

**BIG DATA**

**VAST DATA**

Unstructured data

Semi-Structured data

Volume

Useful Metadata

Data analysis

Petabytes

Unstructured content

People driven

Decision making

Zettabyte

Framework

Smart content database

Text analytics

Semantic Metadata

Concept extraction

Structured data

# Big Data Content Organization, Discovery, and Management

Marjorie M.K. Hlava  
President  
Access Innovations, Inc.  
Mhlava@accessinn.com

12/11/2013

Access Innovations, Inc..



# Outline

- Big Data
- New Government Initiative
- Content Organization
- Discovery (Search)
- Management
- Skills we bring
- Examples of what we can do



# Why do we care about Big Data?

- Data is the new oil – we have to learn how to mine it! Qatar – European Commission Report
- \$ 7 trillion economic value in 7 US sectors alone
- \$90 B annually in sensitive devices
- An insurance firm with 5 terabytes of data on share drives pays \$1.5 m per year
- New McKinsey 4<sup>th</sup> factor of production
- Land, Labor, Capital, + Data

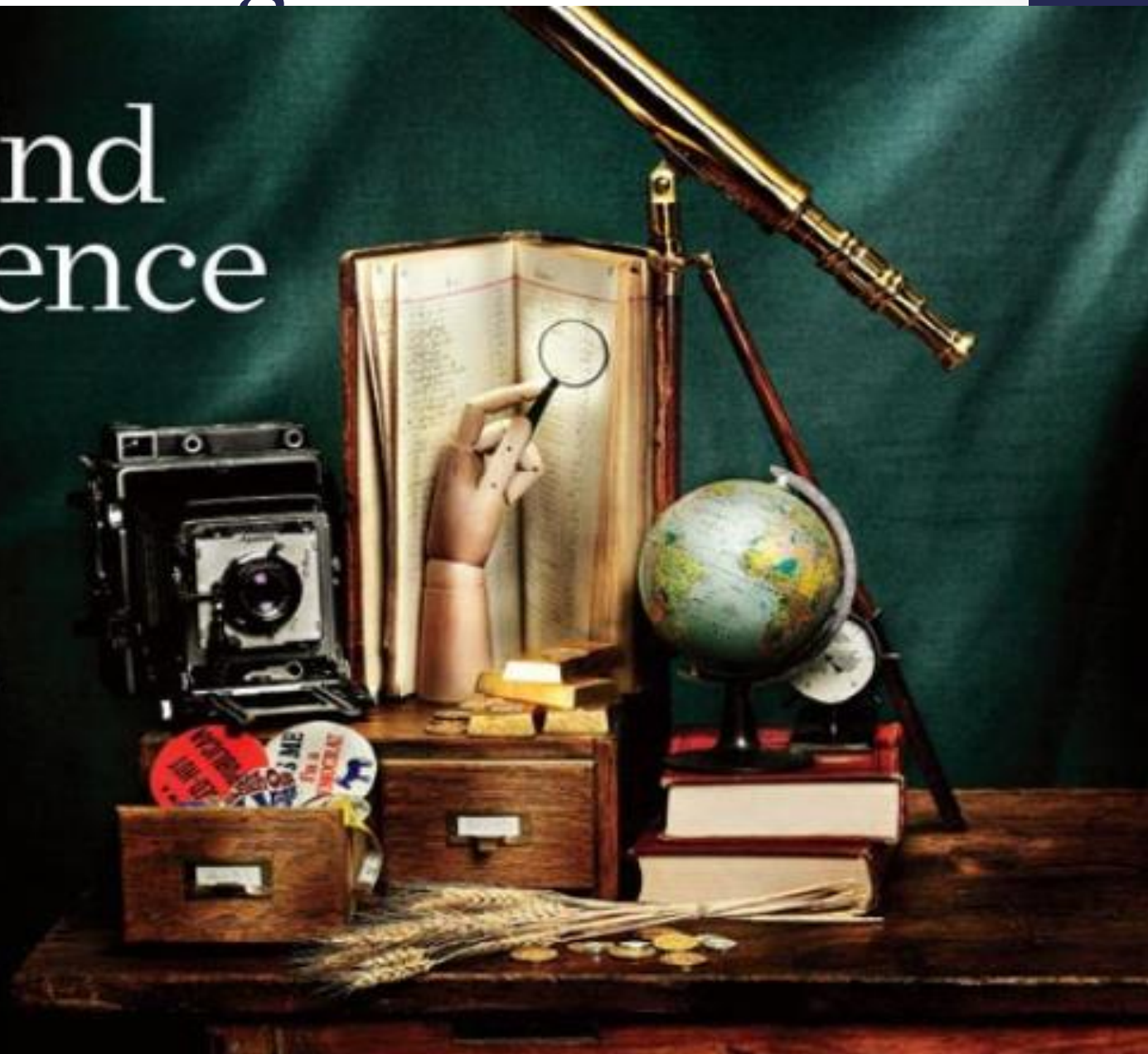




# The Data Deluge – Wired 16.07

## The End of Science

The quest for knowledge used to begin with grand theories. Now it begins with massive amounts of data. Welcome to the Petabyte Age.



# Big Data Born

- Google, eBay, LinkedIn, and Facebook were built around Big Data from the beginning.
- No need to reconcile or integrate Big Data with more traditional sources of data and the analytics performed upon them
- No merging Big Data technologies with their traditional IT infrastructures
- Big Data could stand alone, Big Data analytics could be the only focus of analytics
- Big Data technology architectures could be the only architecture.



# Integrating Big Data

- Large, well-established orgs
- Must be integrated with everything else that's going on in the company.
- Analytics on Big Data have to coexist with analytics on other types of data.
- Hadoop clusters have to do their work alongside IBM mainframes.
- Data scientists must somehow get along and work jointly with mere quantitative analysts.



# What is Big Data?

- *Big Data is a term applied to data sets whose size is beyond the ability of commonly used software tools to capture, manage, and process the data within a tolerable elapsed time. Big Data sizes are a constantly moving target currently ranging from a few dozen terabytes to many petabytes of data in a single data set. – Wikipedia, May 2011*
- “Unstructured”
- Terabytes, petabytes, zettabytes
- Streaming



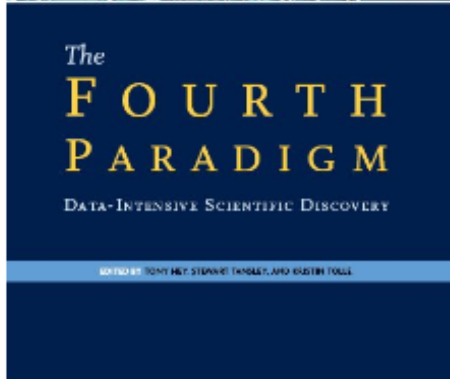


# New kind of science?

2013

## The Fourth Paradigm: Data-Intensive Scientific Discovery

Presenting the first broad look at the rapidly emerging field of data-intensive science



Increasingly, scientific breakthroughs will be powered by advanced computing capabilities that help researchers manipulate and explore massive datasets.

The speed at which any given scientific discipline advances will depend on how well its researchers collaborate with one another, and with technologists, in areas of eScience such as databases, workflow management, visualization, and cloud computing technologies.

In *The Fourth Paradigm: Data-Intensive Scientific Discovery*, the collection of essays expands on the vision of pioneering computer scientist Jim Gray for a new, fourth paradigm of discovery based on data-intensive science and offers insights into how it can be fully realized.

### Download

- Full text, low resolution (6 MB)
- Full text, high resolution (93 MB)
- By chapter and essay

### Purchase from Amazon.com

- Paperback
- Kindle version

### In the news

- Sailing on an Ocean of 0s and 1s (National Geographic Magazine)
- A Deluge of Data Shapes a New Paradigm (New York Times)

# New Special Collections

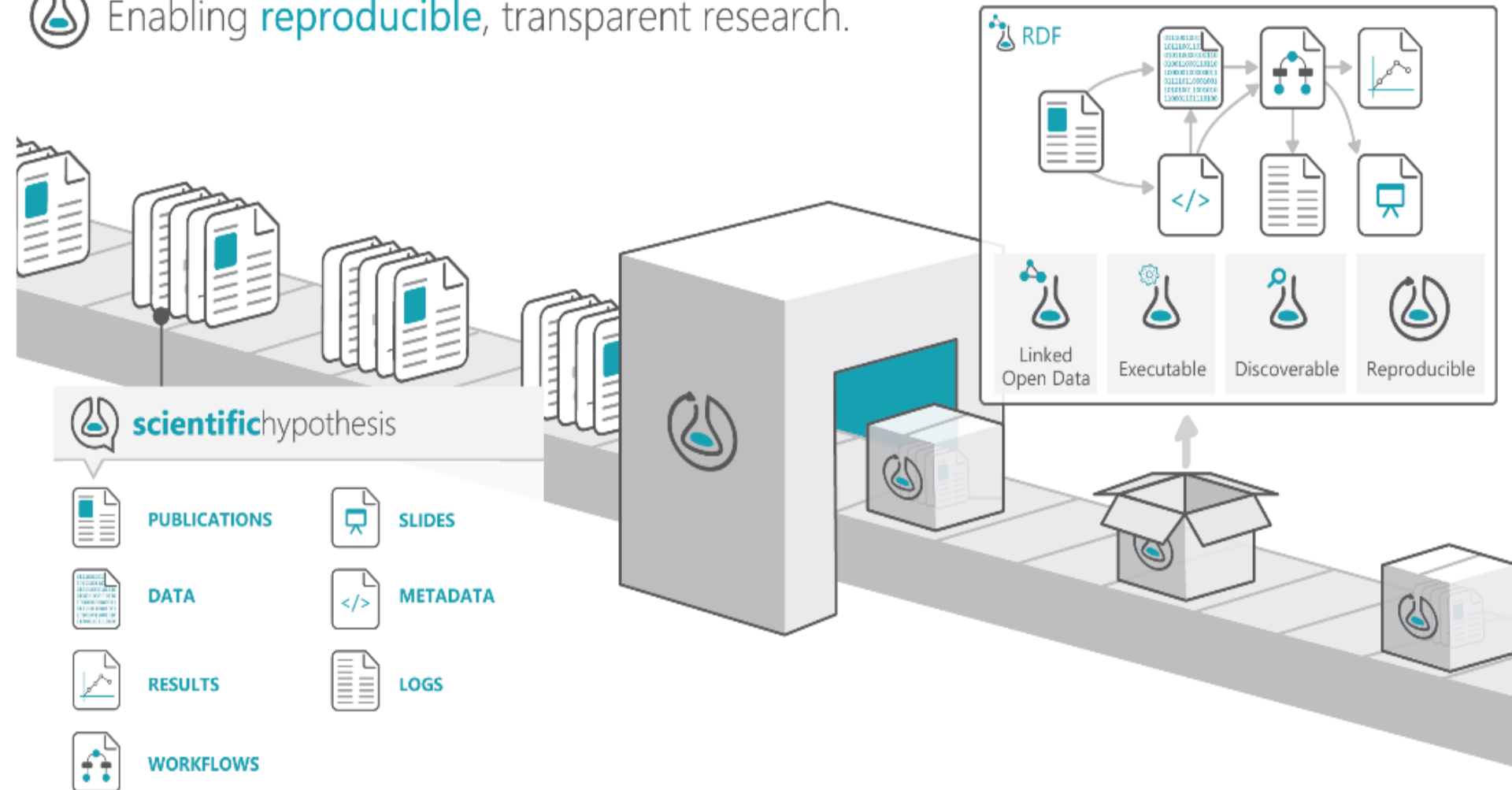
- Volume, Velocity, Variety
- Ability to deal overwhelmed
- More about methods than data
- Location aware data
- Life streaming
- Insurance claims
- Hubble telescope
- CERN Collections
- Flight data



# www.researchobject.org

3

 Enabling **reproducible**, transparent research.



# Unstructured data

- Means untagged or unformatted
- PDF
- Word files
- File shares
- News feeds
- News Data feeds
- Images





# Bit of a misnomer

- All data has some structure and more structure possibilities
- PDF properties
- Word file properties
- File structures

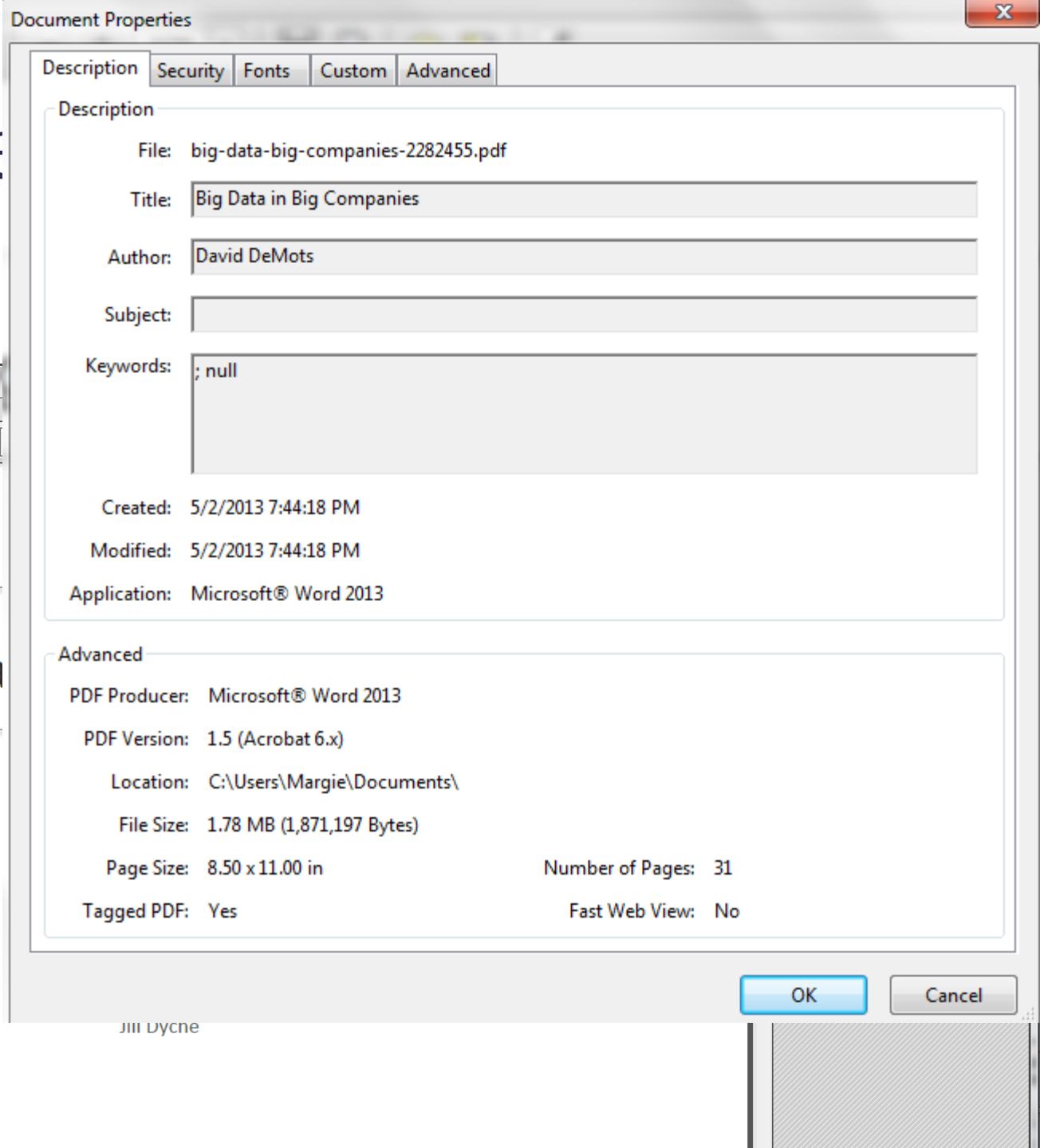
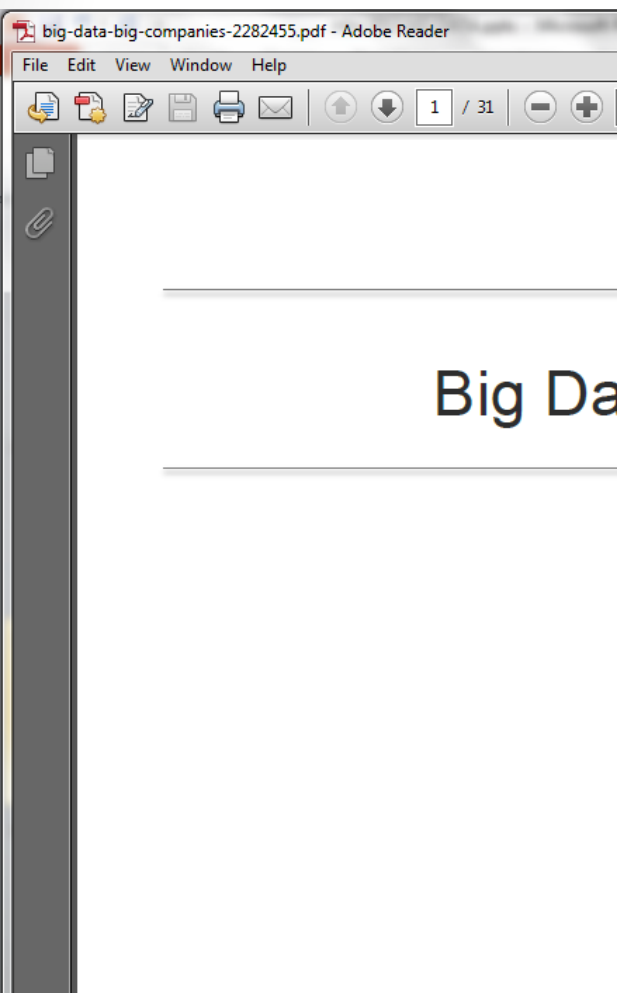


# Property tables

- PDF - property tables
- Word files – property tables
- Files shares – implied structure in file names
- News feeds – headers have metadata
- Hubbell telescope feeds
- Images



# PDF Properties



# File structures

Eudora - [Junk]

File Edit Mailbox

Search Eudora

62  
54  
52  
53  
67  
25/507K/OK

Subject: Taxonomy  
From: TaxoDiary <s...>  
To: mhlava@acces...

