open source ·  
Open access

---

K. Börner (✉) · W. Huang · M. Linnemeier · R. J. Duhon · P. Phillips · N. Ma ·  
A. M. Zoss · H. Guo · M. A. Price  
Cyberinfrastructure for Network Science Center, School of Library and Information Science,  
Indiana University, Bloomington, IN, USA  
e-mail: [katy@indiana.edu](mailto:katy@indiana.edu)

M. Linnemeier  
e-mail: [mwlinnem@indiana.edu](mailto:mwlinnem@indiana.edu)

N. Ma  
e-mail: [nianma@indiana.edu](mailto:nianma@indiana.edu)

## Introduction

About 45 years ago, de Solla Price suggested studying science using scientific methods (de Solla Price 1963). Today, science of science studies draw from diverse fields such as scientometrics, informetrics, webometrics, history of science, sociology of science, psychology of the scientist, operational research on science, the economics of science, the analysis of the flow of scientific information, as well as the planning of science. They gather, handle, interpret, and predict a variety of features of science and technology enterprises, such as scholarly communication, performance, development, and dynamics that are interesting for science (policy) decisions in academia, government, and industry.

Most studies use either Thomson Reuters' Web of Science (hereafter WoS) or Elsevier's Scopus, as they each constitute a multi-disciplinary, objective, internally consistent publication database. A number of recent studies have examined and compared the coverage of WoS, Scopus, Ulrich's Directory, and Google Scholar (hereafter GS) (Meho and Yang 2007; Pauly and Stergiou 2005). It has been shown that the databases have a rather small overlap in records. In one study, the overlap between WoS and Scopus was only 58.2% while the overlap between GS and the union of WoS and Scopus was a mere 30.8%. While Scopus covers almost twice as many journals and conferences as WoS, it covers fewer journals in the arts and humanities. A comprehensive analysis requires access to more than one database, and a growing number of studies also correlate publication output with patent production, funding input, and other datasets.

While diverse tools exist to crawl, pre-process, analyze, or visualize scholarly data, most of the tools used in science of science studies today are proprietary or 'closed source', making it difficult or impossible to replicate results, to compare new and old approaches, or to agree upon standards.

Cyberinfrastructures, i.e., the programs, algorithms and computational resources required to support advanced data acquisition, storage, management, integration, visualization and analysis (Atkins et al. 2003), address the ever-growing need to connect researchers and practitioners to the data, algorithms, massive disk space and computing power that many computational sciences require (Emmott et al. 2006). Researchers need to make sense of massive amounts of data, knowledge, and expertise, and as industry needs assistance to overcome a major recession, access to an effective science of science cyberinfrastructure is highly desirable. The envisioned infrastructure would provide easy access to terabytes of scholarly data as well as advanced algorithms and tools running on powerful computers in support of 'BIG science of science' research—also called 'computational scientometrics', a term coined by C. Lee Giles (2006). Ideally, the infrastructure would be free, i.e., available to anyone, and open source, i.e., anybody could see, improve, and add to the software code. It needs to support the effective communication and discussion of results and should provide means to record analysis workflows so that others can rerun analyses and replicate results.

This paper introduces the beginnings of such a science of science cyberinfrastructure: the Scholarly Database (SDB) and the Network Workbench (NWB) Tool. It starts with a general introduction of the system architecture, functionality, and the user interfaces of SDB and NWB. It then demonstrates their utility in three original research studies that exemplify common workflows for the acquisition and preparation of bibliographic data, temporal data analysis, network analysis and visualization. The paper concludes with a discussion of related work.

### Scholarly Database

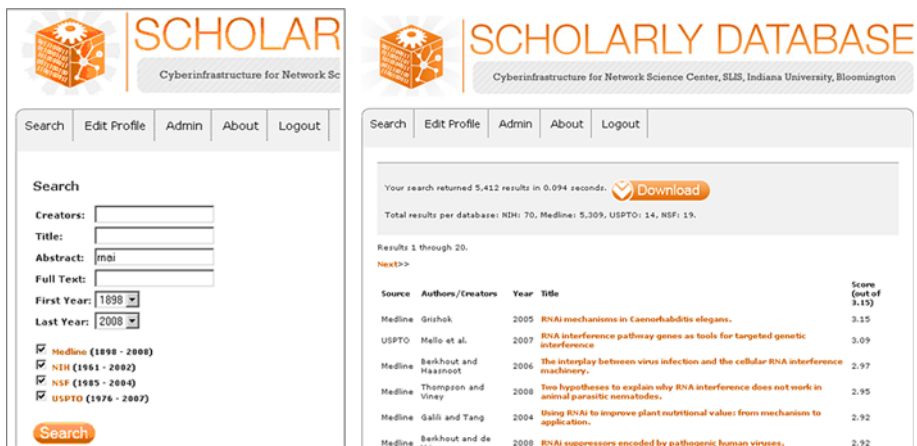
The Scholarly Database (hereafter SDB) at Indiana University aims to serve researchers and practitioners interested in the analysis, modelling, and visualization of large-scale scholarly datasets. The motivation for this database and its previous implementation was presented in (LaRowe et al. 2009). The online interface at <http://sdb.slis.indiana.edu> provides access to four datasets: Medline papers, U.S. Patent and Trademark Office patents (USPTO), National Science Foundation (NSF) funding, and National Institutes of Health (NIH) funding—over 23 million records in total, see Table 1. Users can register for free to cross-search these databases and to download result sets as dumps for scientometrics research and science policy practice.

