



Indiana University
Network Science Institute

Web of Science as a Research Data Set

Valentin Pentchev, Director of IT
Matthew Hutchinson, Data Manager

Indiana University Network Science Institute

- IUNI is cross-campus entity that fosters interdisciplinary, collaborative research in network science.
- Over 150 IU faculty are affiliated with IUNI representing library science, sociology, informatics, neuroscience, public health, and more.
- IUNI Affiliates work with IT and research staff and have access to IUNI IT infrastructure and data.



fostering
**INTERDISCIPLINARY
COLLABORATIVE
RESEARCH**



Web of Science Enclave

A Secure Research Environment

The Data Set

- Over 56,000,000 unique records
- More than 1,000,000,000 references
- Includes records from 1898 -2013 with 2014 and 2015 coming soon.



- Science Citation Index Expanded from 1900-2013
- Social Sciences Citation Index from 1900-2013
- Arts & Humanities Citation Index from 1975-2013
- Book Citation Index -- Science from 2005-2013
- Book Citation Index -- Social Sciences & Humanities from 2005-2013
- Conference Proceedings Citation Index -- Science & Technical from 1990-2013



Balancing Access and Security

- How to allow researchers access to the data without risking the data being released **onto the web**?
- The University classified the data as 'Restricted' adding additional security requirements
- Initial plan had been to store the data on a single machine in a locked office



Reproduced with permission from *Designing Data-Intensive Applications* by Martin Kleppman

The WoS Enclave

1. **WoS Enclave is managed by a dedicated Data Manager, with support from a team of data specialists**
2. **A hardened software environment that does not allow the export of any data by individual users. All exports of data must be performed by the Data Manager**

Hardware:

- A dedicated node on KARST - IU's newest high-throughput computing cluster
- **CPU:** 2x Intel® Xeon® E5-2650 v4 24 Cores 30MB L3 Cache
Memory: 512GB
HDD: 30TB + 1.2 TD SSD

Software:

- **OS:** Redhat Enterprise Linux 6.8
- **Software:** All tools that are available on the University's Karst cluster are available on the enclave and includes standard statistical and development tools including R, SAS, SPSS, MatLab, Python, Stata, LibreOffice and many more



THE DATA - Raw XML

```
<?xml version="1.0" encoding="UTF-8"?>
<!-- Created from PDF via Acrobat SaveAsXML -->
<!-- Mapping Table version: 28-February-2003 -->
<TaggedPDF-doc>
  <?xpacket begin=" id='W5M0MpCehiHzreSzNTczkc9
  <?xpacket begin="" id="W5M0MpCehiHzreSzNTczkc
  - <x:xmpmeta x:xmpmk="Adobe XMP Core 5.2-c001 6:
    - <rdf:RDF xmlns:rdf="http://www.w3.org/1999/
      - <rdf:Description xmlns:xmp="http://ns.adobe
        <xmp:CreateDate>2011-04-06T17:22:05
        <xmp:CreatorTool>ESRI ArcSOC 9.2.0.13
        <xmp:ModifyDate>2011-04-07T08:17:15
        <xmp:MetadataDate>2011-04-07T08:17:
        </rdf:Description>
      - <rdf:Description rdf:about="" xmlns:xmpMM="f
        <xmpMM:DocumentID>uuid:323e04f1-10
        <xmpMM:InstanceID>uuid:6329c275-adc
        </rdf:Description>
      - <rdf:Description rdf:about="" xmlns:dc="http:/
        <dc:format>xml</dc:format>
        </rdf:Description>
      </rdf:RDF>
    </x:xmpmeta>
    <?xpacket end="w"?>
    <?xpacket end='r'?>
  - <Figure>
    <ImageData src="images/WA_Dayton_2011040
    </Figure>
</TaggedPDF-doc>
```

1. Raw XML data from Thomson Reuters/Clarivate Analytics was received and placed in the enclave
2. Unprocessed data available for researchers
3. Parsing software developed by Robert Light and Daniel Halsey at IU's Cyberinfrastructure for Network Science (CNS)
 - Developed in Python and available for use under the Apache License 2.0 from:
[https://github.com/cns-
iu/generic_parser](https://github.com/cns-iu/generic_parser)



