

Realizing the Potential of Research Data

Carole L. Palmer

Information School
University of Washington

Coalition for Networked Information
14 April 2015

- Are we the experts we need to be?
- What are the exemplar for data resources and services?
- Can we learn and lead at the scale and pace needed?

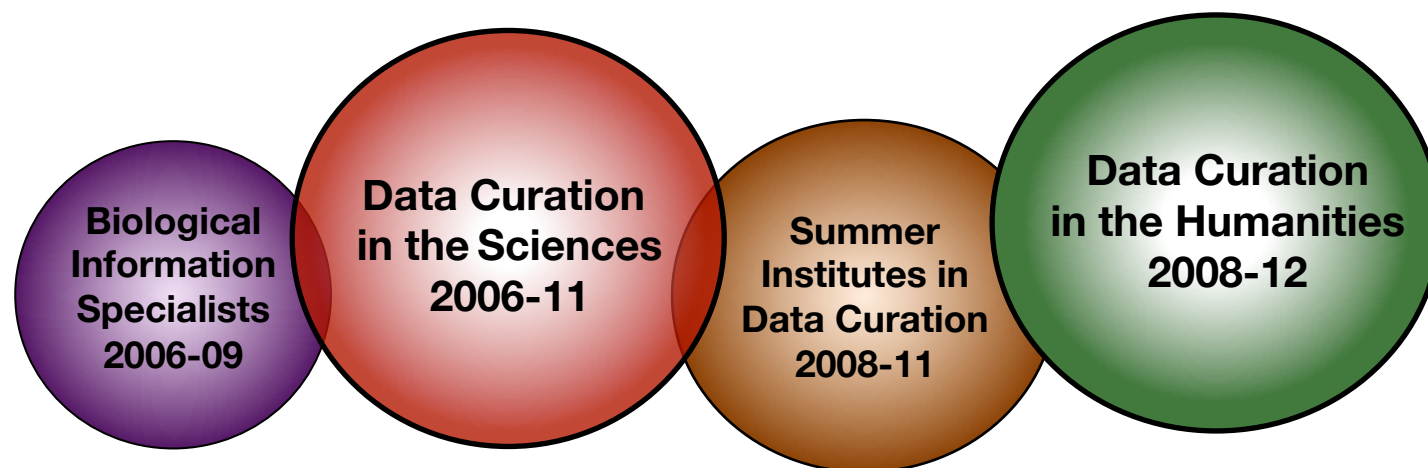
October 16-17, 2008
Arlington, Virginia 22202

Sponsored by [ARL](#) & [CNI](#)

ARL-CNI Forum

2008

Preparing e-Science Information Specialists: New Programs and Professionals



Well positioned—institutional and human infrastructure, expertise, commitment

*Going forward, **must not underestimate challenge.***

Deluge of discourse and directives

2003 at least 11 reports from
NSB, NRC, NSF

Early leadership from
Information Schools

Atkins

Unsworth

Larsen

Building the Infrastructure
for Cyberscholarship

**Understanding Infrastructure:
Dynamics, Tensions, and Design**



Report of a Workshop on "History & Theory of Infrastructure:
Lessons for New Scientific Cyberinfrastructures"