SDB supports search across paper, patent, and funding databases. To initiate a search, enter the search term(s) into creators (author/awardee/inventor), title, abstract, and full text (keywords and other text) fields, select a year range and database(s); see Fig. 1 (left side of the figure).

The importance of a particular term in a query can be increased by putting a ^ and a number after the term. For instance, ‘breast cancer^10’ would increase the importance of matching the term ‘cancer’ by ten compared to matching the term ‘breast’. Custom database queries can be run by contacting the SDB team lead and author of this paper,

**Table 1** Names, number of records, years covered, and update information for datasets currently available via the Scholarly Database

Dataset name	# Records	Years covered	Regular update
Medline Publications	17,764,826	1865–2008	Yes
U.S. Patent and Trademark Office Patents	3,875,694	1976–2008	Yes
National Institutes of Health (NIH) Awards	1,043,804	1961–2002	No
National Science Foundation (NSF) Awards	174,835	1985–2004	No



**Fig. 1** Partial search interface (left) and download interface (right) for the Scholarly Database. Note that the highest scoring record was retrieved from Medline while the second highest record comes from USPTO

Nianli Ma. Search results retrieved from different databases can be downloaded as data dumps in csv file format; see Fig. 1 (right).

SDB stores all data in a PostgreSQL database (PostgreSQL Global Development Group 2009). Full-text search is supported using Solr (The Apache Software Foundation 2007) to index the contents of the database. Solr is an industry-standard, open source search server that can scale to very large amounts of data using replication and sharding. The online interface was developed in Django (Django Software Foundation 2009). Django is a web framework written in the Python (Python Software Foundation 2008) programming language with particularly good support for content-oriented web applications.

### Network Workbench Tool

The Network Workbench (NWB) Tool (<http://nwb.slis.indiana.edu>) is a network analysis, modelling, and visualization toolkit for physics, biomedical, and social science research (Herr et al. 2007). The basic interface comprises a ‘Console’, ‘Data Manager’, and ‘Scheduler’ Window as shown in Fig. 2. The top menu provides easy access to relevant ‘Preprocessing’, ‘Modeling’, ‘Analysis’, ‘Visualization’, and ‘Scientometrics’ algorithms. Information on how to download, install, and run the NWB Tool can be found in the Network Workbench Tool User Manual (Cyberinfrastructure for Network Science Center 2009).

NWB is built on Cyberinfrastructure Shell (CIShell) (Cyberinfrastructure for Network Science Center 2008), an open source software framework for the easy integration and

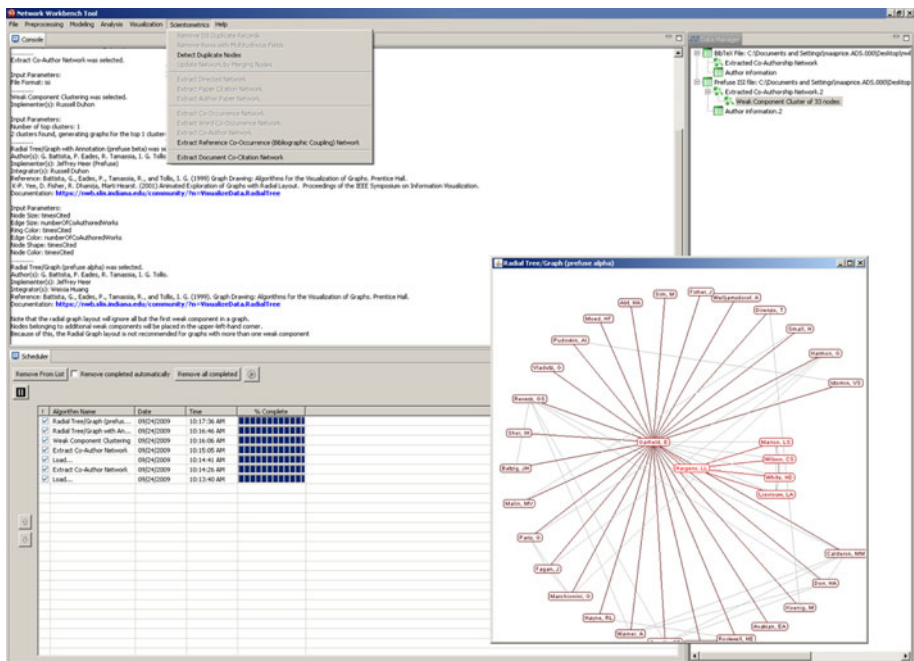


Fig. 2 NWB Tool interface with Console, Data Manager, and Scheduler Windows, the Scientometrics menu, and a Radial Graph visualization of Garfield’s co-author network

utilization of datasets, algorithms, tools, and computing resources. CIShell is based on the OSGi R4 Specification and Equinox implementation (OSGi Alliance 2008).