THE DATA - Relational Databases

1. PostgreSQL database with a fully parsed XML data

- **wos_core**: 46 tables, joined on 'wos_id'
- **wos_simple**: 13 tables with serial
Public and Foreign key relations

2. Database schema and Data Dictionary documentation

- **wos_core**: fully documented and available to users
- **wos_simple**: under development, available for upon request



PostgreSQL



THE DATA – Browser-based Interface

IUNI Web of Science

Browser-Based Query Interface (beta)

This interface allows users who are unfamiliar with SQL to access Web of Science data by using a simple browser-based form to query the database. During Step 1, users can begin their search by entering authors, journals or keywords. Step 2 allows users to filter down their result set before starting the export process. Step 3 lets users customize the output of the file export to suit their specific needs.

Step 3: Export Data

[← Back to Step 2](#)

Export Fields

Select All

WoS ID Title Year Author Full Names

Author Emails [More Fields](#)

Preview

WoS ID wos_id	Title title	Year year	Author Full Names authors_full_name
WOS:000326968800005	Collectively Against the State	2013	Boerner, Stefanie Borner, S
WOS:000304017900003	An Introduction to Modeling Science: Basic Model Types, Key Definitions, and a General	2012	Scharnhorst, A Boerner, Katy vandenBesselaar, P Boerner, Katy Boyack, Kevin W. Scharnhorst, A Milojevic, Stasa Morris, Steven Boyack, Kevin W. Milojevic, Stasa Boerner, Katy vandenBesselaar, P Morrison, Steven Milojevic, Stasa Borner, K Borner, K Morrison, Steven Borner, K vandenBesselaar, P Morrison,

Download Format

CSV JSON

Field Separator

Secondary Separator

[Begin Download](#)

1. Created to allow users unfamiliar with SQL to access data
2. Queries the wos_simple database
3. Allows for text searching on Author Names, Journal Names, Article Names, Abstracts and Years
4. Python backend based on the Flask micro server package



THE DATA – Citation Interface

maahutch@c265.karst.uit.siu.edu - ThinLinc Client

Reference Citation

Enter WoS ID in the box below

WOS:000289266400009

Search

Table Histogram

Show 25 entries

key	id	ref_ctr	ref_id	cited_author	assignee	year	page	volume	cited_title	cited_work	doi
811655143	WOS:000289266400009	30	WOS:000290753500015.38	KIRLIK A		1998	91		MAKING DECISIONS STR		
811655124	WOS:000289266400009	49		WEGNER DM		1987	185		THEORIES GROUP BEHAV		
811655125	WOS:000289266400009	48	WOS:000289266400009.48	VONDEROELSNI T D		2006	54	39	PERSONALFUHRUNG		
811655126	WOS:000289266400009	47	WOS:000244671300008	Vashd, DR		2007	115	46	Boiling-debriefing: Using a reflexive organizational learning model from the military to enhance the performance of surgical teams	HUMAN RESOURCE MANAGEMENT	10.1002/hrm.20148
811655127	WOS:000289266400009	46	WOS:A1965CCL3000002	TUCKMAN, BW		1965	384	63	DEVELOPMENTAL SEQUENCE IN SMALL-GROUPS	PSYCHOLOGICAL BULLETIN	
811655128	WOS:000289266400009	45	000328072500004.138	Tindle, RS		2000	123	3	'Social Sharedness' as a unifying theme for information processing in groups	Group Processes Integr Relat	
811655129	WOS:000289266400009	44	WOS:000267738900006	Smith-Jentsch, KA		2009	181	51	Do Familiar Teammates Request and Accept More Backup? Transactive Memory in Air Traffic Control	HUMAN FACTORS	10.1177/0018720809333367
811655130	WOS:000289266400009	43	WOS:000256262300003	Smith-Jentsch, KA		2008	303	39	Guided team self-correction - Impacts on team mental models, processes, and effectiveness	SMALL GROUP RESEARCH	10.1177/1046496408317794
811655131	WOS:000289266400009	42	CCC:000169689000003	Smith-Jentsch, KA		2001	31		Uncovering differences in team competency requirements: The case of air traffic control teams	IMPROVING TEAMWORK IN ORGANIZATIONS	
811655132	WOS:000289266400009	41	000320893800004.109	Salas, E		1997	249		Methods, tools, and Training for a Rapidly		