Paul N. Edwards
Steven J. Jackson
Geoffrey C. Bowker
Cory P. Knobel

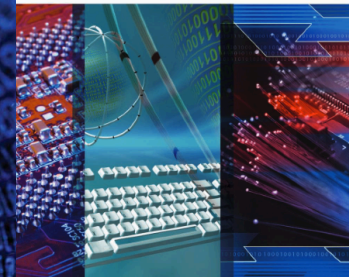
January 2007



**HARNESSING THE POWER
of
DIGITAL DATA
for
SCIENCE and SOCIETY**

Review of e-Science 2009

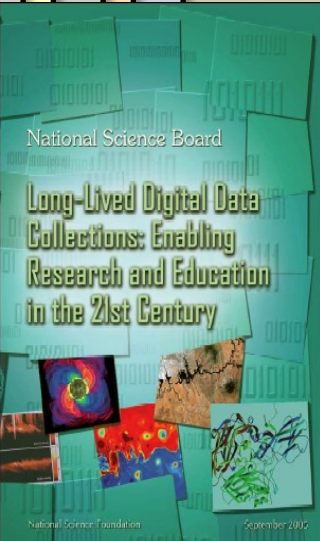
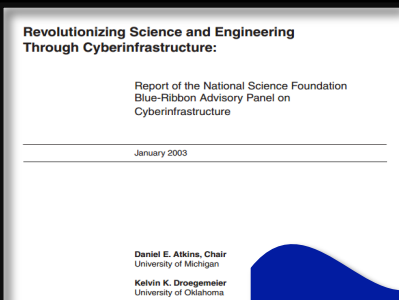
FROM A UK FOUNDATION FOR THE TRANSFORMATIVE
IMPACT OF RESEARCH AND INNOVATION



Riding the wave

How Europe can gain from the rising tide of scientific data

Final report of the High Level Expert Group on Scientific Data
Submission to the European Commission
October 2010



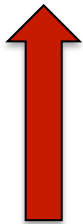
Deluge of repositories and standards

re3data.org
REGISTRY OF RESEARCH DATA REPOSITORIES

Home Search Browse Suggest FAQ About Schem

TAG ARCHIVES: DATABIB

Over 1,000 research data repositories



biosharing

665 “databases”
584 standards

Showing records 1 - 50 of 584.

Standard Type	Publications	Systems	Taxa
TERMINOLOGY ARTIFACT	337	0	No taxa defined.
EXCHANGE FORMAT	176	0	No taxa defined.
REPORTING GUIDELINE	71	0	No taxa defined.

Domains

Domain	Count
DNA	54
ANATOMY	44

Marine Metadata Interoperability

MAIN NAVIGATION

Home
MMI Guides
MMI Semantic Framework
Community
Vocabularies & Standards
Metadata Tools
Projects & Organizations
Search MMI References
Events
Add Content to Site
About MMI

USER LOGIN

Vocabularies & Standards References

215+ standards

MMI provides a set of references to various standards, vocabularies, ontologies, thesauri, services, and formats that are of particular interest to the marine science community or those working on interoperability. MMI has evaluated many of the listed references and provides a description, an informal characterization of its maturity, in some cases, and link to a primary resource where available. If you would like to see a reference added, please [contact us](#) or [sign up for an MMI account](#) and [add it yourself](#). You may also [search all references](#).

MMI has categorized its references in the following sets: Vocabularies, Ontologies and Thesauri, Content Standards, and Services, Protocols and Formats.

Title	Description	Reference Type	Reference Topics
A Universal Ontology for Sensor Networks Data	A prototype ontology using IEEE SUMO	Ontologies and Thesauri	Convention Topics Resource Discovery Sensors, Instruments and Platforms