Taxonomy

- Seman
- Better
- Finding

Semantic Te

Rust@ipiran.ru, 01:...

In Ctrl+1  
Out Ctrl+0  
Junk  
Trash

Recent

New...  
Drafts  
In - OLD  
Jokes  
keyers and subcontractors-old-t  
Meeting and Presentations 2014  
Network Management  
OUT\_OLD  
out-old-s  
out-old-t  
Politics  
Administration  
Business People  
Data Harmony  
dataharmony-old-t  
Family and Personal  
family stuff-old-t  
Frequent Flyer and Hotel points  
HAH  
HFHS  
Hubbell House Alliance  
IT  
Leads  
Leads - OLD  
List serves  
marketing-old-t  
Presentations and meetings 2006  
Presentations and Meetings 2007  
Presentations and Meetings 2008  
Presentations and meetings 2009  
Presentations and Meetings 2010

Presentations and Meetings 2013  
Presentations and Meetings 2014  
presentations-old-t  
prof acts  
Projects  
projects-old-t  
Publications  
Staff Reports  
Standards

DOI  
Eugene Garfield  
FEA Reference Manual  
GDS - IWA  
Homer Hall  
ICSTI  
IMLS Linked data  
Information Systems and Use

New...  
2010 programs  
2012 programs  
Board meetings  
October 2009 Program  
Programs  
RGC/Strategic Planning  
Rio Grande archive / oral histo  
Rio Grande Newsletter  
Strategic Directions  
Website

Places - spaces Advisory Board  
SIA  
SSP  
St Regis  
STM  
SUMmit - ROger  
TRSS Paper  
UNC - IMLS Advisory Board  
ALA  
ASIST  
NFAIS  
Presentations  
SLA

New...  
2003 Annual Meeting  
2005 Annual meeting  
2006 Meeting  
2007 Annual meeting  
2010 - Winter meeting St Louis  
2010 Annual Meeting New Orleans  
2011 Winter meeting  
Archivist  
Award Nominations  
Board chat  
Brainstorming a new SLA  
Bylaws  
Cataloging  
Committee on Committees  
Competitive Scan  
Confluence  
Connections  
Division task force  
Elect Doris Campagin  
Encore  
Executive Director search  
Fellow Nomination  
Fiesta  
Finance  
Government Relations  
Helpful information  
Information Outlook  
ITE Division  
Kentucky  
Leadership list  
Membership  
Nominations Committee  
Oral History Project  
Passwprd  
Photos  
Public Relations Committee  
Research

Web Site and DH  
2004 Nashville  
2008 Seattle  
2009 Washington  
2010  
2011  
2012  
2013 Annual meeting  
2014  
Archives  
Fellows  
Rio Grande Chapter