The Network Workbench Community Wiki (<https://nwb.slis.indiana.edu/community>) provides a one-stop online portal for researchers, educators, and practitioners interested in the study of networks. It is a place for users of the NWB Tool, CIShell, or any other CIShell-based program to get, upload, and request algorithms and datasets to be used in their tool so that it truly meets their needs and the needs of the scientific community at large.

Users of the NWB tool can

- Load sample datasets or their own networks and formats.
- Perform network analysis with some of the most effective algorithms available.
- Generate, run, and validate network models.
- Use different visualizations to interactively explore and understand specific networks.
- Share datasets and algorithms across scientific boundaries.

As of August, 2009, the NWB Tool provides access to over 80 algorithms and 30 sample datasets for the study of networks. The loading, processing, and saving of seven file formats (NWB, GraphML, Pajek.net, Pajek.matrix, XGMML, TreeML, CSV) and an automatic conversion service among those formats is supported. Relevant for science of science studies, the NWB Tool can read data downloaded from SDB, Google Scholar, Thomson Reuters ISI Scientific, Scopus, and the NSF award database as well as EndNote and BibTeX formatted data.

Additional algorithms and data formats can be easily integrated into the NWB Tool using wizard-driven templates. Although the CIShell and the NWB Tool are developed in JAVA, algorithms developed in other programming languages such as FORTRAN, C, and C++ can be integrated. Among others, JUNG (O'Madadhain et al. 2008) and Prefuse libraries (Heer et al. 2005) have been integrated into the NWB as plug-ins. NWB also supplies a plug-in that invokes the GnuPlot application (Williams and Kelley 2008) for plotting data analysis results and the GUESS tool (Adar 2007) for rendering network layouts. Support and advice in algorithm integration and custom tool development is provided by the NWB team lead and author of this paper, Micah Linnemeier.

## Exemplary workflows

This section aims to demonstrate the utility of the SDB and NWB Tool to answer specific research questions in an efficient and repeatable fashion. Detailed, step-by-step instructions on how to run these and many other analyses can be found in the Network Workbench Tool User Manual (Cyberinfrastructure for Network Science Center 2009) and NWB Tutorial Slides (Börner 2008).

NSF funding portfolios and co-investigator networks of U.S. universities

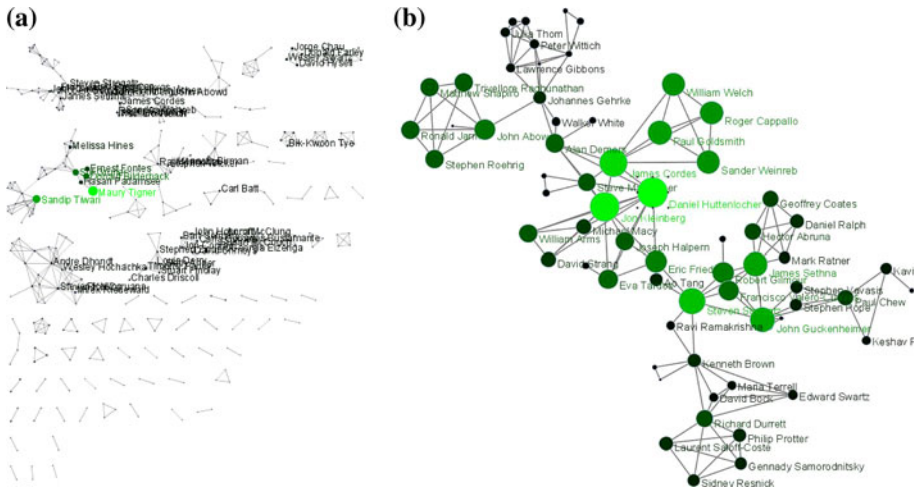
The first study aims to answer: What active funding portfolios do major U.S. universities have, what scholarly co-investigator networks does this funding inspire/support, and what roles do investigators play, e.g., gatekeeper, using betweenness centrality measures (BC) (Freeman 1977), number of collaborators via node degree, total funding amount?

Funding data was downloaded from the *Award Search* site provided by the National Science Foundation (NSF) (<http://www.nsf.gov/awardsearch>). The site supports search by

PI name, institution, and many other fields. Exemplarily, active NSF awards data from Indiana University, Cornell University, and University of Michigan Ann Arbor were downloaded on November 07, 2008. The files were loaded into the NWB Tool and co-investigator networks were extracted and visualized in GUESS. In these networks, nodes represent investigators and edges denote their co-occurrence as an award, i.e., co-authorship. The co-investigator network of Cornell University without isolate nodes is shown in Fig. 3a, left. The largest connected component of this network from the top left corner of Fig. 3a is shown in Fig. 3b, right. In both networks, the node area size and color corresponds to the total award amount with smaller, darker nodes representing less money and larger, light green nodes denoting more funding and the top-50 nodes with the highest funding amounts are labeled.