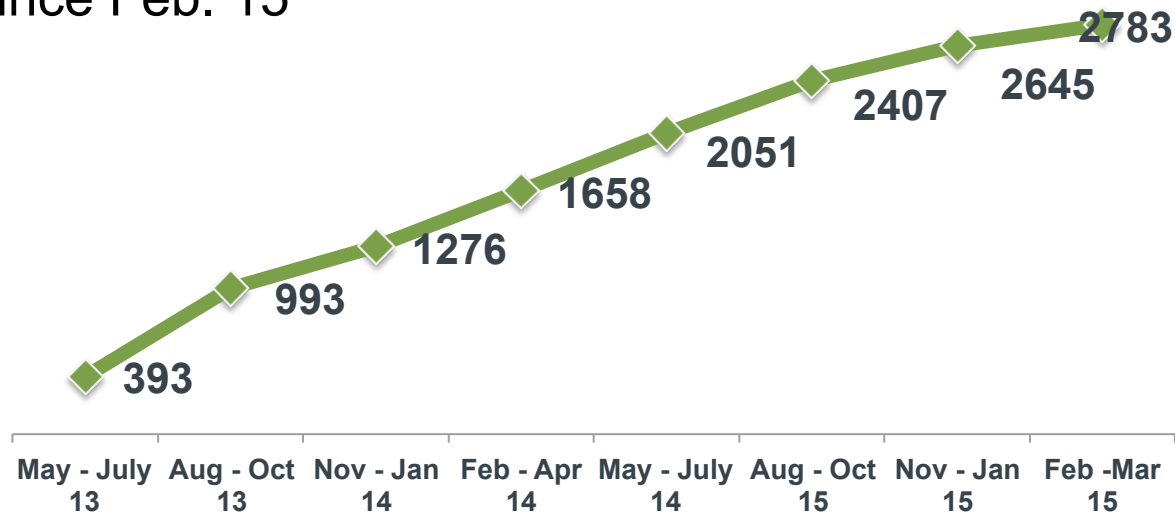


Reduce barriers to data sharing

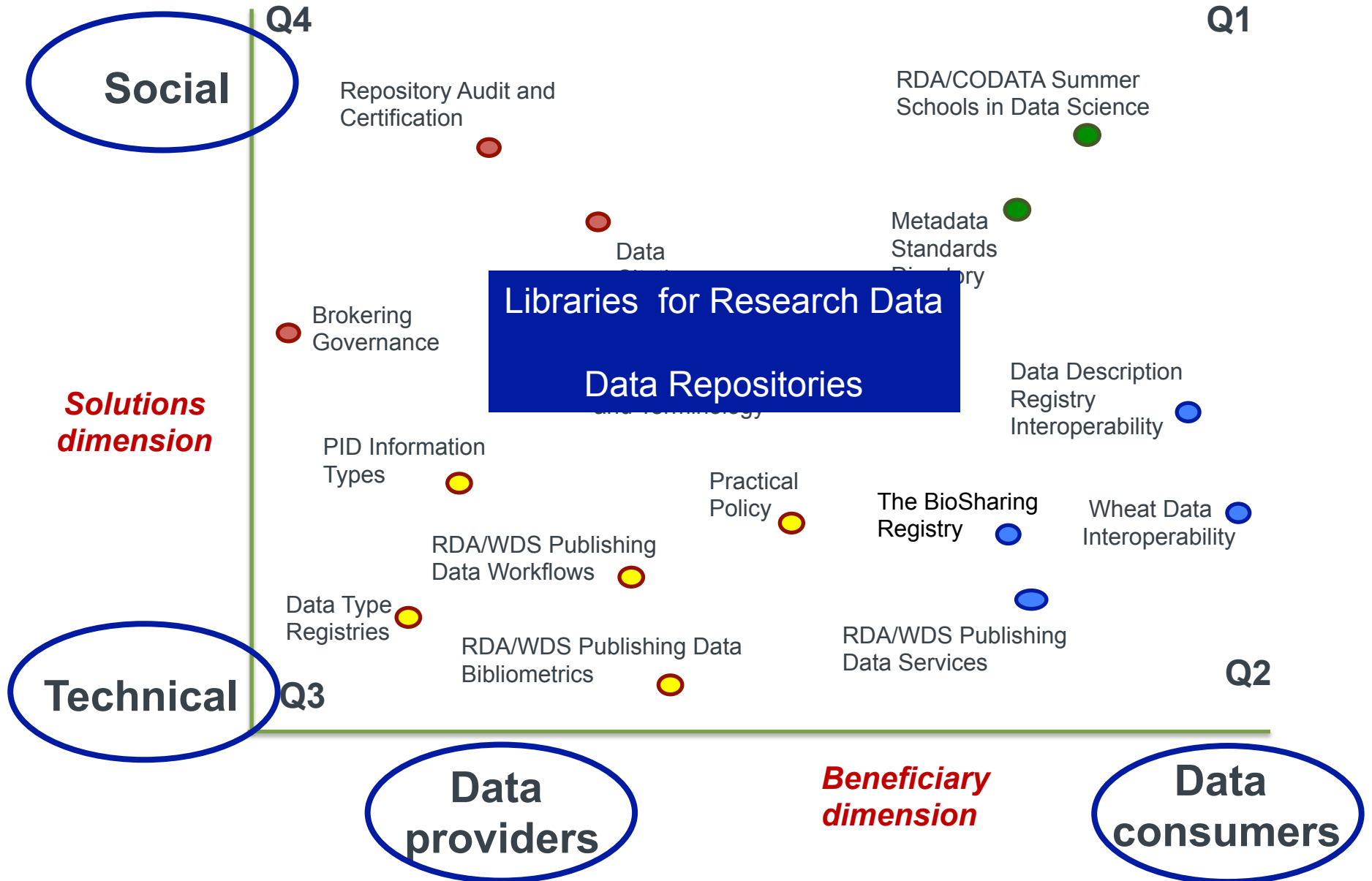
Accelerate coordinated global data infrastructure