1. Created to allow users unfamiliar with SQL to explore the reference and citations networks for individual records
2. Also queries the wos_simple database
3. Allows searching both citations and references
4. Created using R Shiny



Super Computing and the Web of Science

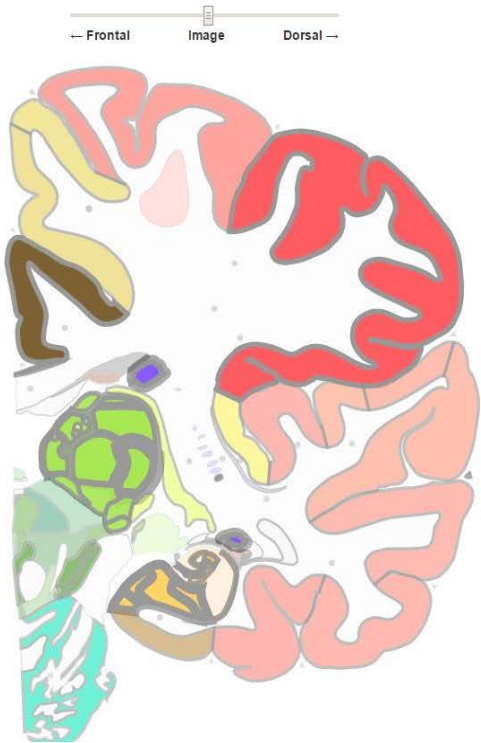
1. Scripts can be run against the raw WoS XML data on both Karst and Big Red II
2. Simply submit your script along with any necessary instructions to the IUNI Data Manager



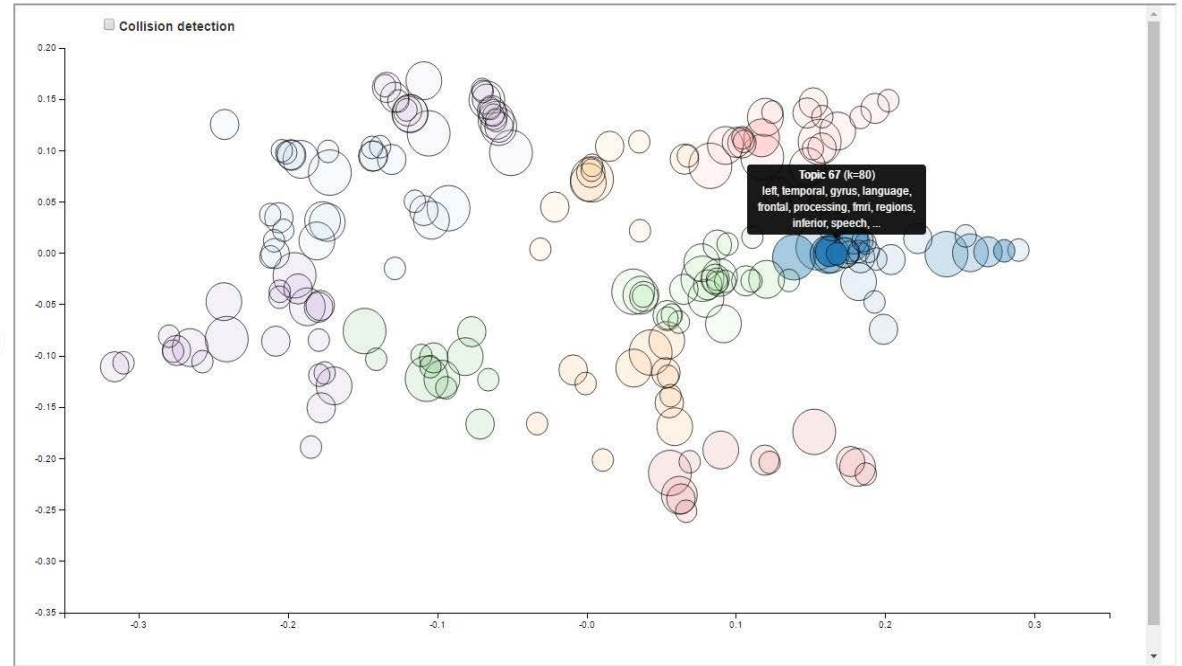
Research:

1. **Not just Bibliometrics and Scientometrics!**
2. **Cognitive Science – Hyperbrain project**
3. **Network Science – Correlating Transportation Networks with Academic Collaborations**
4. **And many more...**



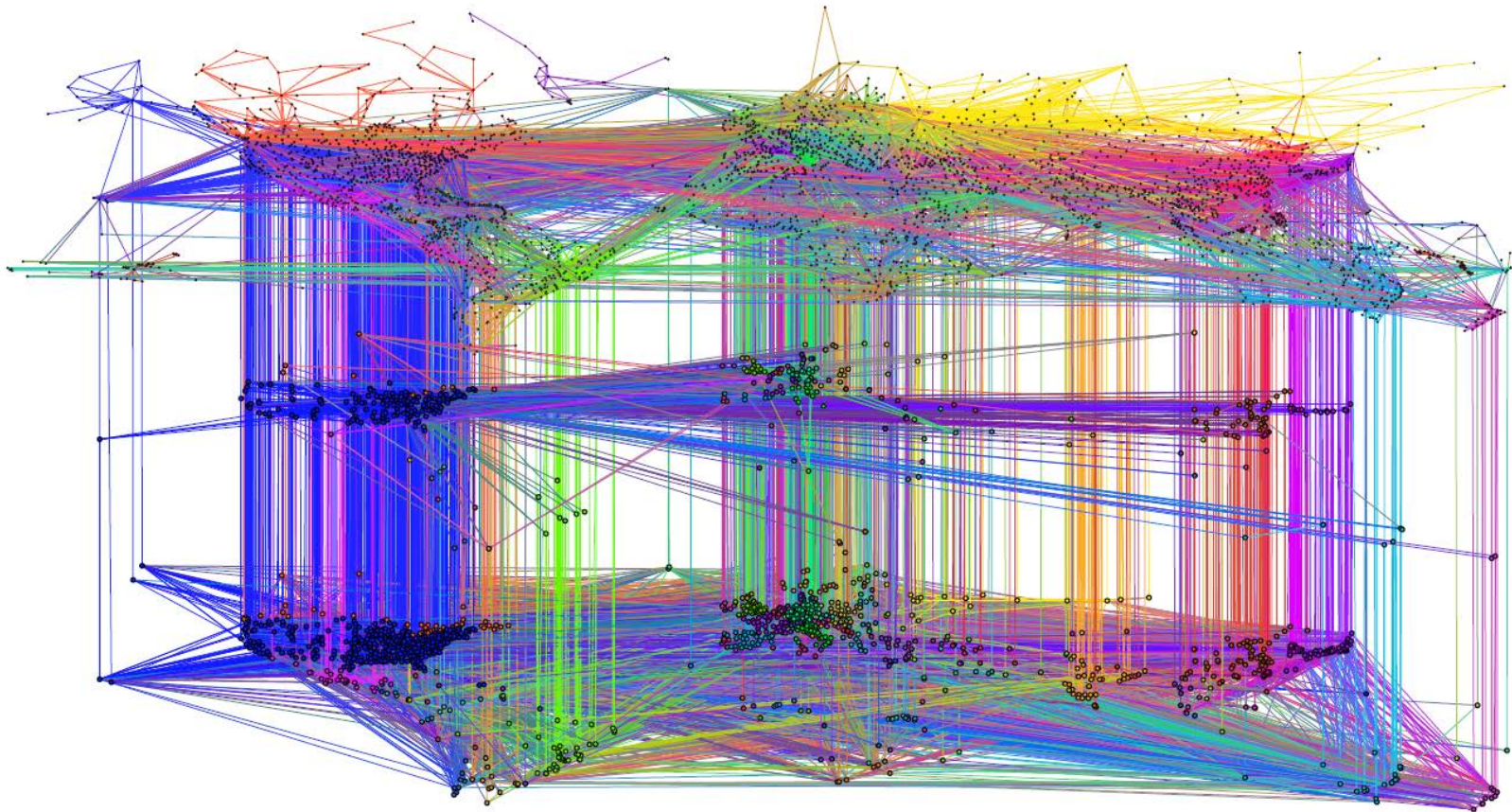


supramarginal gyrus, left



Hyperbrain.org, reproduced with permission by Jaimie Murdock & Franco Pestilli



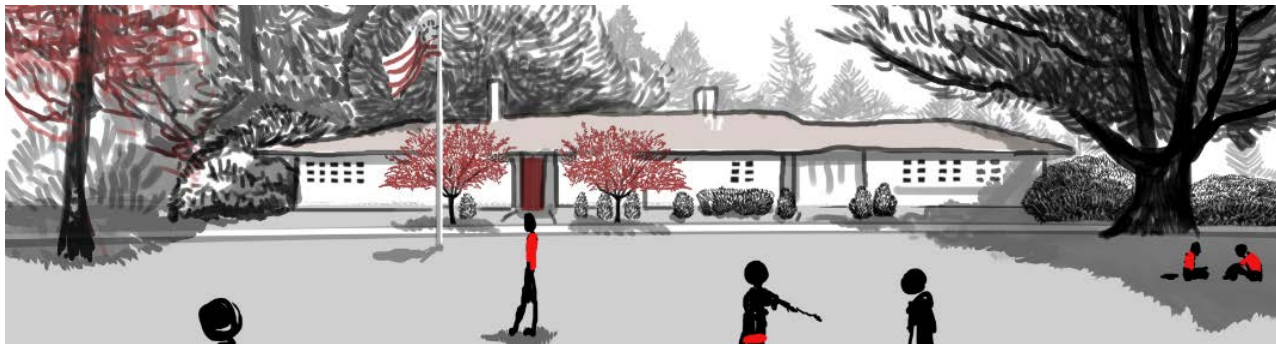


“Mapping Six Years of Research by Three Universities” Reproduced with permission by Xiaoran Yan, Indiana University



Interested in Working with IUNI's copy of the WoS data?

1. Simply go to <http://iuni.iu.edu/resources/web-of-science>
2. Read the details of the data available and the WoS Access Policy
3. Complete the application for 'General Data User'
4. Request an account on Karst if you do not already have one
5. Wait to be contacted by IUNI staff (which should occur within 24 hrs)



The Team



Valentin Pentchev
Director of Information Technology
IUNI



Katy Börner
Faculty, Center Director
CNS



Matt Hutchinson
Data Manager
IUNI



Robert Light
Senior Systems Analyst, Database
Administrator
CNS



Ben Serrette
Web Developer
IUNI



Daniel Halsey
Senior System Architect, Project Manager
CNS



Patricia Mabry
Executive Director
IUNI



Joe Rinkovsky
Lead Systems Administrator
High Performance Systems, UITS

Acknowledgements: Jennifer Adams, IU General Counsel, Sara Chambers and the IU Committee of Data Stewards and Timothy Otto, Clarivate Analytics





Indiana University
Network Science Institute

Thank You

Valentin Pentchev
Director of IT

1001 E SR 45/46 Bypass | Bloomington,
IN 47408-1415
Email: vpentche@iu.edu |
Phone: (812) 856-7087 | Fax: (812) 856-1192

Matthew Hutchinson
Data Manager

1001 E SR 45/46 Bypass | Bloomington,
IN 47408-1415
Email: maahutch@iu.edu |
Phone: (812) 855-1404 | Fax: (812) 856-1192



INDIANA UNIVERSITY BLOOMINGTON
FULFILLING *the* PROMISE