The general characteristics of all three networks are given in Table 2. Note that the total award amount is attributed to the main investigator and his/her institution exclusively.

There are interesting differences in the funding portfolios of these universities. Michigan has clearly the largest number of currently active NSF awards totalling 497. With \$546 million, Cornell has the highest total award amount. Cornell also has the largest giant



**Fig. 3** Complete network (a) and largest component (b) of Cornell University's co-investigator network (67 nodes)

**Table 2** Award properties and co-investigator network features for the three universities

University	#Records/total award amount	Co-investigator network: #nodes/#edges/# components	Largest component: #nodes/#edges	Name, department, and \$ amount of top investigator
Indiana University	257/\$100 million	223/312/52	19/37	Curtis Lively, Biology \$7,436,828
Cornell University	501/\$546 million	375/573/78	67/143	Maury Tigner, Physics \$107,216,976
Michigan University	619/\$305 million	497/672/117	55/105	Khalil Najafi, EECS \$32,541,158

component with 67 investigator nodes and 143 collaboration links, indicating much cross-fertilization across different disciplines. Cornell also happens to employ the investigator who currently has the highest total award amount: Maury Tigner. Note that being the main investigator on one major center grant and several campus equipment grants can easily result in multi-millions to spend over many years. Note also that the funding portfolios, networks, and top-investigators from other agencies, e.g., NIH, might look very different.

A closer examination of the largest component of the Cornell co-investigator network shown in Fig. 3b, reveals that Steven Strogatz has the highest betweenness centrality (BC), effectively bridging between several disciplines, and Daniel Huttenlocher has the highest degree, i.e., the most collaborations with others in this network.

Future work should consider different means to associate award amounts to investigators and institutions. An analysis of the distribution of funding over scientific disciplines and departments is desirable. Co-investigator linkages among institutions deserve further attention.

### Paper-citation and co-author networks of major network science researchers

The second study addresses the questions: Do researchers which come from different domains of science but make major contributions to one and the same domain, e.g., network science, grow different collaboration networks? How much do their publication, citation, and H-index dynamics differ?

Exemplarily, four major network science researchers were selected: Eugene Garfield and Stanley Wasserman, Alessandro Vespignani and Albert-László Barabási of the Network Workbench project: Data for all four male researchers was downloaded from Thomson Reuters in December 2007. Their names, research domains, ages (retrieved from Library of Congress), number of citations for highest cited paper, H-index (Bornmann 2006), and number of papers and citations over time as calculated by Thomson Reuters' WoS (Thomson Reuters 2009) are given in Table 3. Note that this dataset does not capture any books or Conference proceedings by these authors.

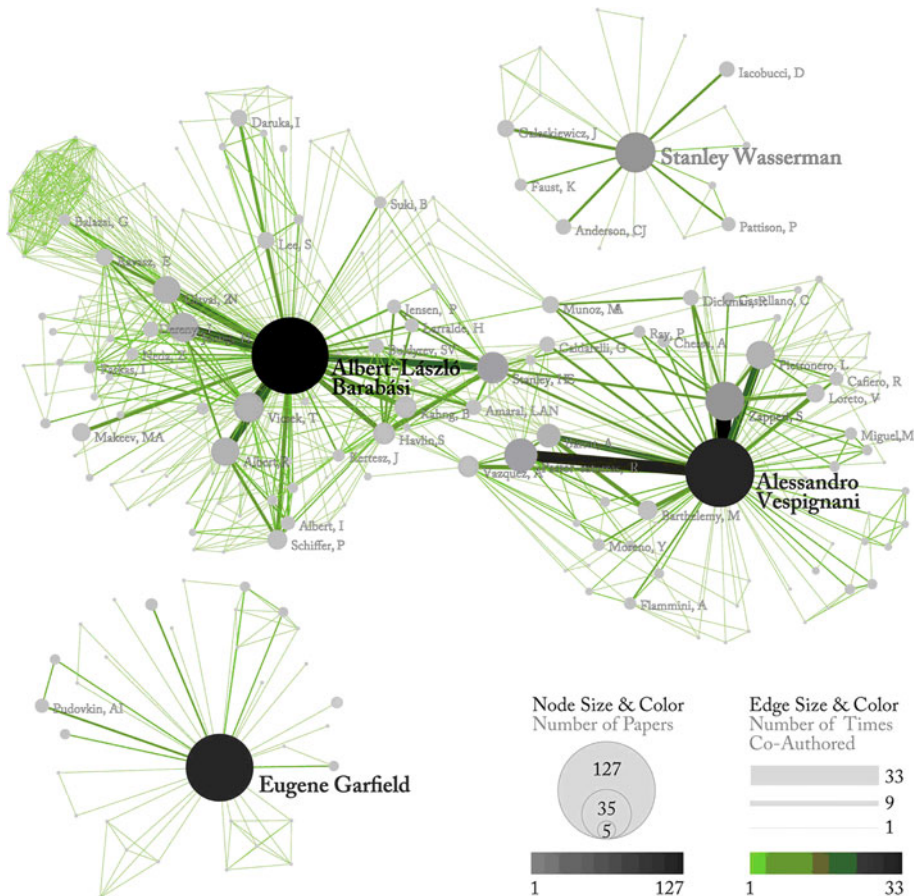
The older an author the more papers, citations, and the higher an H-index are expected. Yet, Vespignani and Barabási publishing in physics and biology manage to attract citation counts and have H-indexes that are impossible to achieve in social science domains in such a short time frame. To give a concrete example, in Dec. 2007, Garfield's highest cited paper on "Citation Analysis as a Tool in Journal Evaluation" published in 1972 had 672 counts. Barabási's highest cited paper, "Emergence of Scaling in Random Networks," published in 1999, has 2,218 counts; in December 2008 the same paper has 3,488 citation counts. Within one single year, Barabási's H-index increased by 5 to 52.