95 countries - 50% Europe, 37% US

2783 RDA Community Members
+138 since Feb. 15



56 Working and Interest Groups



National data services



**The National
DATA SERVICE**



UK Data Service



CSC-IT CENTER FOR SCIENCE

Data Archiving and Networked Services



Abundance of data science initiatives



Stanford Data Science Initiative

UNIVERSITY OF ROCHESTER



INSTITUTE FOR DATA SCIENCE

Data Science Initiative

Data Sciences Initiative



Data Science Institute



UCI Data Science Initiative

2014 - UW eScience Celebration of data intensive research



Data Science Kickoff Session:
137 posters from 30+ departments and units

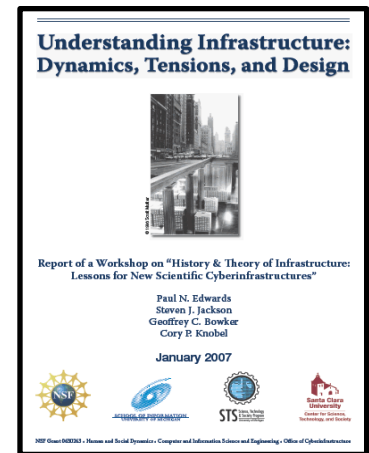
<http://escience.washington.edu/>

Problem: dynamics of systems and networks

from homogeneous, centralized, local

to heterogeneous, distributed, coordinated

- consolidation
- gateways for interoperation



(Edwards, et al., 2007)

"make-or-break" phase (Parsons & Berman, 2013)

Early choices constrain options

Institutions as intellectual habitat

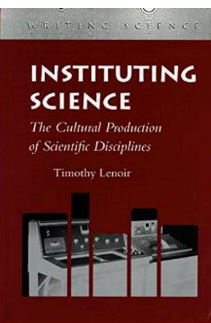
Research programs of researchers

- extend and legitimate products of work
- dominate cycles of credit and resources

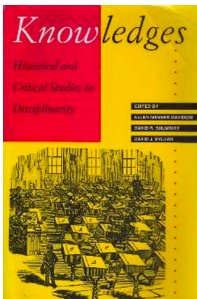
Institutions support

new routines long enough for distinctive types of work to emerge

- establish service roles
- facilitate links with other disciplines
- enable transmission of techniques and information



Lenoir, Timothy. 1993. "The Discipline of Nature and the Nature of Disciplines." In *Knowledges: Historical and Critical Studies in Disciplinarity*, edited by Ellen Messer-Davidow, David R. Shumway, and David J. Sylvan. Charlottesville: University Press of Virginia.