# Structu

- XML tagg

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE us-patent-grant SYSTEM "us-patent-grant-v42-2006-08-23.dtd">
<us-patent-grant lang="EN" country="US" dtd-version="v4.2 2006-08-23">
  - <us-bibliographic-data-grant>
    - <publication-reference>
      - <document-id>
        <country>US</country>
        <doc-number>6170828</doc-number>
        <date>20010109</date>
      </document-id>
    </publication-reference>
    - <application-reference appl-type="utility">
      - <document-id>
        <country>US</country>
        <doc-number>09492395</doc-number>
        <date>20000127</date>
      </document-id>
    </application-reference>
    <us-application-series-code>09</us-application-series-code>
  - <classifications-ipcr>
    - <classification-ipcr>
      - <ipc-version-indicator>
        <date>20000101</date>
      </ipc-version-indicator>
      <section>A</section>
      <class>63</class>
      <subclass>F</subclass>
      <main-group>1</main-group>
      <subgroup>00</subgroup>
    </classification-ipcr>
  </classifications-ipcr>
  - <classification-national>
    <country>US</country>
    <main-classification>273292</main-classification>
    <further-classification>273292</further-classification>
```

# What are the problems?

- Data infrastructure challenges
- “taking diverse and heterogeneous data sets and making them more homogeneous and usable”
- An opportunity?
- All that data – what can it tell us?
- Privacy
- Copyright
- Neurological impact
- Data collection methods



# New Government Initiative

The Big Data Senior Steering Group (BDSSG) was formed to identify current Big Data research and development activities across the Federal government, offer opportunities for coordination, and identify what the goal of a national initiative in this area would look like. Subsequently, in March 2012, The White House Big Data R&D Initiative was launched and the BDSSG continues to work in four main areas to facilitate and further the goals of the Initiative.



# The National Big Data R&D Initiative

- Fast-growing volume of digital data of digital data
- Advance state-of-the-art core technologies needed to collect, store, preserve, manage, analyze, and share huge quantities of data.
- Harness these technologies to accelerate the pace of discovery in science and engineering, strengthen our national security, and transform teaching and learning; and
- Expand the workforce needed to develop and use Big Data technologies.





# Data to Knowledge to Action

- Advance supporting technologies
  - Big Data
  - Data analytics;
- Educate and expand the Big Data workforce
- Improve key outcomes in economic growth, job creation, education, health, energy, sustainability, public safety, advanced manufacturing, science and engineering, and global development

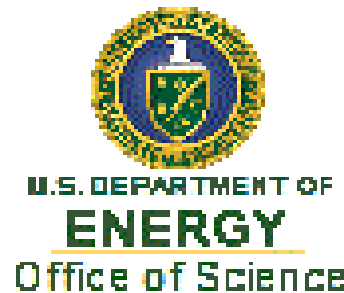


# Data on Data

- January 22, 2013 - Data on Data: Presenting Stakeholder Alignment Data on the Cyberinfrastructure for Earth System Science
- Presentation and discussion with Professor Joel Cutcher-Gershenfeld. Professor Cutcher-Gershenfeld presented information on the NSF EarthCube initiative including stakeholder survey data (approximately 850 responses).



# Who is involved?



# Groups breakdown

- INTERAGENCY WORKING GROUPS
  - Cyber Security and Information Assurance
  - High End Computing
- COMMUNITY OF PRACTICE (CoP)
  - Faster Administration of Science and Technology Education and Research
- COORDINATING GROUPS
  - Human Computer Interaction and Information Management
  - High Confidence Software and Systems
  - Large Scale Networking
  - Software Design and Productivity
  - Social, Economic, and Workforce Implications of IT





# SENIOR STEERING GROUPS (SSGs)

- Big Data
- Cyber Physical Systems
- Cyber Security and Information Assurance Research and Development
- Health Information Technology Research and Development
- Wireless Spectrum Research and Development
- SUBGROUP
- Health Information Technology Innovation and Development Environments Subgroup
- TEAMS
- Joint Engineering Team
- Middleware and Grid Interagency Coordinating Team

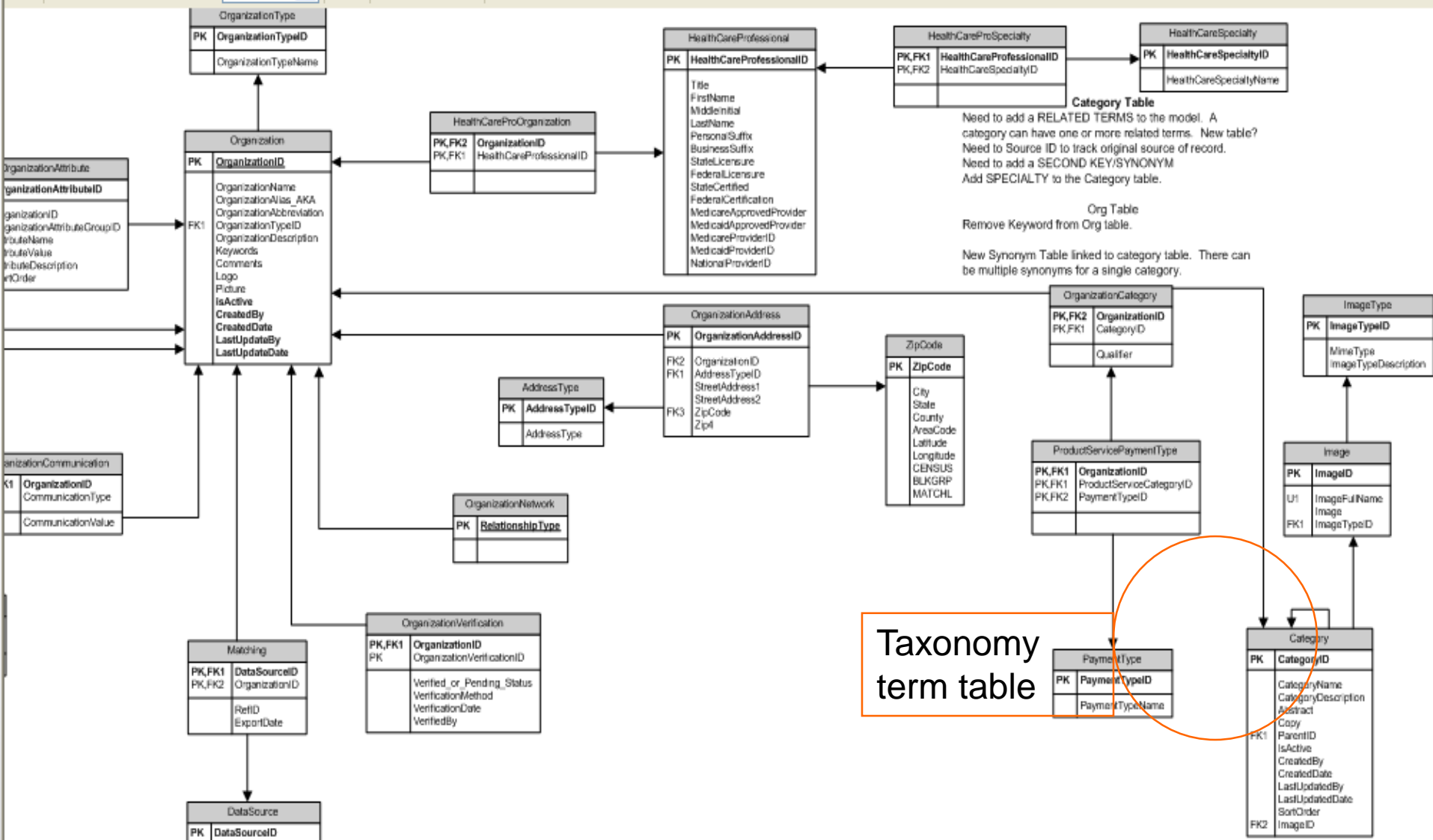


# Content Organization

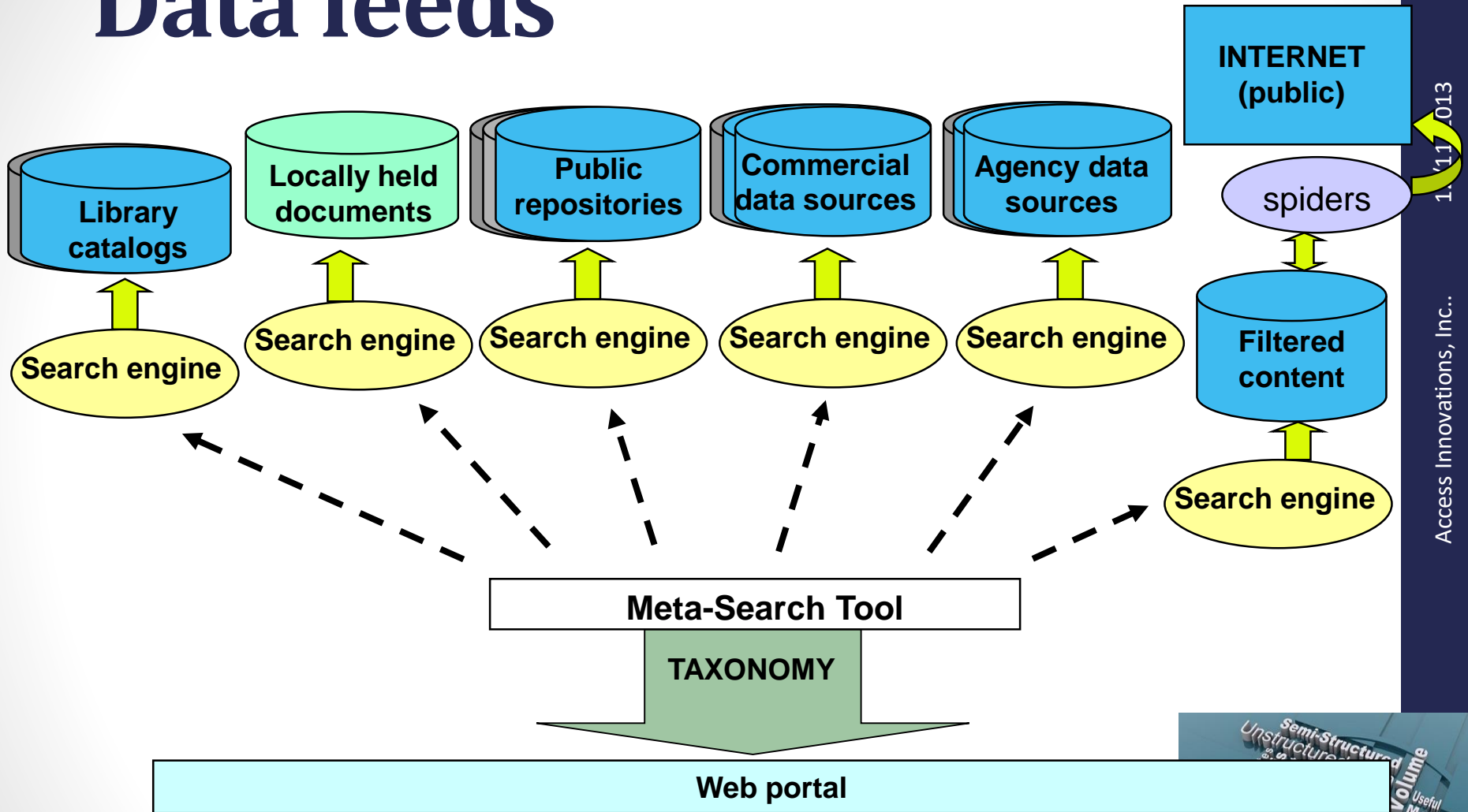
- Data on machines
  - Local
  - Cloud
  - Remote
  - Streaming
- Undifferentiated
- Unstructured
- Needs organization
- Type of database structure
  - RDBMS
  - Object oriented



# RDBMS Connection



# Data feeds



# Metadata options

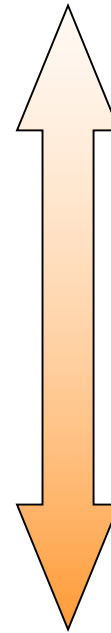
- Structure unstructured data
- Create metadata
- Where to put it?
- Store metadata with the records
  - HTML header
  - Properties tables
  - XML files
- Store metadata in a separate file
  - Database
  - Metadata repository
  - Search system
  - File structure
  - SharePoint application
  - Web interface



# Information retrieval starts with a knowledge organization system

- Uncontrolled list
- Name authority file
- Synonym set / ring
- Controlled vocabulary
- *Taxonomy*
- *Thesaurus*
- Ontology
- Topic Map
- Semantic Network

***Not complex - \$***



***Highly complex - \$\$\$\$***

***LOTS OF OVERLAP!***





# Structure of controlled vocabularies

List of words    Synonyms    Taxonomy    Thesaurus

INCREASING COMPLEXITY / RICHNESS

Ambiguity control

Synonym control

Ambiguity control

Synonym control

Hierarchical rel's

Ambiguity control

Synonym control

Hierarchical rel's

Associative rel's



# Taxonomy / thesaurus

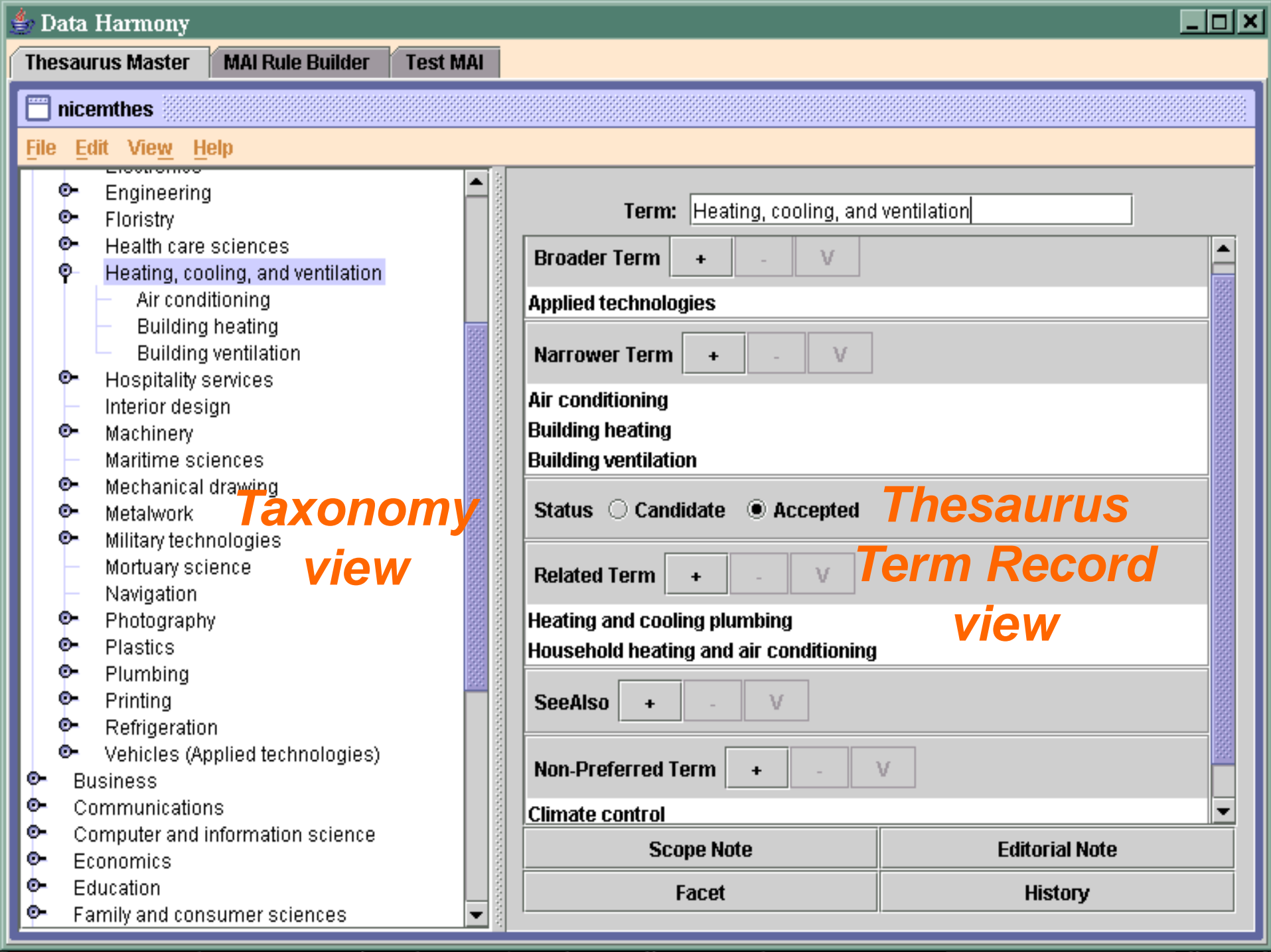
- Main Term (MT) \*
- Top Term (TT)
- Broader Terms (BT)
- Narrower Terms (NT)
- Narrower Term Instance
- Related Terms (RT)
  - See also (SA)
- Non-Preferred Term (NP)
  - Used for (UF), See (S)
- Scope Note (SN)
- History (H)

**THESAURUS**

**TAXONOMY**  
**ONTOLOGY**

\*a.k.a. subject term, heading, node, index term, keyword, category, descriptor, class