**Table 3** Names, age, number of papers, number of citations for highest cited paper, and *H*-index for four major network science researchers

Author name	Department	Age in 2007	# Papers	# Cites	<i>H</i> -Index
Eugene Garfield	IS, Scientometrics	82	672	672	31
Stanley Wasserman	Sociology, Psychology, Statistics	56	35	122	17
Alessandro Vespignani	Physics, IS	42	101	451	33
Albert-László Barabási	Physics, Biology	40	126	2,218	47 (Dec 07)
		41	159	3,488	52 (Dec 08)

Similarly, there are major differences in the structure of the collaboration networks in which these four authors are embedded in. Figure 4 shows the joint co-author network of all records retrieved for the four authors as rendered in GUESS. Each node represents an author and is color and size coded by the number of papers per author. Edges represent co-author relationships and are color and thickness coded by the number of times two authors wrote a paper together. The top-50 authors with the most papers are labelled.

While Barabási's and Vespignani's co-author networks are strongly interlinked with Stanley and Vazquez as major connectors with high betweenness centrality values, Garfield's and Wasserman's networks are unconnected to the networks of any of the three other researchers. The size and density of the networks differs considerably. When extracting the co-author network from all of Barabási's papers, a co-author network with 128 nodes and density of 0.075 results; for Vespignani the network has 72 nodes and 0.086 density, Wasserman has 18 and 0.202, and Garfield has 33 and 0.119. The top-3 strongest co-author linkages for Barabási are: Vicsek, Jeong, Albert; for Vespignani: Zapperi, Pastor-Satorras, Pietronero; for Wasserman: Galaskiewicz, Iacobucci, Anderson; and for Garfield: Pudovkin, Welljamsdorof, and Sher is tied for third place with four other authors. The paper with the most authors is entitled "Experimental determination and system level



**Fig. 4** Joint co-author network of all four network science experts



**Table 4** Words and their burst length and weight as well as start and end years

Word	Length	Weight	Start	End
Elegan	5	105.4307	1998	2002
Strand	4	70.72182	1999	2002
Doubl	5	62.88262	1999	2003
Embryo	6	42.70386	1998	2003
Caenorhabd	5	36.59776	1998	2002
Drosophila	5	33.98245	1999	2003
Phenotyp	3	31.08153	1999	2001

analysis of essential genes in *Escherichia coli* MG1655" (2003) and has 21 authors, one of them being Albert-László Barabási.

### Bursts of activity in streams of text

The third study aims to answer: What topic bursts exist in an emerging area of research, e.g., RNA interference (RNAi) research? Exactly what topics are active and when?

Using the Scholarly Database, 5,319 Medline papers with 'rna' in the abstract field were retrieved and downloaded; see Fig. 1, left, for interface snapshots. The first paper was published in 1978 and the number of papers on this topic increases considerably in recent years. The sequence of papers published over time can be seen as a discrete time series of words. Kleinberg's burst detection algorithm (Kleinberg 2002) identifies sudden increases in the usage frequency of words over time. Rather than using plain frequencies of the occurrences of words, the algorithm employs a probabilistic automaton whose states correspond to the frequencies of individual words. State transitions correspond to points in time around which the frequency of a word changes significantly. The algorithm generates a ranked list of the word bursts in the document stream, together with the intervals of time in which they occurred. Using the burst detection algorithm available via the NWB Tool, the abstracts of the 5,309 Medline papers were analyzed. First, all words occurring in the abstract were normalized: they were stemmed, i.e., words such as 'scientific', 'science' were reduced to 'scien' and stop words such as 'a' or 'the' were removed. The result is a 1,224 row table with all bursting words, their burst length, weight, and the start and end years of the bursts. The table was sorted by burst weight and four words: 'protein', 'result', 'use', and 'function' top the list with an infinite burst weight. The subsequent top-7 most highly bursting words are given in Table 4. As the burst detection algorithm was run with 'bursting state = 1', i.e., only one burst per word, the burst weight is identical to the burst strength in this output and only the burst weight is shown in Table 4. Interestingly, all these words burst rather early in 1998 or 1999. Many of them have a rather long burst duration with 'embryo' bursting over 6 years.