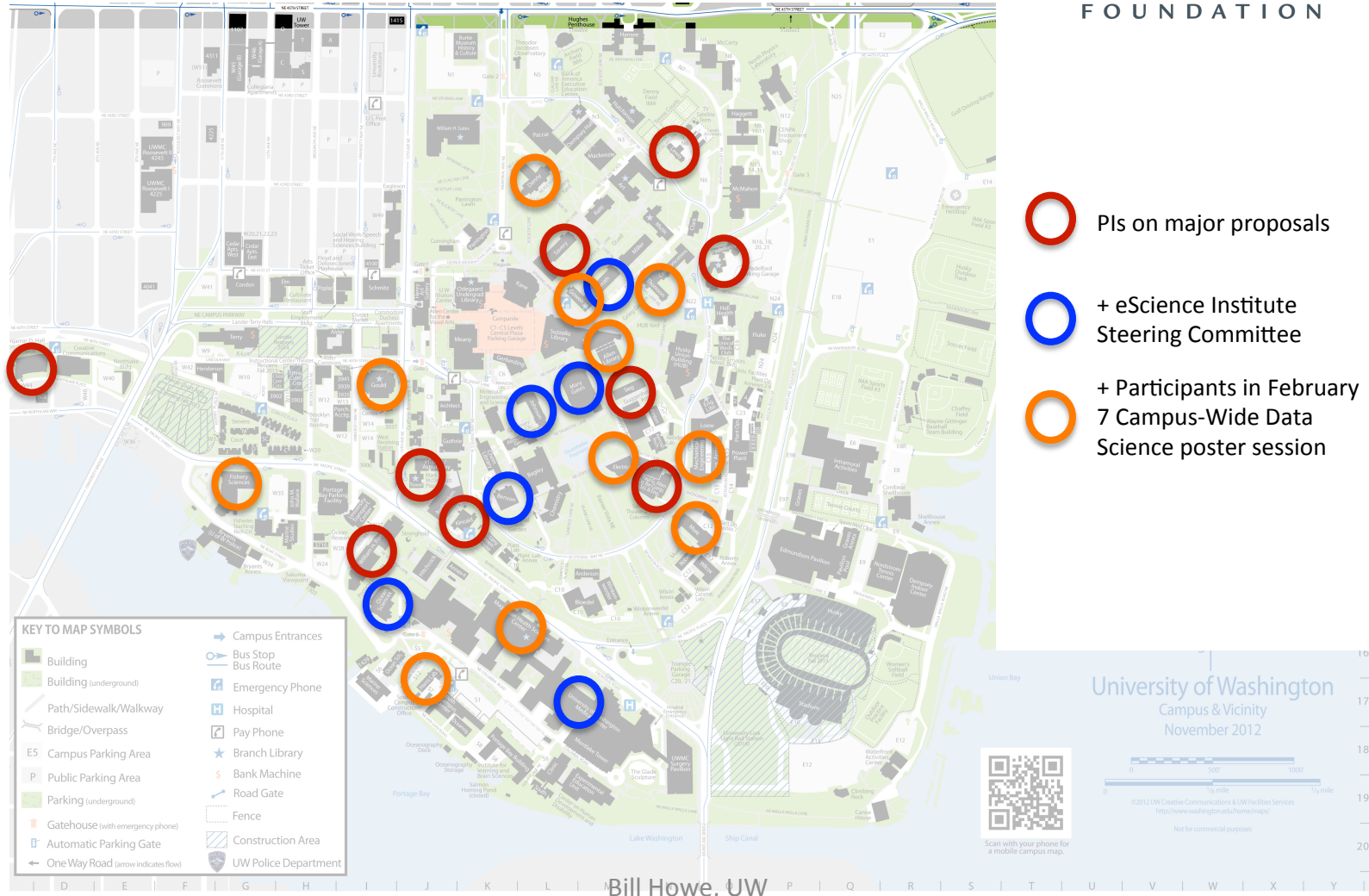


5-year, \$37.8 million cross-institutional collaboration to create a *data science environment*



ALFRED P. SLOAN
FOUNDATION

GORDON AND BETTY
MOORE
FOUNDATION





Data Science Studio

6th floor Physics Astronomy Building

Partnership among:

- Provost
- UW Libraries
- Physics, Astronomy, Arts & Sciences
- eScience Institute



Revised the library focus on working spaces and culture



PHYSICS-ASTRONOMY READING ROOM TRANSITION TO DATA SCIENCE STUDIO

TIMELINE

June 13, 2014, 5pm

The Physics-Astronomy Reading Room (the entire 6th floor of the Physics/Astronomy Building) and book drop will close permanently.