# Taxonomies in Context

A taxonomy aspires to be:

- a correlation of the different functional, regional and (possibly) national languages used by a community of practice
- a support mechanism for navigation
- a support tool for search engines and knowledge maps
- an authority for tagging documents and other information objects
- a knowledge base in its own right

*Reference: "Taxonomies: the vital tool of information architecture", [www.tfpl.com](http://www.tfpl.com).*

*Courtesy of Lillian Gassie, Naval Postgraduate School, Monterey, CA*



# Where to use a taxonomy

- Link the taxonomy and indexing
- Always in sync with the industry
- Keep up to date with terminology
- Automatically index the old data
- Filter newsfeeds
- Search using the taxonomy
- File using the taxonomy
- Spell check using the taxonomy
- Link to translation system
- Catalog using the taxonomy



# Value of a Taxonomy

- Improves organization & structure
- Facilitates navigation
- Facilitates knowledge discovery
- Reduces effort
- Saves time

“Taxonomies are better created by professional indexers or librarians than by domain experts.”

Courtesy of Lillian Gassie, Naval Postgraduate School, Monterey, CA



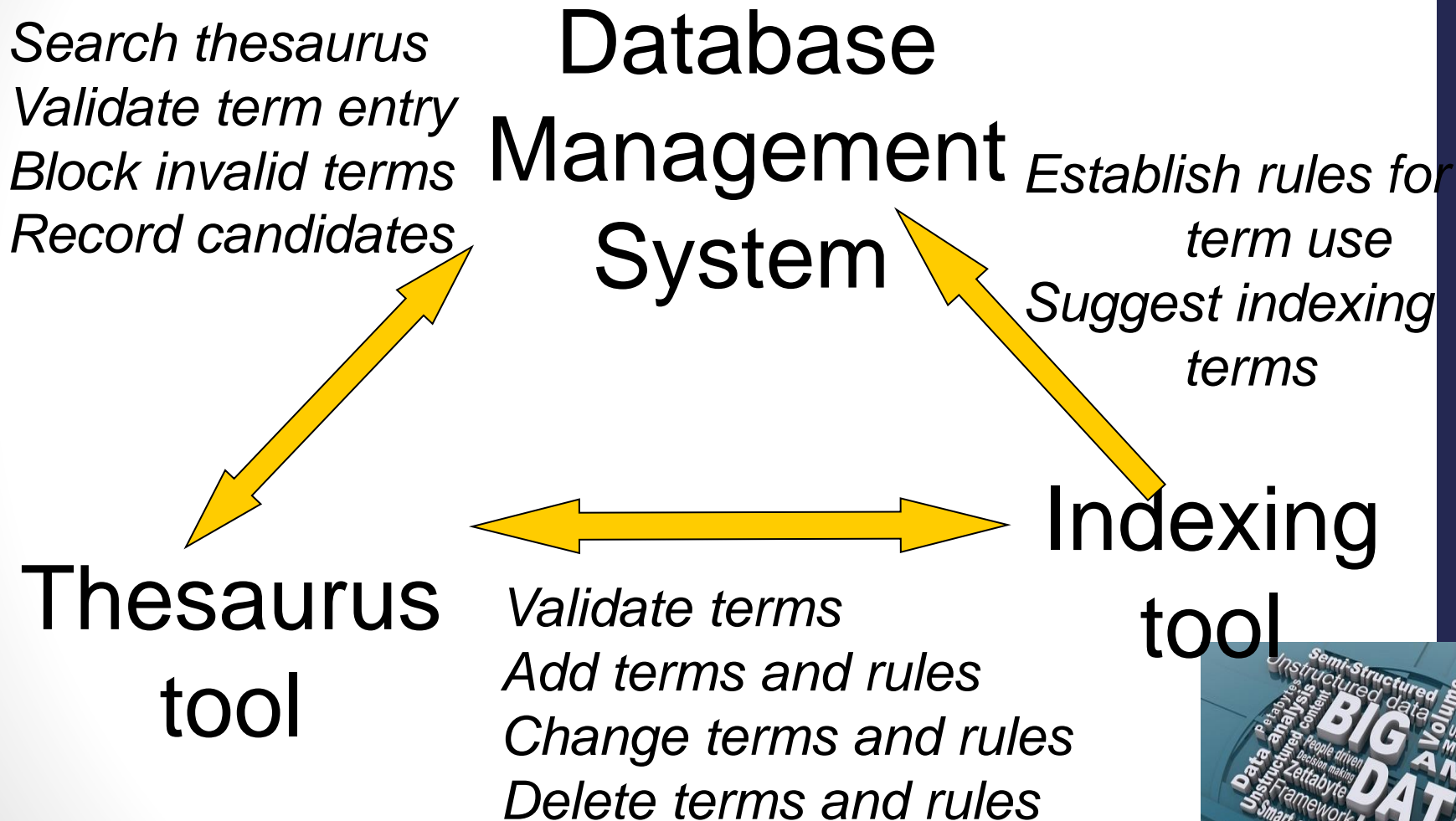


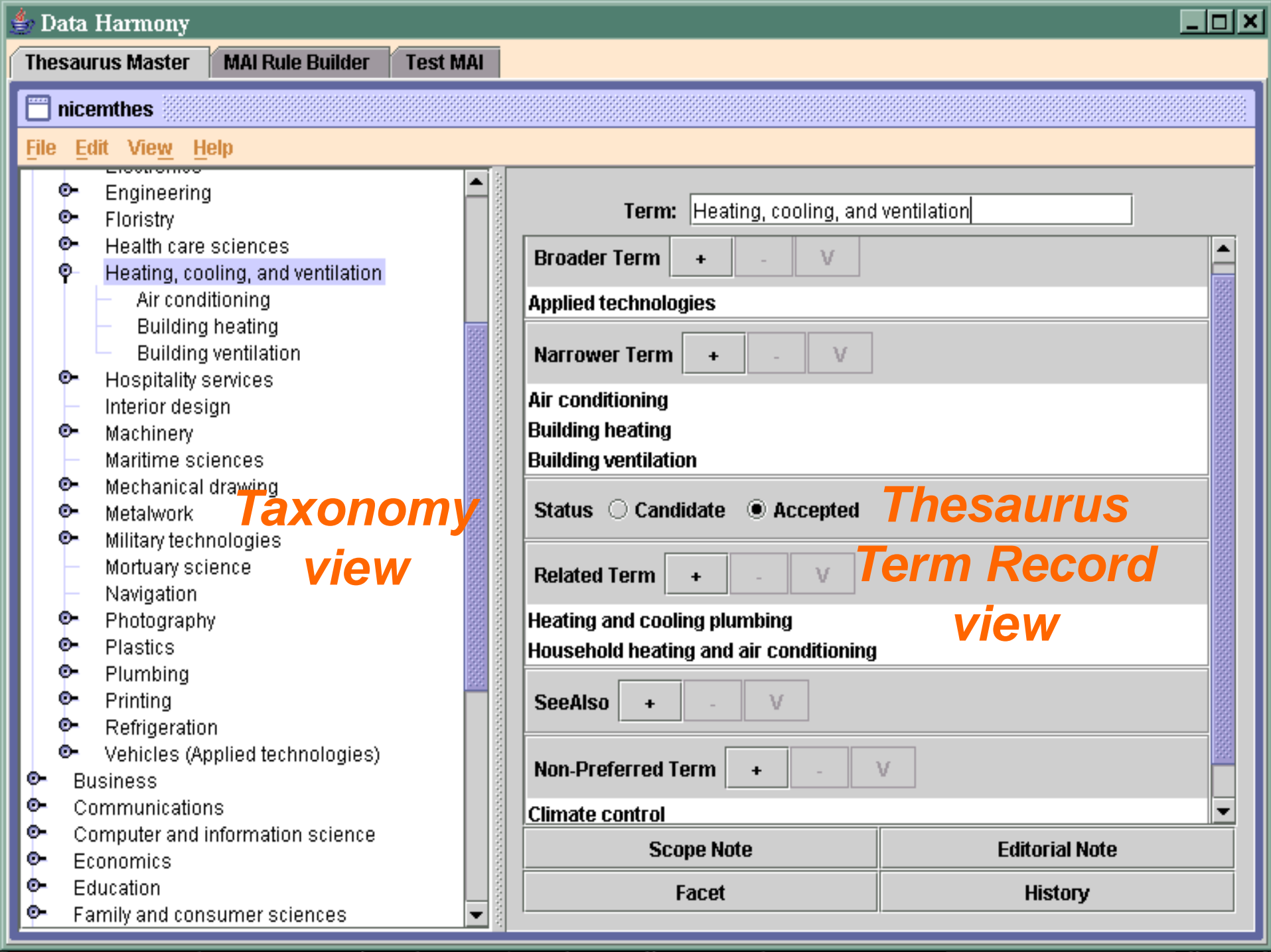
# Taxonomies

- A well formed taxonomy is based on a thesaurus
- Provides a flexible platform for many views of the taxonomy
- Allows fast deployment
- Is the basis for a good
  - knowledge management system
  - search retrieval system
  - portal interface



# A quick look behind the scenes





Edit nicem 0464973 1 of 3

**File**

Original Title				Filing	
Germs and Disease				N	
Original Language	Country of Origin	Production Year	Production Date Misc	Producer Code	
ENGLISH		1999			
General Audience		Specific Audience	Distributor		
			REVID		
P					

**Abstract**

Shows what bacteria and viruses are and how they enter the body. Teaches about common childhood diseases, viruses why vaccinations can play an important role in prevention. Cleanliness is emphasized.

Thesaurus Descriptors	Additional Descriptors	Other Field
		0
MAI		

**MAI Suggested Terms**

- Viruses|(4) viruses(4)
- Disease prevention|(2) prevent\*(2)
- Pharmaceutical drugs|(2) vaccin\*(2)
- Microorganisms|(2) germs(2)
- Personal health and hygiene|(2) germs(2)
- Childhood (Age groups)|(2) childhood(2)
- Bacteria|(2) bacteria(2)
- School learning and achievement|(1) students(1)
- Students|(1) students(1)
- Antibiotics|(1) antibiotics(1)
- Youth (Age groups)|(1) students(1)
- Animated photography|(1) animation(1)

Select

Suggested taxonomy descriptors

# Workflow order

- Create the metadata structure
- Gather the locations of the records
- Index in place?
- Point to the data?
- Add metadata to the record?
- Store metadata in separate “table”?
- Use in full text?
- Use in search?
- Link databases and data sets with APIs



# APIs and Web Calls

- Link data cross platforms
- Not federated search
- Examples
  - Was: Card catalog or OPACs to books on shelf
  - Now: EBSCO host to Mendeley, Zotero, ResearchGate, Academic.edu, etc.
  - Use an API or web call
- Web calls
  - Call to another web platform
- APIs, written handshakes
  - Z39.50 is one standard for libraries





# Discovery (Search)

- Search
  - Free text / full text
  - Fielded / Faceted
- Parts of Search
  - Presentation layer
  - Caching
  - Inverted index



# Kinds of search

- Keyword –
  - Autonomy / Verity
- Bayesian –
  - FAST
  - Lucene
- Faceted
  - Endeca
- Boolean
  - Dialog
- Ranking Algorithms
  - Google



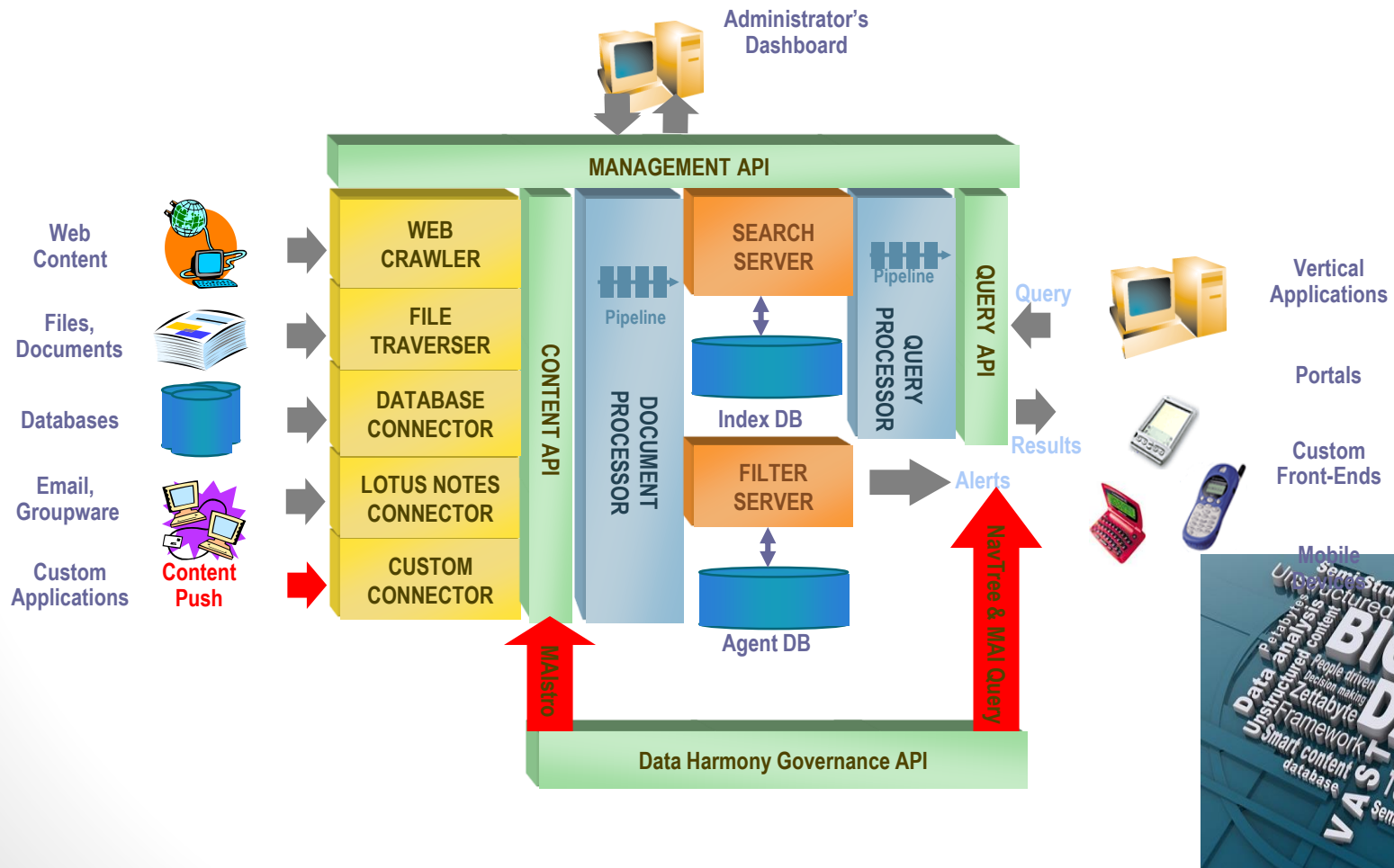
# Parts of search

- Query language
- Search technology
- Inverted index
- Ranking algorithms
- Other enhancements

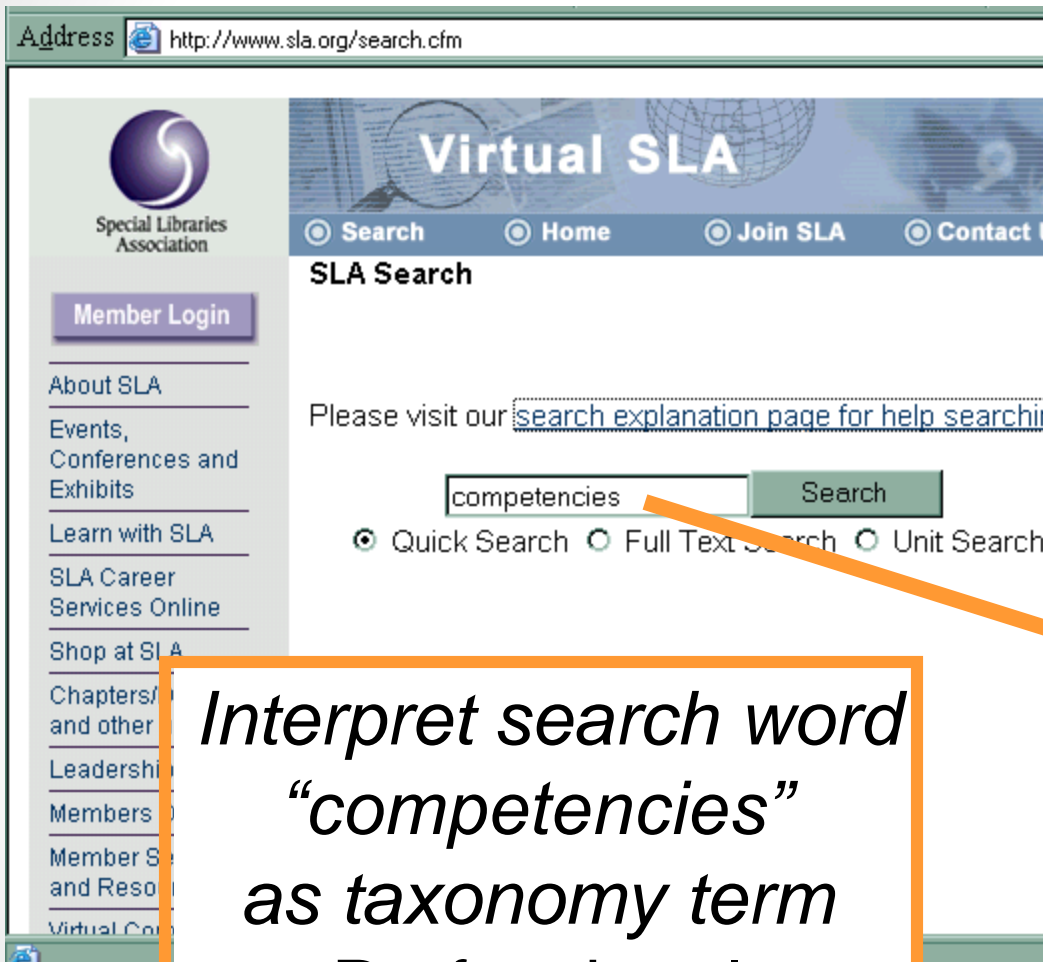


# FAST ESP and Data Harmony Architecture

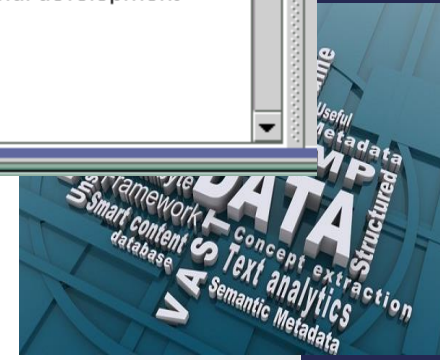
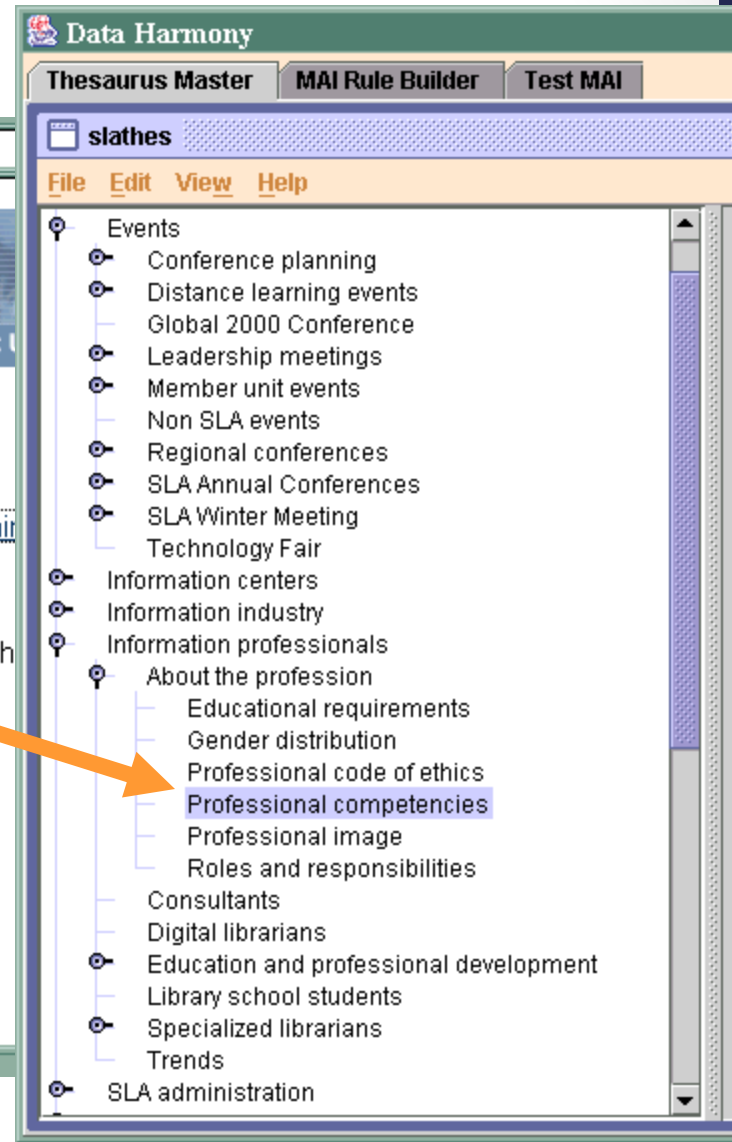
## Core Architectural Components



# SLA search



*Interpret search word  
“competencies”  
as taxonomy term  
Professional  
competencies*





ABOUT US MEMBERSHIP

SEARCH SLA'S WEB SITE

SLA Search

Taxonomy powered by:



DATA HARMONY

Searched using taxonomy  
27 matches found

- [Competencies for Information Professionals](#)  
<http://www.sla.org/content/competencies/index.cfm>
- [Speaker Bios](#) (09/01/2004)  
<http://www.sla.org/content/speakerbios/index.cfm>
- [Resources for Emerging Librarians](#)  
<http://www.sla.org/content/resources/emerginglibrarians/index.cfm>
- [July 2005 - SLA Calendar](#)  
<http://www.sla.org/content/calendar/index.cfm>

Search SLA's Web Site - Special Libraries Association - Microsoft Internet Explorer

File Edit View Favorites Tools Help

- [SLA Press Release 2003-02](#) (02/11/2004)  
<http://www.sla.org/content/SLA/pressroom/pressrelease/2003pressrelease/2302.cfm>
- [2003 Virtual Seminar Series](#) (02/06/2004)  
Learn more about **XML** and **The Value of the Information Professional** through upcoming virtual seminars!  
<http://www.sla.org/content/Learn/Learnmore/distance/virtsem2003/index.cfm>
- [January 24, 2003 Virtual Seminar](#) (02/04/2004)  
Need to know what tools are available to **Build, Maintain, and Grow** your information services? Find out during SLA's first Virtual Seminar of 2003  
<http://www.sla.org/content/Learn/Learnmore/distance/virtsem2003/jan24.cfm>
- [July 23, 2003 Virtual Seminar](#) (02/04/2004)  
<http://www.sla.org/content/Learn/Learnmore/distance/virtsem2003/july23.cfm>
- [Value of the Information Professional](#) (01/30/2004)  
<http://www.sla.org/content/Learn/ivalue/index.cfm>

Searched all header fields using entered keyword 'taxonomy'  
1 match found

- [Taxonomy](#) (02/07/2005)  
Information Portal on taxonomies.  
<http://www.sla.org/content/resources/infoportals/taxonomies.cfm>