The result shows the words and topics important to the events being studied that increased in usage, were more active for a period of time, and then faded away.

### Related work and discussion

A discussion of the unique features of the Scholarly Database and its relation to similar efforts was provided in (LaRowe et al. 2009). Here features and related work of the NWB Tool are discussed. Table 5 provides an overview of existing tools used in scientometrics

**Table 5** Network analysis and visualization tools commonly used in scientometrics research

Tool	Year	Domain	Description	Open source	Operating system	References
S&T Dynam. Toolbox	1985	Scientom.	Tools from Loet Leydesdorff for organization analysis, and visualization of scholarly data	No	Windows	(Leydesdorff 2008)
In Flow	1987	SocSci	Social network analysis software for organizations with support for what-if analysis	No	Windows	(Krebs 2008)
Pajek	1996	SocSci	A network analysis and visualization program with many analysis algorithms, particularly for social network analysis	No	Windows	(Batagelj and Mrvar 1998)
BibExcel	2000	Scientom	Transforms bibliographic data into forms usable in Excel, Pajek, NetDraw, and other programs	No	Windows	(Persson 2008)
Boost Graph Library	2000	CS	Extremely efficient and flexible C++ library for large networks	Yes	All Major	(Siek et al. 2002)
UCINET	2000	SocSci	Social network analysis software particularly useful for exploratory analysis	No	Windows	(Borgatti et al. 2002)
Visone	2001	SocSci	Social network analysis tool for research and teaching, with a focus on innovative and advanced visual methods	No	All Major	(Brandes and Wagner 2008)
Cytoscape	2002	Bio	Network visualization and analysis tool focusing on biological networks, with well-designed visualizations	Yes	All Major	(Cytoscape Consortium 2008)
GeoVISTA	2002	Geo	GIS software that can be used to lay out networks on geospatial substrates	Yes	All Major	(Takatsuka and Gabegan 2002)
iGraph	2003	CS	A library for classic and cutting edge network analysis usable with many programming languages	Yes	All Major	(Csárdi and Nepusz 2006)
Tulip	2003	CS	Graph visualization software for networks over 1,000,000 elements	Yes	All Major	(Auber 2003)
CiteSpace	2004	Scientom	A tool to analyze and visualize scientific literature and frontiers, particularly co-citation structures	Yes	All Major	(Chen 2006)
GraphViz	2004	Networks	Flexible graph visualization software	Yes	All Major	(AT&T Research Group 2008)
Hittite	2004	Scientom	Analysis and visualization tool for data from the Web of Science	No	Windows	(Garfield 2008)

**Table 5** continued

Tool	Year	Domain	Description	Open source	Operating system	References
R	2004	Statistics	A statistical computing language with many libraries for sophisticated network analyses	Yes	All Major	(Ihaka and Gentleman 1996)
Prefuse	2005	Visualiz.	A general visualization framework with many capabilities to support network visualization and analysis	Yes	All Major	(Heer et al. 2005)
NWB Tool	2006	Bio, IS, SocSci, Scientom	Network analysis & visualization tool conducive to new algorithms supportive of many data formats	Yes	All Major	(Huang et al. 2007)
GUESS	2007	Networks	A tool for visual graph exploration that integrates a scripting environment	Yes	All Major	(Adar 2007)
Publish or Perish	2007	Scientom	Harvests and analyzes data from Google Scholar, focusing on measures of research impact	No	Windows, Linux	(Harzing 2008)
VOSviewer	2009	Scientom	Constructs and views maps based on co-occurrence data	No	All Major	(van Eck and Waltman 2009)

research; see also (Börner et al. 2007; Fekete and Börner-chairs 2004). The tools are sorted by the date of their creation. Domain refers to the field in which they were originally developed, such as social science (SocSci), scientometrics (Scientom), biology (Bio), geography (Geo), and computer science (CS).

Many of these tools are specialized and very capable. For instance, BibExcel and Publish or Perish are great tools for bibliometric data acquisition and analysis. HistCite and CiteSpace each support very specific insight needs—from studying the history of science to the identification of scientific research frontiers. The S&T Dynamics Toolbox provides many algorithms commonly used in scientometrics research and it provides bridges to more general tools. Pajek and UCINET are very versatile, powerful network analysis tools that are widely used in social network analysis. Cytoscape is excellent for working with biological data and visualizing networks.

The NWB Tool has fewer analysis algorithms than Pajek and UCINET and simpler visualizations than Cytoscape. However, it is open source, highly flexible, and scalable to very large networks. NWB Tool makes it much easier for researchers and algorithm authors to integrate new and existing algorithms and tools that take in diverse data formats. This is made possible by the OSGi component architecture and CShell algorithm architecture built on top of OSGi.