Access will be limited to construction activities beginning at 5pm.

beginning June 16 through Summer 2014

Library materials will be moved out of the space and relocated to other libraries.*
 Interior construction and remodeling will take place. Access for construction only.

Bill Howe, UW

LIBRARIES RESOURCES AFTER JUNE 13, 2014

BOOKS

- The books from the Physics-Astronomy Reading Room will be moved to the Suzallo and Allen Libraries main (open) stacks and shelved with the astronomy and physics research books already located there.
- A few dozen books will be moving to the Mathematics Research Library (Paddelford Hall), Engineering Library, and Odegaard Undergraduate Library, where the content of those books will be closely associated with the subject coverage and clientele of those units.

JOURNALS:

- Journals available online will be moved to the Libraries Auxiliary Storage, and can be requested through the Library catalog
- Journals unavailable online will be moved to the Suzallo and Allen Library main stacks, and will be shelved alongside related material

REFERENCE

- The core Reference collection, i.e., encyclopedias, dictionaries, handbooks, and directories will be moved to the Suzallo Library first floor reference collection. This open browsing area near the elevators has room to read and study, and is conveniently located near several bookscan stations.
- The remainder of the Reference collection will move to the Suzallo and Allen Library main stacks and will be available for checkout
- A few basic and introductory encyclopedias will be moving to Odegaard Undergraduate Library.

COURSE RESERVES

- Odegaard will be the default location for physics and astronomy course reserves after June 13, 2014.
- Starting Summer Quarter 2014, instructors will have the option to have their graduate course reserves at Odegaard Undergraduate Library, Built Environments Library, Health Sciences Library, Engineering Library, or Mathematics Research Library
- Additionally, instructors may choose to have their course materials in the Suzallo Library first floor reference collection, which cannot be checked out, but would be available for browsing or scanning in the library.