Search again? Please visit our [search explanation page](#) for more information.

Search "taxonomy"  
in XML descriptor field  
returns all documents  
in that category → 27  
Search in original  
metadata → 1  
**Solution:**  
*Include descriptors  
with metadata!*



# Management

- Data management layers
- Curation
- Preservation
- Archiving
- Storage
- Choose tools wisely
- Data mashups



# Go – No Go

- Reach 85% precision to launch for productivity - assisted
- Reach 85% for filtering or categorization
  - Sorting for production
- Level of effort to get to 85%
- Integration into the workflow is efficient



# Learn to understand the parts

- Information infrastructure
- Information dynamics
- Tensions on data
- Design elements
- How and when to set the data framework
- Identifiers (taxonomies) Glue to hold data together



# Services

- Core cross disciplines
- Combo of humans and machines
- What is open
- What needs a gatekeeper
- Store locally
- Store in cloud
- Some things people do better



# Benchmarks

- 15 – 20% irrelevant returns / noise
- Amount of work needed to achieve 85% level
- How good is good enough?
  - Satisfice = satisfaction + suffice
  - How good is good enough?
  - How much error can you put up with?
- Hits, misses, noise



# ROI Calculations

- Assume – 5000 term thesaurus
  - 1.5 synonyms per terms
  - 7500 terms total
- Assume 85% accuracy
  - Use assisted for indexing
  - Use automatically for filtering
- Assume \$55 per hour for staff
- Assume 10000 records for test batch



# Co-occurrence - Training sets

- Collect 20 items per thesaurus entry
  - Preferred and non-preferred terms
  - Find records – could be programmatic
  - Need to collect and review about 60 to get 20 appropriate ones
  - Review records – ensure they are correct
  - 3 minutes average per final selected record
  - = 1 hour per term entry
  - 1 minute to review a record (20 resulting records)
  - 7500 terms at one hour per term = 7500 hours
  - Estimated at 7500 hours \$55 / hour = \$412,500





# Rulebuilding - rules

- 80% of rules built automatically
- $7500 \times .8 = 6000$
- 20% require complex rules
  - Average rule takes 5 minutes
  - (Actually MUCH faster using a rule building tool)
  - $5 \times 1500 = 7500$  minutes
  - $125 \text{ hours} \times \$55 = \$6875$



# ROI - Segments

- Cost of auto system
- Cost of getting system ready
- Ongoing maintenance
- Increased efficiency
- Increased quality of retrieval
- Cost of legacy system maintenance



# ROI – Co-occurrence

- Cost of auto system- \$500,000 (could be less, but the one used for this study cost this amount)
- Cost of getting system ready
  - Programming support and integration
    - Estimated at 2 weeks programming \$100 / hour = \$8000
  - Training sets
    - Estimated at 7500 hours \$55 / hour = \$412,500
    - Possible need to re-run training set several times
- Ongoing maintenance
  - Estimated at 15% of purchase price for license = \$75,000
  - Quarterly retraining to keep up with new terms –
  - Training sets for new terms 50 terms per quarter = 200 x \$55=\$11,000
- Increased efficiency
- Expected accuracy at 60%
- Increased quality of retrieval
- Cost of legacy system maintenance



# ROI – Rules system

- Cost of auto system- \$60,000
- Cost of getting system ready
  - Programming support and integration
    - Estimated at 2 weeks programming at \$100 / hour = \$8000
  - Rule building
    - Estimated at 125 hours at \$55 / hour = \$6850
    - Possible need to re run training set several times
- Ongoing maintenance
  - Estimated at 15% of purchase price for license = \$9000
  - Rule building for new terms, 50 terms per quarter
    - 200 terms x .8 = 160 automatic
    - 40 at 5 minutes per term = 200 minutes /60 = 3.33 hours x \$55 = \$183
- Increased efficiency
- Expected accuracy at 60%
- Increased quality of retrieval
- Cost of legacy system maintenance



# ROI - Rules

- Year one
  - $\$60,000 + \$8,000 + \$6,850 = \$74,850$
- Years thereafter
  - $\$9000 + \$183 = \$9183$
- 85% accuracy



# ROI – Co-occurrence

- Year one
  - $\$500,000 + \$8,000 + \$412,500 = \$920,500$
- Years thereafter
  - $\$75,000 + \$11,000 = \$86,000$
- 60% accuracy



# Summary on ROI

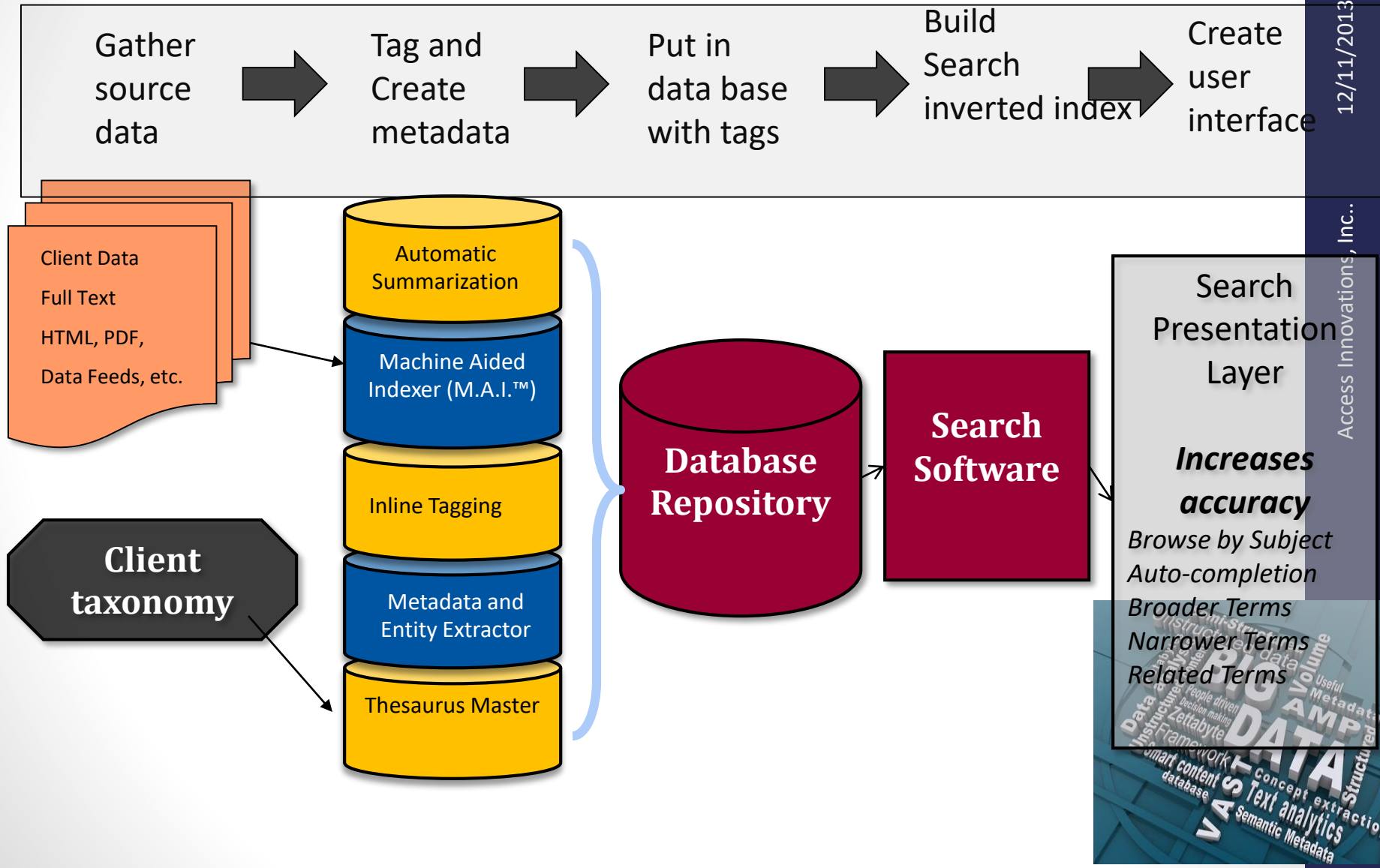
- Novelty detection requires inference techniques
- Is needed to discover new terms
- Controlled vocabulary / thesaurus can be more controlled and accurate approach
- Combination of the two would be best
  - Controlled – use rules
  - New terms - use inference





# Fully integrated with MOSS

## The Workflow



# Inline Tagging

oil rises as U.S. stimulus hopes outweigh weak demand oil prices rose on Thursday as hopes that the White House would move quickly on an economic stimulus package outweighed flagging demand and rising inventories in the world's top consumers. U.S. crude settled 12 cents higher at \$43.67 a barrel, after falling as low as \$40.41 earlier. London Brent settled at \$45.39 a barrel, up 37 cents. Earlier in the day, crude prices had dropped after a U.S. government report showed that crude oil, gasoline and distillate fuels rose last week as demand for fuels weakened again. The U.S. stock market pared early losses after a spokesman said President Obama's administration is committed to moving as quickly as possible on an economic stimulus plan. "crude is still resilient, despite the big build you have seen in the EIA data. I have a feeling that there will soon be a rebound in the stock market and that will spill over to the energy markets," said Mark Waggoner, president of Excel Futures in Huntington Beach, California. "In the last week, I've seen investor confidence improving, with the new Obama administration now installed ... I feel that if this confidence continues, it will spill into more (consumer) buying and that will improve demand for gasoline and other energy products," he added. oil fell in early trade after U.S. energy Information Administration data showed that crude oil inventories jumped 6.1 million barrels last week, well above expectations for a build of 1.4 million barrels.