The Cytoscape team recently decided to adopt an architecture based on OSGi. This will make it possible for Cytoscape to use many of the NWB analysis and modelling algorithms while the NWB Tool can benefit from Cytoscape's visualization capabilities. Other software development teams are exploring an adoption of OSGi. Ultimately, a true marketplace-like cyberinfrastructure might result that makes it easy to share and use datasets, algorithms, and tools across scientific boundaries. Ideally, this marketplace is free for anybody to use and contribute to, enabling and harnessing the power of millions of minds for studies in biology, physics, social science, and many other disciplines but also for the study of science itself (Börner, submitted).

**Acknowledgments** We would like to acknowledge the contributions and support by the NWB team and advisory board. This work is funded by the School of Library and Information Science and the Cyberinfrastructure for Network Science Center at Indiana University, the James S. McDonnell Foundation, and the National Science Foundation under Grants No. IIS-0715303, IIS-0534909, and IIS-0513650. 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.

## References

- Adar, E. (2007). *Guess: The graph exploration system*. Accessed April 22, 2008, from <http://graphexploration.cond.org/>.
- AT&T Research Group. (2008). *Graphviz-graph visualization software*. Accessed July 17, 2008, from <http://www.graphviz.org/Credits.php>.
- Atkins, D. E., Drogemeier, K. K., Feldman, S. I., Garcia-Molina, H., Klein, M. L., Messerschmitt, D. G., et al. (2003). *Revolutionizing science and engineering through cyberinfrastructure: Report of the national science foundation blue-ribbon advisory panel on cyberinfrastructure*. Arlington: National Science Foundation.
- Auber, D. (2003). Tulip: A huge graph visualisation framework. In P. Mutzel & M. Jünger (Eds.), *Graph drawing softwares, mathematics and visualization* (pp. 105–126). Berlin: Springer-Verlag.
- Batagelj, V., & Mrvar, A. (1998). Pajek-program for large network analysis. *Connections*, 21(2), 47–57.
- Borgatti, S. P., Everett, M. G., & Freeman, L. C. (2002). *Ucinet for windows: Software for social network analysis*. Accessed July 15, 2008, from [http://www.analytictech.com/ucinet/ucinet\\_5\\_description.htm](http://www.analytictech.com/ucinet/ucinet_5_description.htm).
- Börner, K. (2008). *Network Workbench Tool: For large scale network analysis, modeling, and visualization*. (unpublished) <http://ivl.slis.indiana.edu/km/pres/2008-borner-nwb-ws.pdf>.

- Börner, K. (submitted). Plug-and-Play macroscopes. *Communications of the ACM*.
- Börner, K., Sanyal, S., & Vespignani, A. (2007). Network science. In B. Cronin (Ed.), *Annual review of information science & technology (ARIST)* (Vol. 41, pp. 537–607). Medford, NJ: Information Today, Inc./American Society for Information Science and Technology.
- Bormmann, L. (2006). *H index: A new measure to quantify the research output of individual scientists*. Accessed July 17, 2008, from [http://www.forschungsinform.de/iq/agora/H\\_Index/h\\_index.asp](http://www.forschungsinform.de/iq/agora/H_Index/h_index.asp).
- Brandes, U., & Wagner, D. (2008). *Analysis and visualization of social networks*. Accessed July 15, 2008, from <http://visone.info/>.
- Chen, C. (2006). CiteSpace II: Detecting and visualizing emerging trends and transient patterns in scientific literature. *Journal of the American Society for Information Science and Technology*, 54(5), 359–377.
- Csárdi, G., & Nepusz, T. (2006). *The igraph software package for complex network research*. Accessed July 17, 2008, from <http://necsi.org/events/iccs6/papers/c1602a3c126ba822d0bc4293371c.pdf>.
- Cyberinfrastructure for Network Science Center. (2008). *Cyberinfrastructure shell*. Accessed July 17, 2008, from <http://cishell.org/>.
- Cyberinfrastructure for Network Science Center. (2009). *Network Workbench Tool: User manual, 1.0.0. (9/16)*. Accessed September 23, 2009, from <http://nwb.slis.indiana.edu/Docs/NWBTTool-Manual.pdf>.
- Cytoscape Consortium. (2008). *Cytoscape*. Accessed September 14, 2008, from <http://www.cytoscape.org/index.php>.
- de Solla Price, D. J. (1963). *Little science, big science*. Unpublished Manuscript.
- Django Software Foundation. (2009). *Django: The web framework for perfectionists with deadlines*. Accessed January 13, 2008, from <http://www.djangoproject.com/contact/foundation/>.
- Emmott, S., Rison, S., Abiteboul, S., Bishop, C., Blakeley, J., Brun, R., et al. (2006). *Towards 2020 science. The Microsoft Research Group and the 2020 Science Group*. Accessed January 13, 2008, from [http://research.microsoft.com/en-us/um/cambridge/projects/towards2020science/downloads/T2020S\\_ReportA4.pdf](http://research.microsoft.com/en-us/um/cambridge/projects/towards2020science/downloads/T2020S_ReportA4.pdf).
- Fekete, J.-D., & Börner-chairs, K. (2004). Workshop on information visualization software infrastructures Austin, Texas.
- Freeman, L. C. (1977). A set of measuring centrality based on betweenness. *Sociometry*, 40, 35–41.
- Garfield, E. (2008). *HistCite: bibliometric analysis and visualization software (Version 8.5.26)*. Bala Cynwyd, PA: HistCite Software LLC. Accessed July 15, 2008, from <http://www.histcite.com/>.
- Giles, C. L. (2006). *The future of CiteSeer: CiteSeerx*, Lecture Notes in Computer Science (Vol. 4213). Berlin/Heidelberg: Springer.
- Harzing, A.-W. (2008). *Publish or Perish: A citation analysis software program*. Accessed April 22, 2008, from <http://www.harzing.com/resources.htm>.
- Heer, J., Card, S. K., & Landay, J. A. (2005). Prefuse: A toolkit for interactive information visualization. *Conference on human factors in computing systems* (pp. 421–430). Portland, OR/New York: ACM Press.
- Herr II, Bruce W., Huang, W. (Bonnie), Penumathy, S., & Börner, K. (2007). Designing highly flexible and usable cyberinfrastructures for convergence. In W. S. Bainbridge & M. C. Roco (Eds.), *Progress in convergence: Technologies for human wellbeing* (Vol. 1093, pp. 161–179). Boston, MA: Annals of the New York Academy of Sciences.
- Huang, W. (Bonnie), Herr, B., Duhon, R., & Börner, K. (2007). Network Workbench—Using service-oriented architecture and component-based development to build a tool for network scientists. Presented at International Workshop and Conference on Network Science, Queens, NY.
- Ihaka, R., & Gentleman, R. (1996). R: A language for data analysis and graphics. *Journal of Computational and Graphical Statistics*, 5(3), 299–314. Accessed July 17, 2008, from <http://www.amstat.org/publications/jcgs/>.
- Kleinberg, J. M. (2002). Bursty and hierarchical structure in streams. *8th ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 91–101). ACM Press.
- Krebs, V. (2008). *Orgnet.com: Software for social network analysis and organizational network analysis*. Accessed July 17, 2008, from <http://www.orgnet.com/inflow3.html>.
- LaRowe, G., Ambre, S. A., Burgoon, J. W., Ke, W., & Börner, K. (2009). The Scholarly Database and its utility for scientometrics research. *Scientometrics*, 79(2), 219–234. Accessed September 23, 2008, from <http://ivl.slis.indiana.edu/km/pub/2009-larowe-sdb.pdf>.
- Leydesdorff, L. (2008). *Software and data of Loet Leydesdorff*. Accessed July 15, 2008, from <http://users.fmg.uva.nl/leydesdorff/software.htm>.
- Meho, L. I., & Yang, K. (2007). Impact of data sources on citation counts and rankings of LIS faculty: Web of science versus scopus and google scholar. *Journal of the American Society for Information Science and Technology*, 58(13), 2105–2125. Accessed April 1, 2008, from <http://dlist.sir.arizona.edu/1733/01/meho-yang-03.pdf>.

- O'Madadhain, J., Fisher, D., & Nelson, T. (2008). *Jung: Java universal network/graph framework*. University of California, Irvine. Accessed from <http://jung.sourceforge.net/>.
- OSGi Alliance (2008). *OSGi Alliance*. Accessed July 15, 2008, from <http://www.osgi.org/Main/HomePage>.
- Pauly, D., & Stergiou, K. I. (2005). Equivalence of results from two citation analyses: Thomson ISI's citation index and Google scholar's service. *Ethics in Science and Environmental Politics*, 2005, 33–35.
- Persson, O. (2008). *Bibexcel*. Umeå, Sweden: Umeå University. Accessed July 15, 2008, from <http://www.umu.se/inforsk/Bibexcel/>.
- PostgreSQL Global Development Group (2009). *PostgreSQL: The world's most advanced open source database*. Accessed January 13, 2008, from <http://www.postgresql.org/about/>.
- Python Software Foundation. (2008). *Python programming language—Official website*. Accessed January 13, 2008, from <http://www.python.org/>.
- Siek, J., Lee, L.-Q., & Lumsdaine, A. (2002). *The boost graph library: User guide and reference manual*. New York: Addison-Wesley.
- Takatsuka, M., & Gahegan, M. (2002). GeoVISTA studio: A codeless visual programming environment for geoscientific data analysis and visualization. *The Journal of Computers & Geosciences*, 28(10), 1131–1144.
- Thomson Reuters. (2009). *Web of Science*. Accessed September 23, 2009, from <http://scientific.thomsonreuters.com/products/wos/>.
- The Apache Software Foundation. (2007). *Apache Solr*. Accessed January 13, 2008, from <http://lucene.apache.org/solr/>.
- van Eck, N. J., & Waltman, L. (2009). *VOSviewer*. Accessed December 7, 2008, from <http://www.vosviewer.com/>.
- Williams, T., & Kelley, C. (2008). *gnuplot homepage*. Accessed July 17, 2008, from <http://www.gnuplot.info/>.