## MAIstroTermList

Term Name="Crude oil (Commodity markets)" NumberOfTimes="3"  
Term Name="Petroleum resources" NumberOfTimes="3"  
Term Name="Oil (Fuels)" NumberOfTimes="3"  
Term Name="Energy resources (Commodity markets)" NumberOfTimes="3"  
Term Name="Product inventories" NumberOfTimes="3"  
Term Name="Gasoline" NumberOfTimes="2"  
Term Name="Prices" NumberOfTimes="2"  
Term Name="Commodity market prices" NumberOfTimes="2"  
Term Name="Stock market" NumberOfTimes="2"  
Term Name="California" NumberOfTimes="1"  
Term Name="Futures (Investments)" NumberOfTimes="1"  
Term Name="Commerce" NumberOfTimes="1"  
Term Name="Money and banking" NumberOfTimes="1"

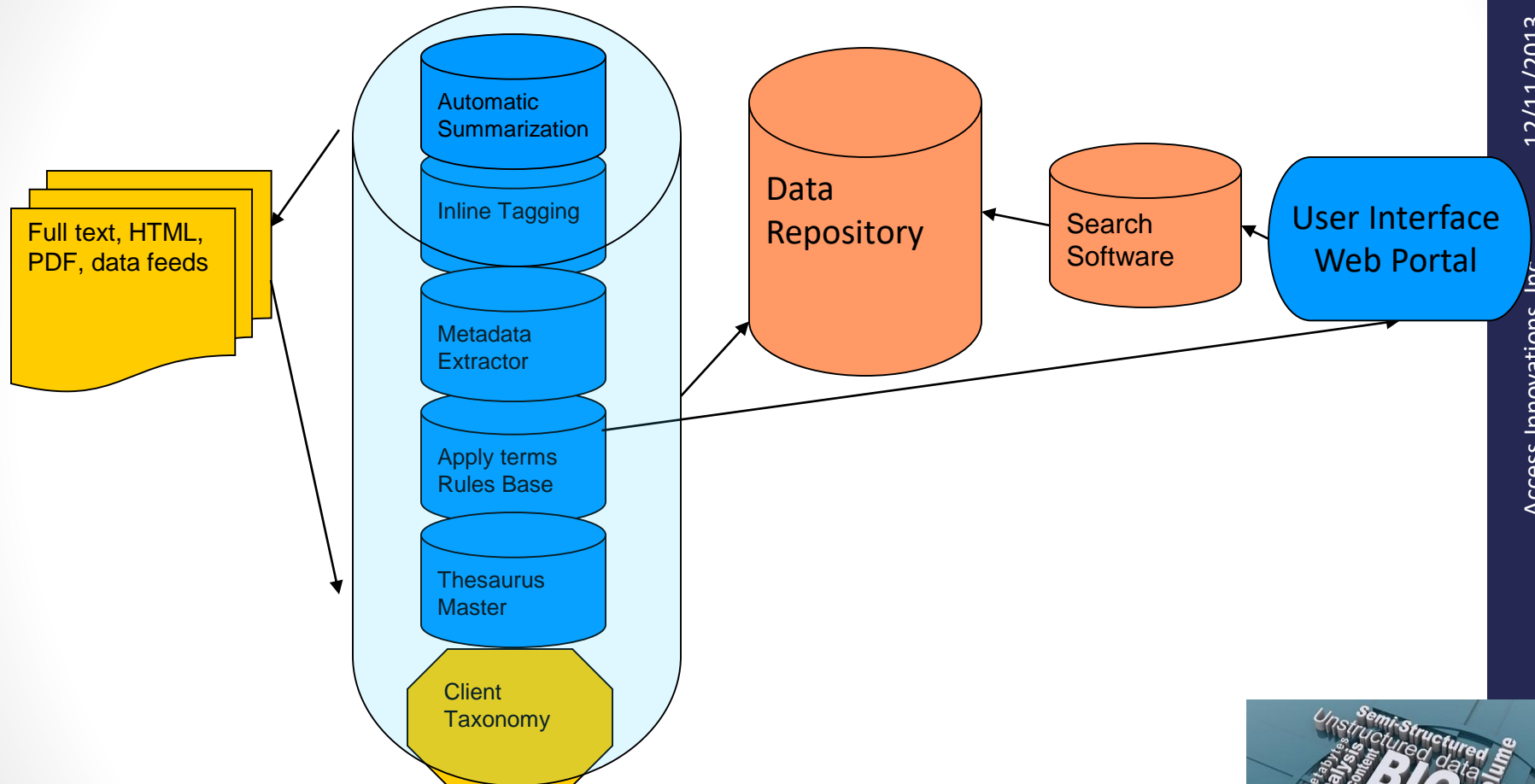
MAIstroTerm : Crude oil (Commodity markets)  
MAIstroTerm : Petroleum resources  
MAIstroTerm : Oil (Fuels)

Shows the exact point where the concept is mentioned

Mouse-over to view the term record

Statistical summary, showing the number of times each term is mentioned in the article

# Semantic Process



12/11/2013

Access Innovations, Inc.:



## Access Innovations, Inc..



# *Sample DOCUMENT*

## Creating an Inverted File Index

### Outline of Presentation

- 1 Define key terminology
- 2 Thesaurus tools
  - Features
  - Functions
- 3 Costs
  - Thesaurus construction
  - Thesaurus tools
- 4 Why & when?





## The terms from the “outline”

&

1

2

3

4

construction

costs

define

features

functions

key

of

outline

presentation

terminology

thesaurus

tools

when

why



# Complex inverted file index

## Placement location

& - Stop

1 - Stop

2 - Stop

3 - Stop

4 - Stop

construction - L7, P2, SH

costs - L6, P1, H

define - L2, P1, H

features - L4, P1, SH

functions - L5, P1, SH

key - L2, P2, H

of - Stop

outline - L1, P1, T

presentation - L1, P3, T

terminology - L2, P3, H

thesaurus - (1) - L3, P1, H

(2) - L7, P1, SH

(3) - L8, P1, SH

tools - (1) - L3, P2, H

(2) - L8, P2, SH

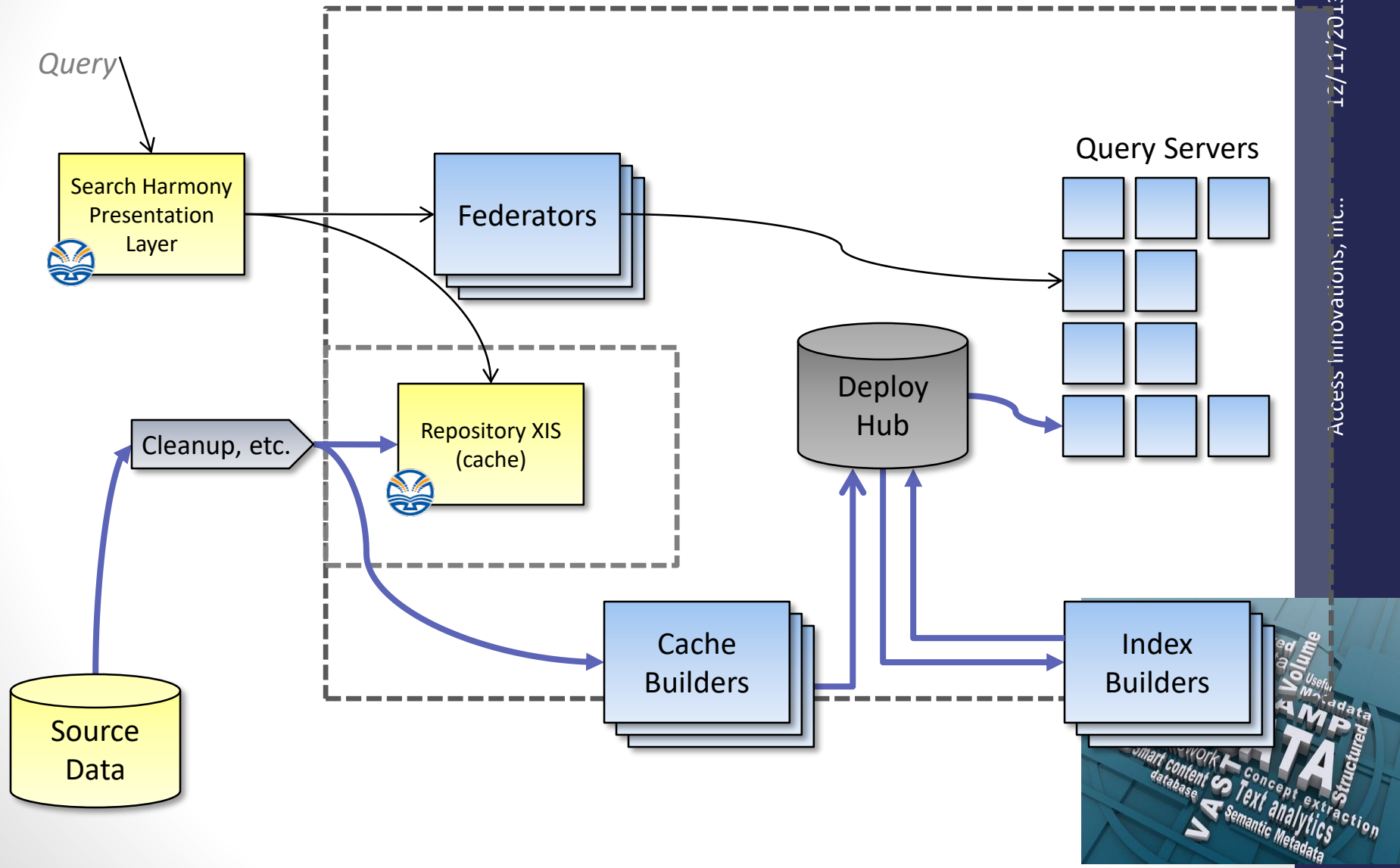
when - L9, P3, H

why - L9, P1, H





# Complex Search Farm



# What happens at the search presentation layer?

- That is what librarians usually look at.
- What are the options coming?
- How can we encourage useful changes?



# Semantic search options

Autocompletion  
using controlled  
terms

Full taxonomy  
view for Navigation  
and Browsing

Related, Broader,  
and Narrower  
Terms

**media sleuth**

Home | Company | Buyer's Guide | Contact

- \* Agriculture (895)
- \* Applied technologies (83)
- \* Business (2246)
- \* Computer and information science (2027)
- \* Economics (1177)
- \* Education (5211)
- \* Family and consumer sciences (13)
- \* Geography (1486)
- \* Health and wellness (449)
  - \* **Alternative health care (20)**
  - \* Care of persons with disabilities (19)
  - \* Child care (110)
  - \* Dental hygiene (22)
  - \* Disease prevention (142)
  - \* Elder care (65)
  - \* First aid (262)
  - \* Hazardous substances (157)
  - \* Health education (135)
  - \* Health facilities (26)
  - \* Lifestyle (110)
  - \* Men's health (66)
- \* Nutrition (650)
- \* Personal grooming (186)
- \* Personal health and hygiene (404)
- \* Physical fitness (392)
- \* Public health (126)
- \* Rest (32)
- \* Safety (2732)
- \* Sexual hygiene (39)
- \* Sleep (61)
- \* Smoking (209)
- \* Substance abuse prevention (318)
- \* Women's health (129)
- \* History (7801)

**MediaSleuth**

Results for **Alternative health care**

**Conversation**  
addressing the  
Descriptors: **Alternative health care** - Medicine  
<http://www.media-sleuth.com/revmail/bibrec.php?accnum=1>

**Homeopathy**  
Introduces views  
the founder, Sam  
homeopathy, in  
Descriptors: **Alternative health care** - Medicine  
<http://www.media-sleuth.com/revmail/bibrec.php?accnum=1>

Explores the mixed population of homeopathic practitioners and profiles the ideal patient. Looks at the internal healing mechanism, the three dimensions of belief for cure, and the responsibilities of the patient.  
Descriptors: **Alternative health care** - Medicine  
<http://www.media-sleuth.com/revmail/bibrec.php?accnum=1>

**Homeopathy in America**  
Offers a perspective on the rise and decline of homeopathic medicine in the United States during the 19th century. Looks at homeopathy's resurgence and visits the National Center for Homeopathy in Washington D.C.  
Descriptors: **Alternative health care** - Medicine  
<http://www.media-sleuth.com/revmail/bibrec.php?accnum=1>

**Movement**  
Demonstrates the benefits of movement therapy for various conditions, including arthritis and stress

**psy**

- Psychiatric hospitals (nonpreferred)
- Psychic phenomena (nonpreferred)
- Psychological dependence (nonpreferred)
- Psychological disorders (nonpreferred)
- Psychopathology (nonpreferred)
- Psychotropic drugs (nonpreferred)
- Abnormal psychology
- Adolescent developmental psychology
- Adult developmental psychology
- Animal behavioral psychology
- Animal psychology
- Child developmental psychology
- Clinical psychology
- Cognitive psychology
- Community psychology
- Clinical psychopharmacology

99 docs (0.0 seconds)

**> Expand your search**  
-Thesaurus Related Terms

[Folk medicine](#)  
[Medical care](#)  
[Medical treatment](#)  
[Western medicine](#)

so features  
users of

**> Target your search**  
-Thesaurus Narrower Terms

[Acupressure](#)  
[Acupuncture](#)  
[Aromatherapy](#)  
[Ayurveda](#)  
[Chiropractic medicine](#)  
[Color therapy](#)  
[Faith healing](#)  
[Herbal medicine](#)  
[Holistic medicine](#)  
[Homeopathy](#)  
[Massage and therapeutic touch](#)  
[Naturopathy](#)  
[Osteopathy](#)  
[Self healing](#)

OPEN ACCESS | PEER-REVIEWED  
RESEARCH ARTICLE

3,579  
VIEWS

10  
CITATIONS

79  
SAVES

# Predation by Bears Drives Senescence in Natural Populations of Salmon

Stephanie M. Carlson, Ray Hilborn, Andrew P. Hendry, Thomas P. Quinn

Published: Dec 12, 2007 • DOI: 10.1371/journal.pone.0001286

Article

About the Authors

Metrics

Comments

Related Content

Download PDF

Print

Share

## Abstract

Introduction

Results

Discussion

Materials and Methods

Supporting Information

Acknowledgments

Author Contributions

References

Reader Comments (2)

Figures

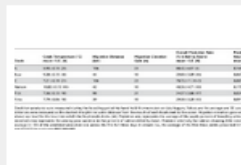
## Abstract

Classic evolutionary theory predicts that populations experiencing higher rates of environmentally caused ("extrinsic") mortality should senesce more rapidly, but this theory usually neglects plausible relationships between an individual's senescent condition and its susceptibility to extrinsic mortality. We tested for the evolutionary importance of this condition dependence by comparing senescence rates among natural populations of sockeye salmon (*Oncorhynchus nerka*) subject to varying degrees of predation by brown bears (*Ursus arctos*). We related senescence rates in six populations to (1) the overall rate of extrinsic mortality, and (2) the degree of condition dependence in this mortality. Senescence rates were determined by modeling the mortality of individually-tagged breeding salmon at each site. The overall rate of extrinsic mortality was estimated as the long-term average of the annual percentage of salmon killed by bears. The degree of condition dependence was estimated as the extent to which bears killed salmon that exhibited varying degrees of senescence. We found that the degree of condition dependence in extrinsic mortality was very important in driving senescence: populations where bears selectively killed fish showing advanced senescence were those that senesced least rapidly. The overall rate of extrinsic mortality also contributed to among-population variation in senescence-but to a lesser extent. Condition-dependent susceptibility to extrinsic mortality should be incorporated more often into theoretical models and should be explicitly tested in natural populations.

## Figures



Site	Year	Senescence Rate	Extrinsic Mortality	Condition Dependence
1	2000	0.15	0.25	0.10
2	2001	0.20	0.30	0.15
3	2002	0.25	0.35	0.20
4	2003	0.30	0.40	0.25
5	2004	0.35	0.45	0.30
6	2005	0.40	0.50	0.35



## Subject Areas

- Bears
- Senescence
- Predation
- Fishes
- Freshwater fish
- Death rates
- Marine fish
- Animal sexual behav...

ADVERTISEMENT





Firefox

W Plasmid - Wikipedia, the free encyclo... x PLOS PLOS ONE: Predation by Bears Drives ... x

www.plosone.org/article/info%3Adoi%2F10.1371%2Fjournal.pone.0001286

Google

Wikipedia Webmail JSTOR Dictionary.com Alb. Weather nmroads Sony OCIS codes Google Drive AAAS

ADVERTISEMENT

The easiest NGS Library preparation you will ever find!  
Just 1 tube, 2 hours, and 3 easy steps  
Learn more

**diagenode**  
Innovating Epigenetic Solutions

plos.org create account sign in

**PLOS** | ONE

Subject Areas For Authors About Us

Search advanced search

Browse Subject Areas: **Biology and life sciences / Zoology / Animal behavior**

<

Taxonomy

Theoretical biology (5)

Toxicology

Veterinary science

**Zoology**

View All Articles (4054)

**Animal behavior**

Animal cognition (97)

Animal phylogenetics (359)

Animal physiology

^

Collective animal behavior (30)

Foraging (466)

Grazing (114)

Hibernation (45)

Hunting behavior (37)

v

OPEN ACCESS PEER-REVIEWED

RESEARCH ARTICLE

3,579

10

79

VIEWS

CITATIONS

SAVES

# Predation by Bears Drives Senescence in Natural Populations of Salmon

Stephanie M. Carlson , Ray Hilborn, Andrew P. Hendry, Thomas P. Quinn

Published: Dec 12, 2007 • DOI: 10.1371/journal.pone.0001286

Article

About the Authors

Metrics

Comments

Related Content

Download PDF

Print

Share

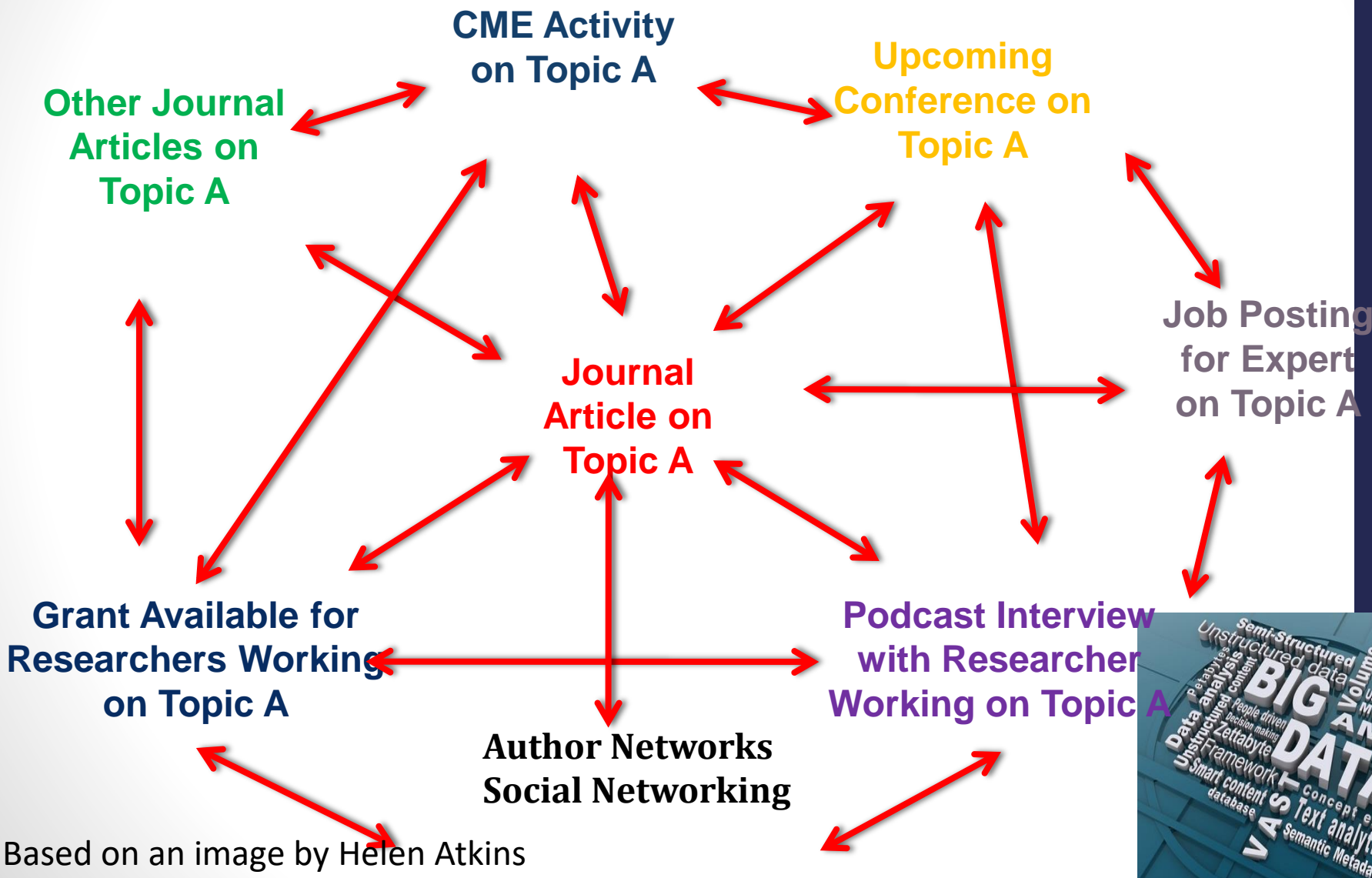
► Abstract

Abstract

# Data linked and served by metadata

12/11/2013

Access Innovations, Inc..



## **Cancer Epidemiology Biomarkers & Prevention**

Vol. 12, 161-164,

February 2003

© 2003 [American Association for Cancer Research](#)

### **Short Communications**

#### **Alcohol, Folate, Methionine, and Risk of Incident Breast Cancer in the American Cancer Society Cancer Prevention Study II Nutrition Cohort**

**Heather Spencer Feigelson<sup>1</sup>, Carolyn R. Jonas, Andreas S. Robertson, Marjorie L. McCullough, Michael J. Thun and Eugenia E. Calle**

Department of Epidemiology and Surveillance Research, American Cancer Society, National Home Office, Atlanta, Georgia 30329-4251

Recent studies suggest that the increased risk of breast cancer associated with alcohol consumption may be reduced by adequate folate intake. We examined this question among 66,561 postmenopausal women in the American Cancer Society Cancer Prevention Study II Nutrition Cohort.

[Prevention Working Groups](#)

• [Finance](#)

• [Charter](#)

• [Molecular Epidemiology](#)

#### **Related Awards**

• [AACR-GlaxoSmithKline Clinical Cancer Research Scholar](#)

#### **Awards**

• [ACS Award](#)

• [Weinstein Distinguished Lecture](#)

#### **Think Tank Report**

[Related Think Tank Report Content](#)

#### **Webcasts**

[Related Webcasts](#)

#### **Related Press Releases**

- [How What and How Much We Eat \(And Drink\) Affects Our Risk of Cancer](#)
- [Novel COX-2 Combination Treatment May Reduce Colon Cancer Risk Combination Regimen of COX-2 Inhibitor and Fish Oil Causes Cell Death](#)
- [COX-2 Levels Are Elevated in Smokers](#)

#### **Related AACR Workshops and Conferences**

- [Frontiers in Cancer Prevention Research](#)
- [Continuing Medical Education \(CME\)](#)
- [Molecular Targets and Cancer Therapeutics](#)

#### **Related Meeting Abstracts**

- [Association between dietary folate intake, alcohol intake, and methylenetetrahydrofolate reductase C677T and A1298C polymorphisms and subsequent breast](#)
- [Folate, folate cofactor, and alcohol intakes and risk for colorectal adenoma](#)
- [Dietary folate intake and risk of prostate cancer in a large prospective cohort study](#)

#### **Related Education Book Content**

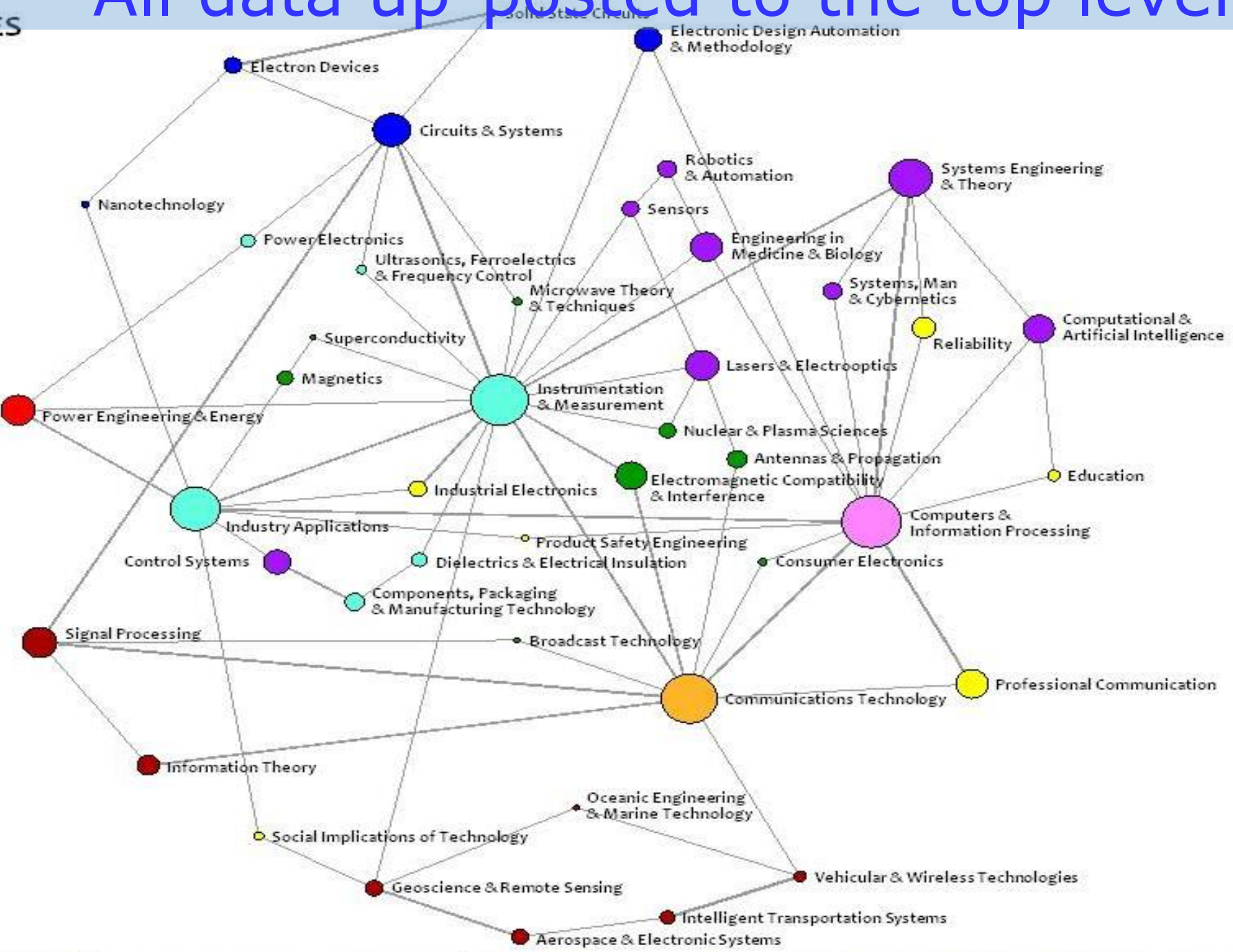
[Oral Contraceptives, Postmenopausal Hormones, and Breast Cancer](#)

[Physical Activity and Cancer](#)

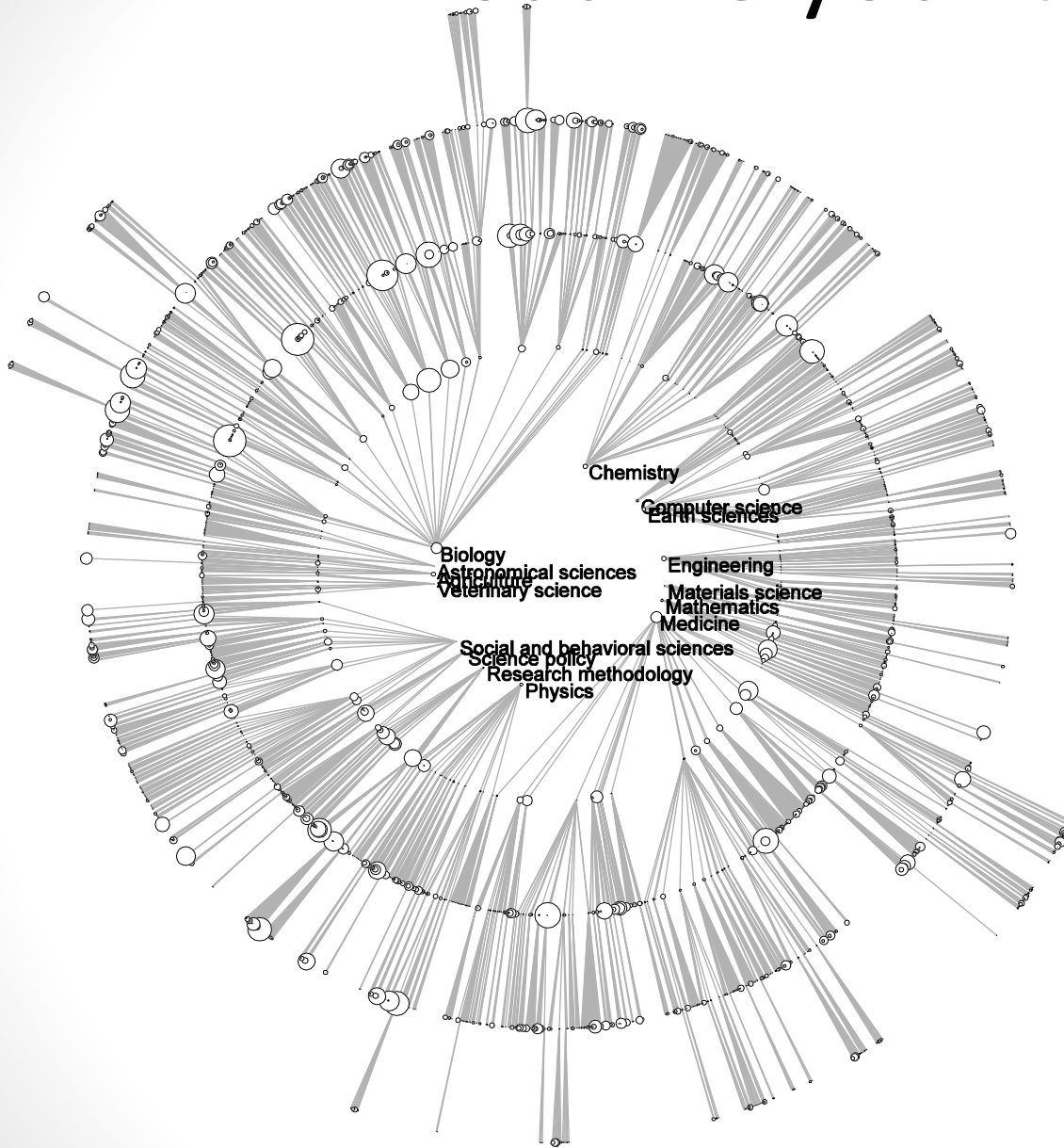
[Hormonal Interventions: From Adjuvant Therapy to Breast Cancer Prevention](#)



# All data up-posted to the top level



# Visualize your tagged data



This is a radial graph of “plothes”. The number of records for which each index term occurs is reflected by circle sizes.

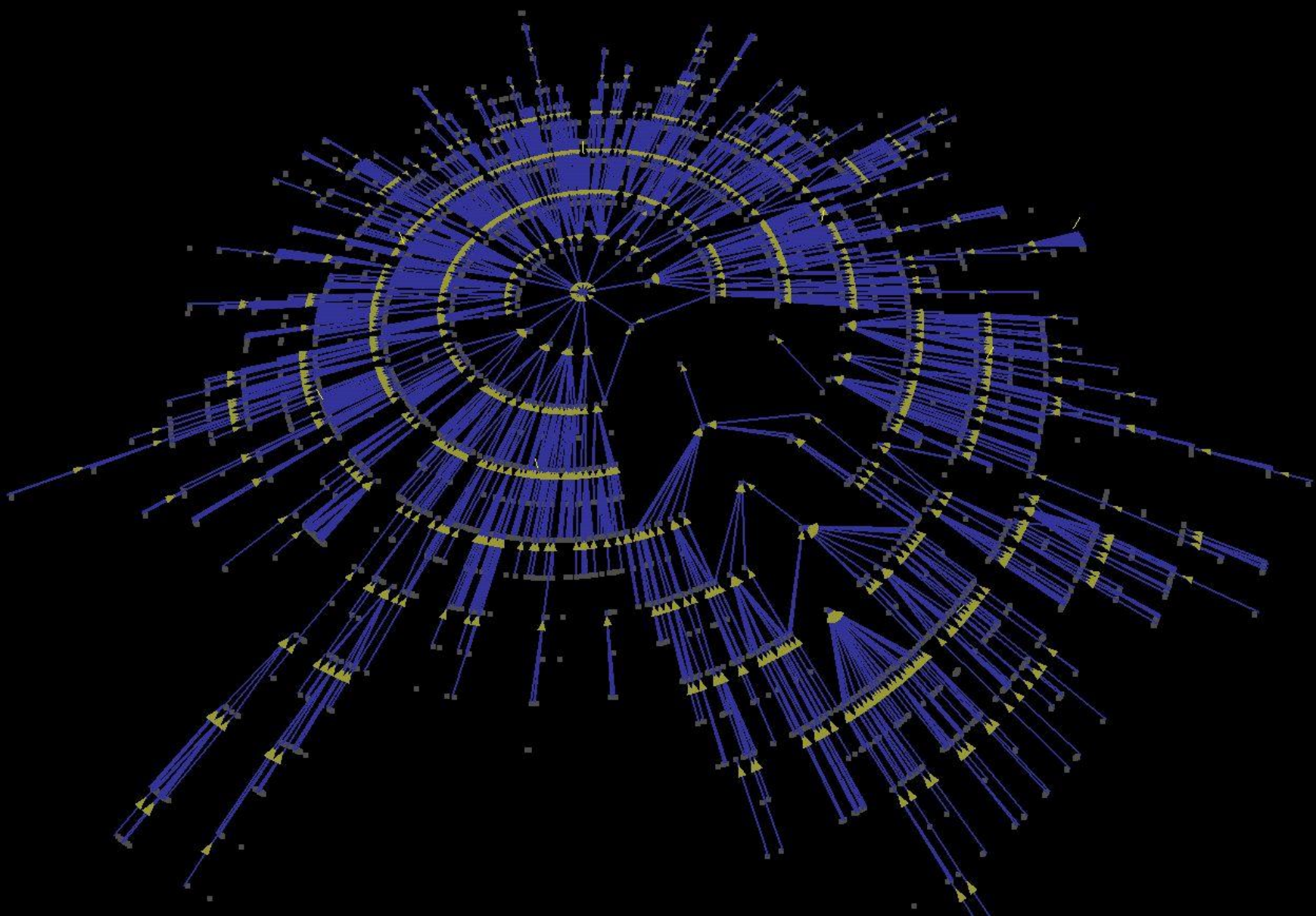




# Semantic Hierarchy



browsable tree

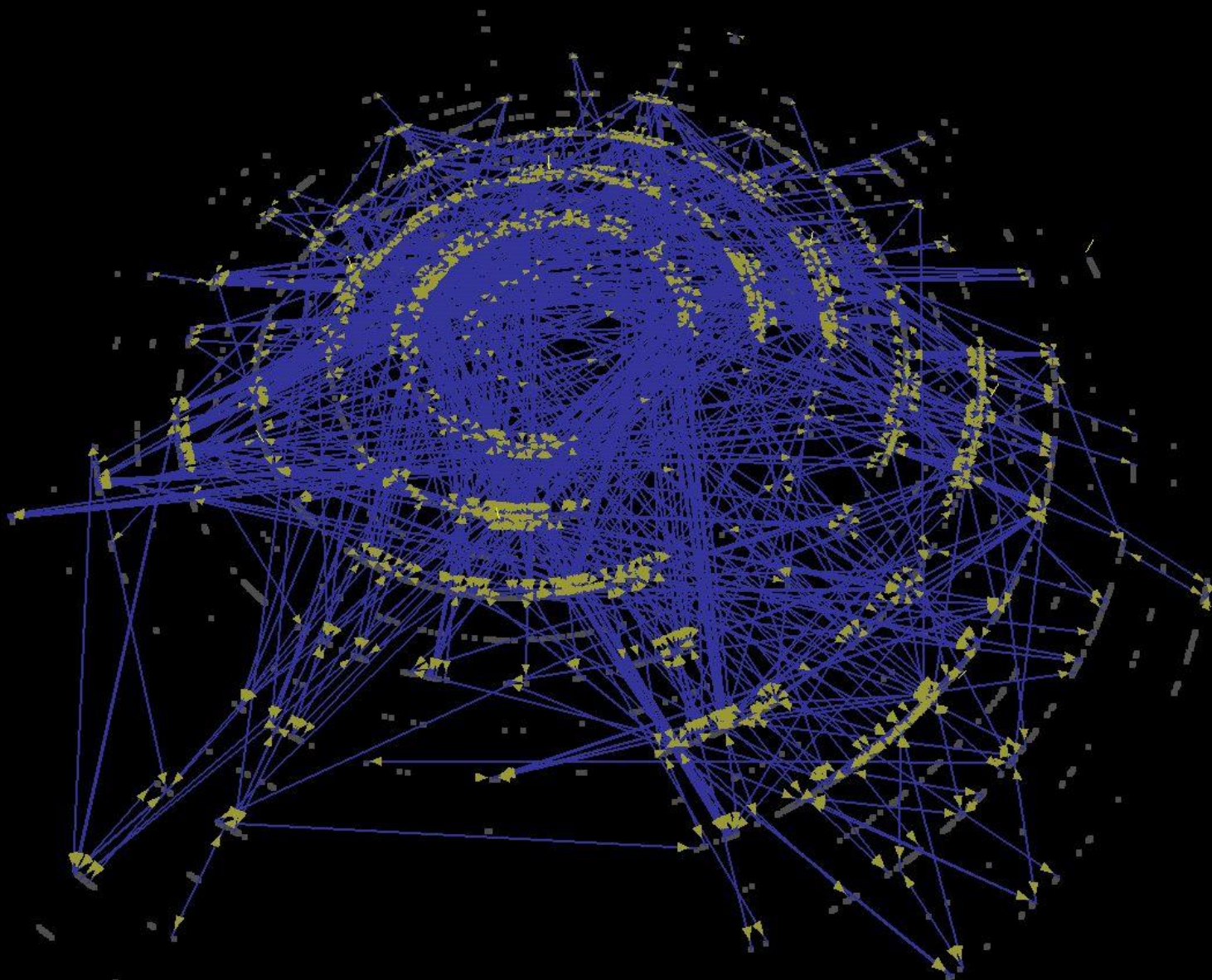


12/11/2013

Access Innovations, Inc..



# Related Terms → semantic web



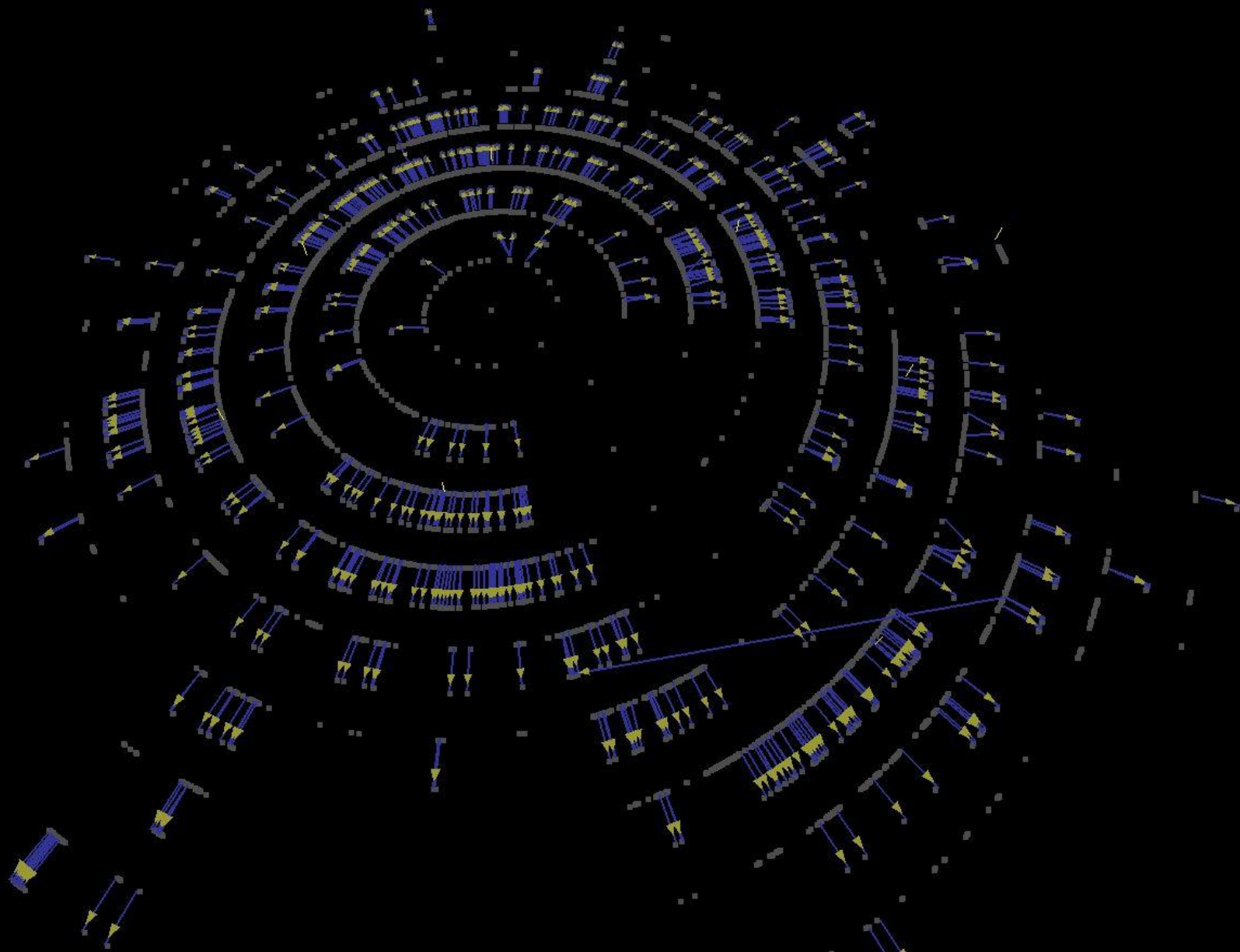
12/11/2013

Access Innovations, Inc..





# Synonyms → search foundation



# Load to a visualization program such as Prefuse

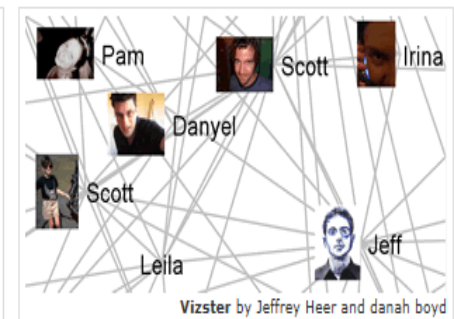
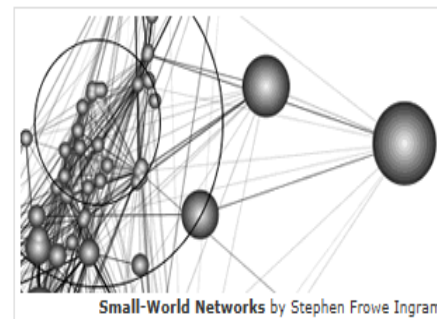
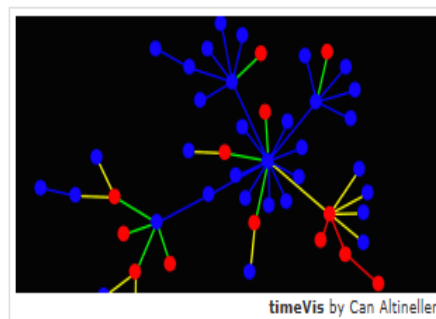
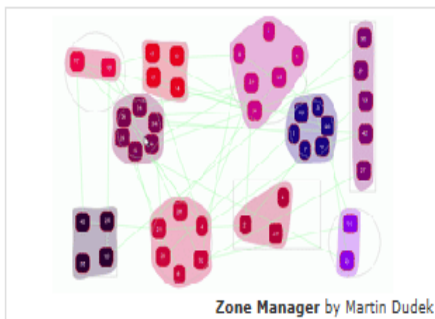
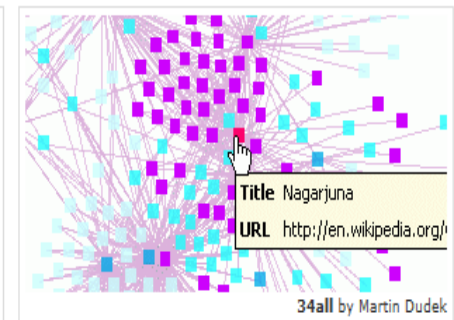
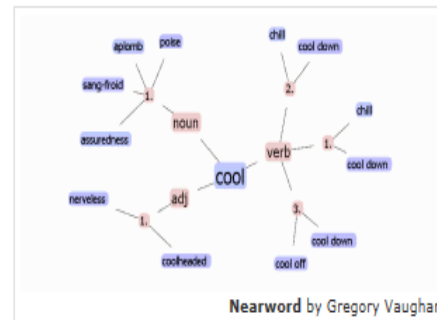
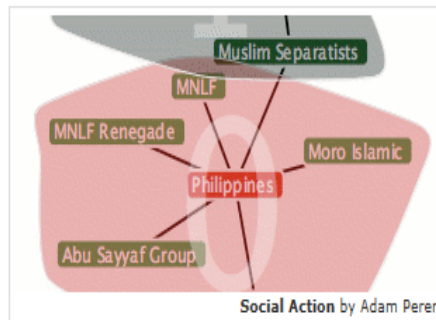
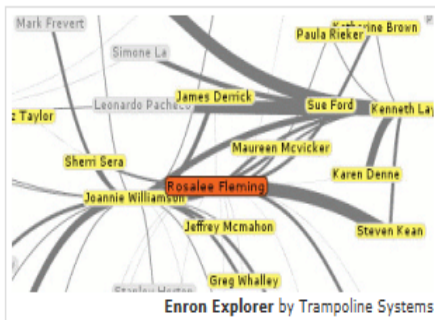
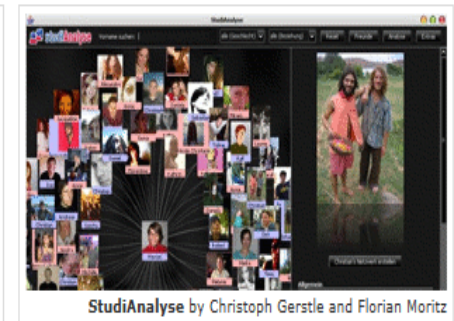
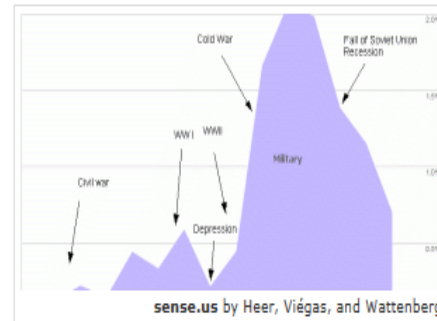
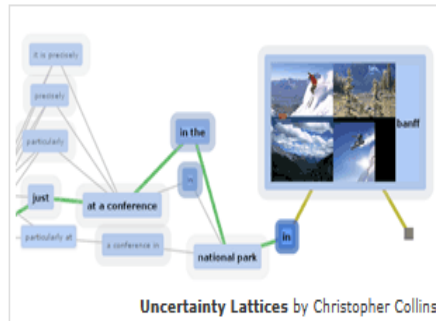
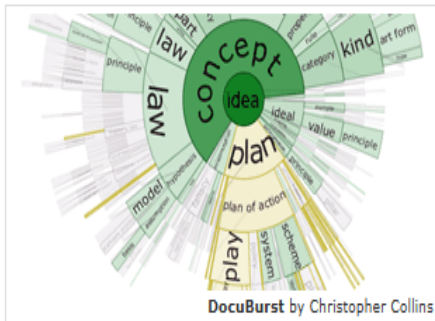
**prefuse**

INFORMATION VISUALIZATION TOOLKIT

visualization gallery

[Home](#) | [Download](#) | [Gallery](#) | [Documentation](#) | [FAQ](#)

Search



# Data Mashups

- Drawing together information from several sources
- Overlay on additional surfaces – like maps
- Fun distribution of data
- Cornell School of Ornithology
  - Migration patterns
  - YouTube
  - Citizen Science





## Access Innovations, Inc..

- Access Innovations, Inc..

Access Innovations, Inc..



Access Innovations, Inc..



Access Innovations, Inc..

# Authors at a Place

Fly To Find Businesses Directions

Fly to e.g., 37 25' 19.1"N, 122 05' 06"W

Cancer researchers

Cancer researchers (1 - 10)

- ☒ Sponsored Links  
[Diabetes research](#)  
[www.informit.com.au/health](http://www.informit.com.au/health)  
& other core health research.  
Search & download online at
- ☒ **A** [Memorial Sloan-Kettering Cancer](#)  
East Drive, New York, NY 10024  
(212) 639-2000
- ☒ **B** [Cancer Epidemiology Biomarkers](#)  
615 Chestnut St # 17,  
Philadelphia, PA 19106-4406
- ☒ **C** [University Of Kentucky MD: Com](#)  
300 North Broadway Road,  
Lexington, Kentucky 40508

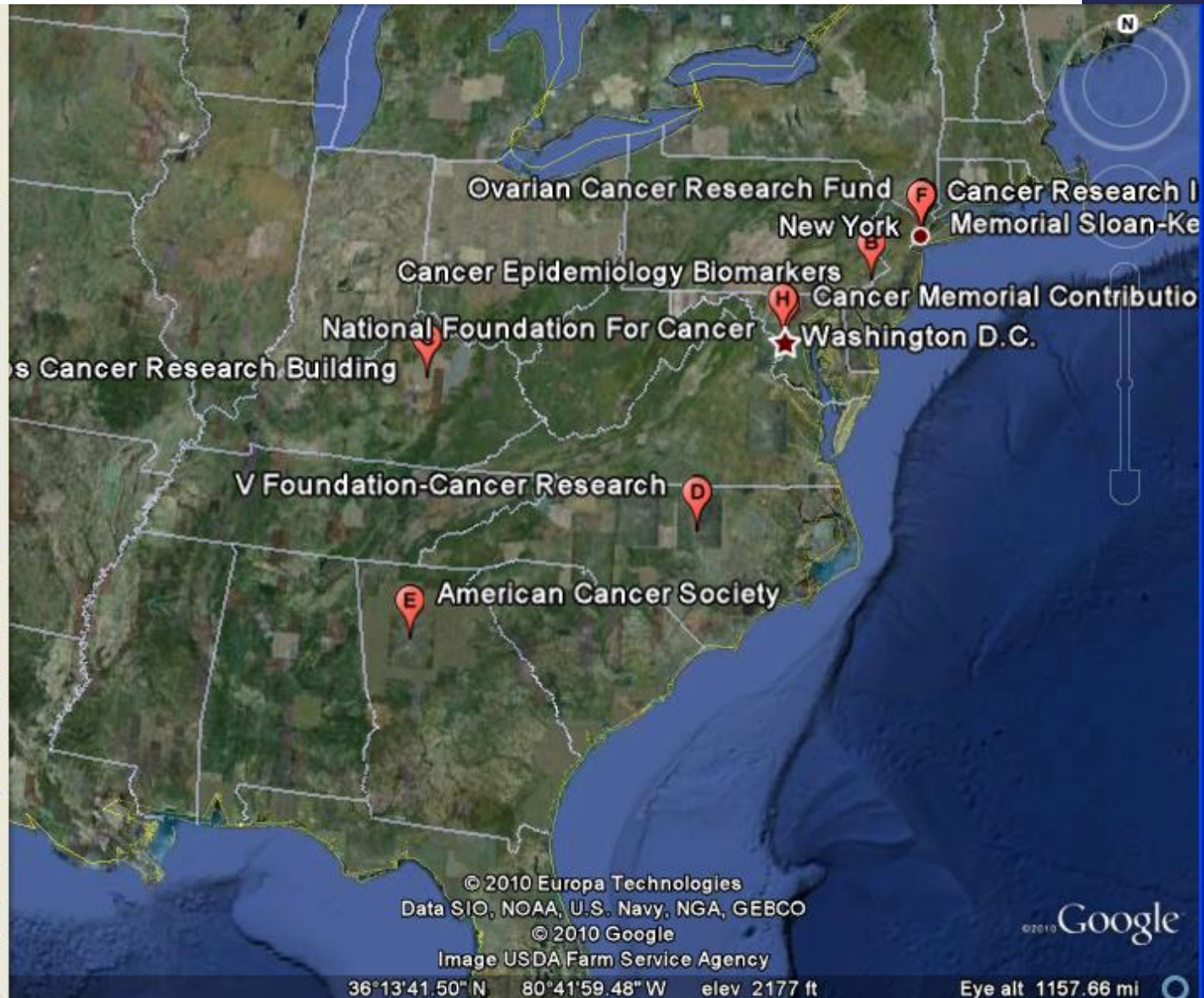
Places

- ☒ My Places
- ☒ [Sightseeing Tour](#)  
Make sure 3D Buildings  
layer is checked
- ☐ Temporary Places

Layers

Earth Gallery >>

- ☒ Primary Database
- ☒ Borders and Labels
- ☒ Places



# Future?

- Document systems replaced
- New infrastructure
- Institutional repositories not good long term for Big Data
- Need to operate at scale
- Integrated, ecosystem to infrastructure
- Replace customized human mediated
  - With interpretive layer – computer assisted
- Links to data
- Open data / supplemental data
- Too much is not enough!





# Differentiation

- Not by size of collection
- By the services they offer
- More services more competitive
- Spreadsheet science
  - Go down hall and ask
- Big Data
  - Too many people to ask
  - Lose provenance
  - Hard to differentiate the layers



# New ways of doing things

- People, data usage
- Preservation and discovery
- Publishers' content, not format
- Visualization of Big Data sets
- Reproducibility of research
- Shared knowledge
- Open peer review
  - Researcher.org
  - Galaxy zoo
  - Zooniverse
  - 900,000 citizen scientists
  - Wikipedia



# Few have witnessed what you're about to see

Experience a privileged glimpse of the distant universe as observed by the SDSS, the Hubble Space Telescope, and UKIRT



We are trying something new! Come help us understand a very specific type of galaxy and experience science from start to end. [Take part](#)

## Classify Galaxies

To understand how galaxies formed we need your help to classify them according to their shapes. If you're quick, you may even be the first person to see the galaxies you're asked to classify.

[Begin Classifying](#)



# Skills we bring

- Librarians
  - Weeding
  - Redefining collections
  - Collection development
- Skills to apply to Big Data
  - Vocabulary development
  - Search expertise
  - Reference ability
- What to throw away
- What to keep for discoverability
- Need metadata to preserve
  - Better discoverability
  - Easier preservation
  - Keep the data provenance





# We covered

- Big Data – What is it?
- New Government Initiative
- Content organization
- Discovery (Search)
- Management
- Skills we bring
- Examples of what we can do



# Who are we?

- Access Innovations –
  - Semantic Enrichment Services Provider
  - We change search to find!
- Creator of Data Harmony tools
- More than 600 taxonomies created
- More than 2000 engagements
- Financed by Sweat, Persistence, and Good Cash Flow Management
- Accurate, on time, under budget!



Slides at  
[www.accessinn.com/presentations](http://www.accessinn.com/presentations)

# Thank you

- Marjorie Hlava
  - [mhlava@accessinn.com](mailto:mhlava@accessinn.com)
- +1-505-998-0800
- Access Innovations/Data Harmony/Access Integrity
- Taxodiary.com (blog)
- Taxobank.org (reference tool for taxonomists)