SCIENCE FICTION BROWSING COLLECTION

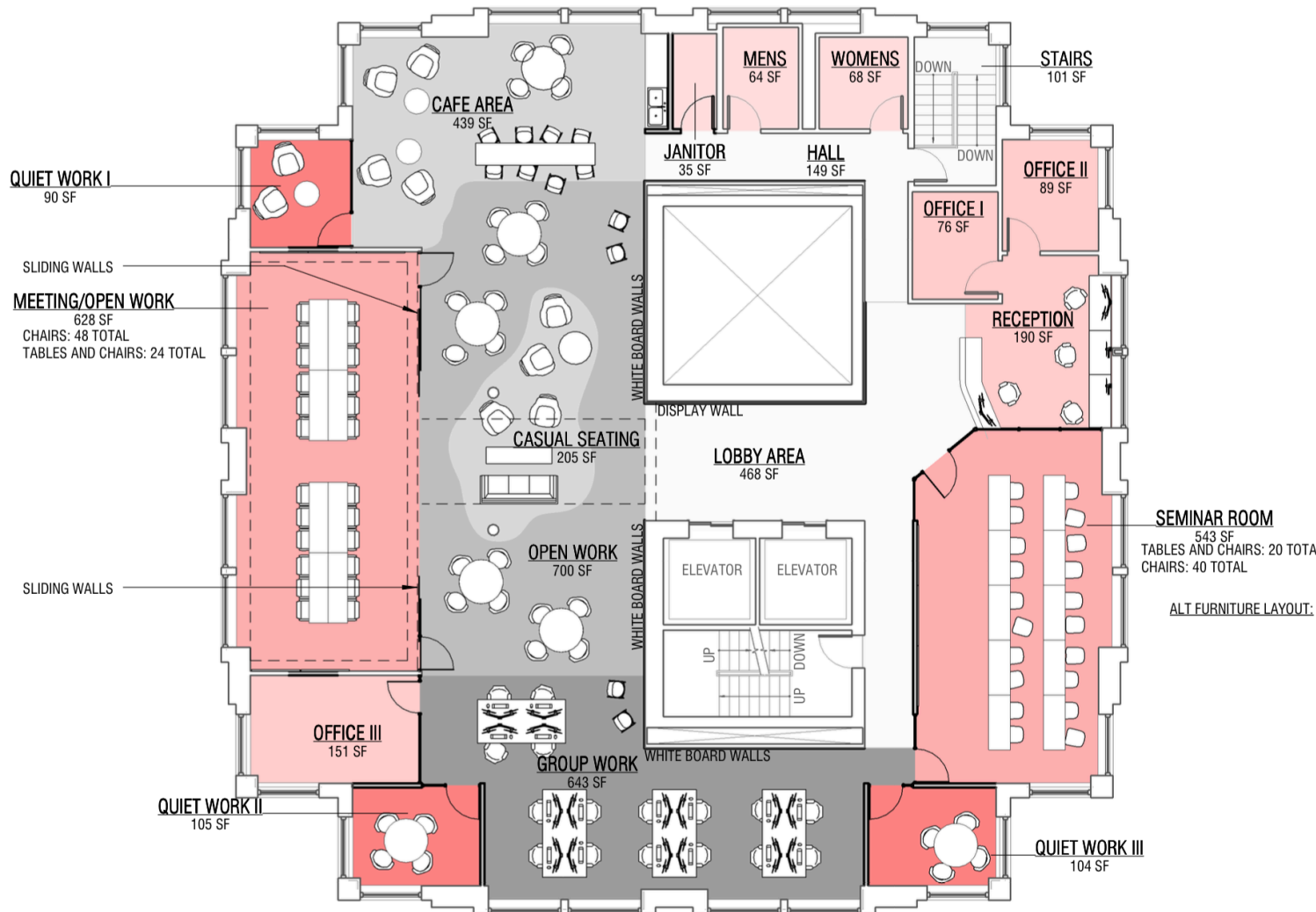
- The browsing collection of science fiction books from the former Chemistry Library will be moving to Odegaard Undergraduate Library and the Suzallo and Allen Libraries.

HOLDS AND REQUESTS

**Casual &
Open work**

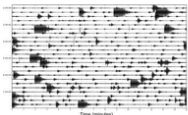
**Seminar &
Group work**

Quiet work





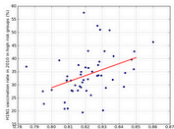
Spring 2014 Incubation Projects



Automated Detection and Analysis of Repeating Earthquakes

Alicia Hotovec-Ellis, Kate Allstadt, Jon Connolly, and John Vidale — Earth and Space Sciences

eScience Contact: Jake Vanderplas



Using social media data to identify geographic clustering of anti-vaccination sentiments

Benjamin Brooks, Abraham Flaxman — Institute for Health Metrics and Evaluation

eScience Contact: Andrew Whitaker



Analysis of Kenya's Routine Health Information

Gregoire Lurton, Abraham Flaxman, Emmanuela Gakidou — Institute for Health Metrics and Evaluation

eScience Contact: Daniel Halperin

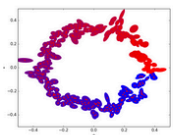
Project leads must physically co-locate with the incubator staff.



Efficient Computation on Large Spatiotemporal Network Data

Ian Kelley, Josh Blumenstock — Information School

eScience Contact: Andrew Whitaker



Scalable Manifold Learning for Large Astronomical Survey Data

Marina Meilă — Statistics

eScience Contact: Jake Vanderplas



ASPASIA: Adult Service Providers and Some Incidental Addenda

Sam Henly — Economics

eScience Contact: Andrew Whitaker

Resident data science team

- Permanent staff of ~5 *data scientists* – applied research and development
- Drop-in open workspace
- Studio “Office Hours”
- Incubation Program

“Don’t see how you do it
without the library.”

...plus seminars, sponsored lunches, workshops, bootcamps...



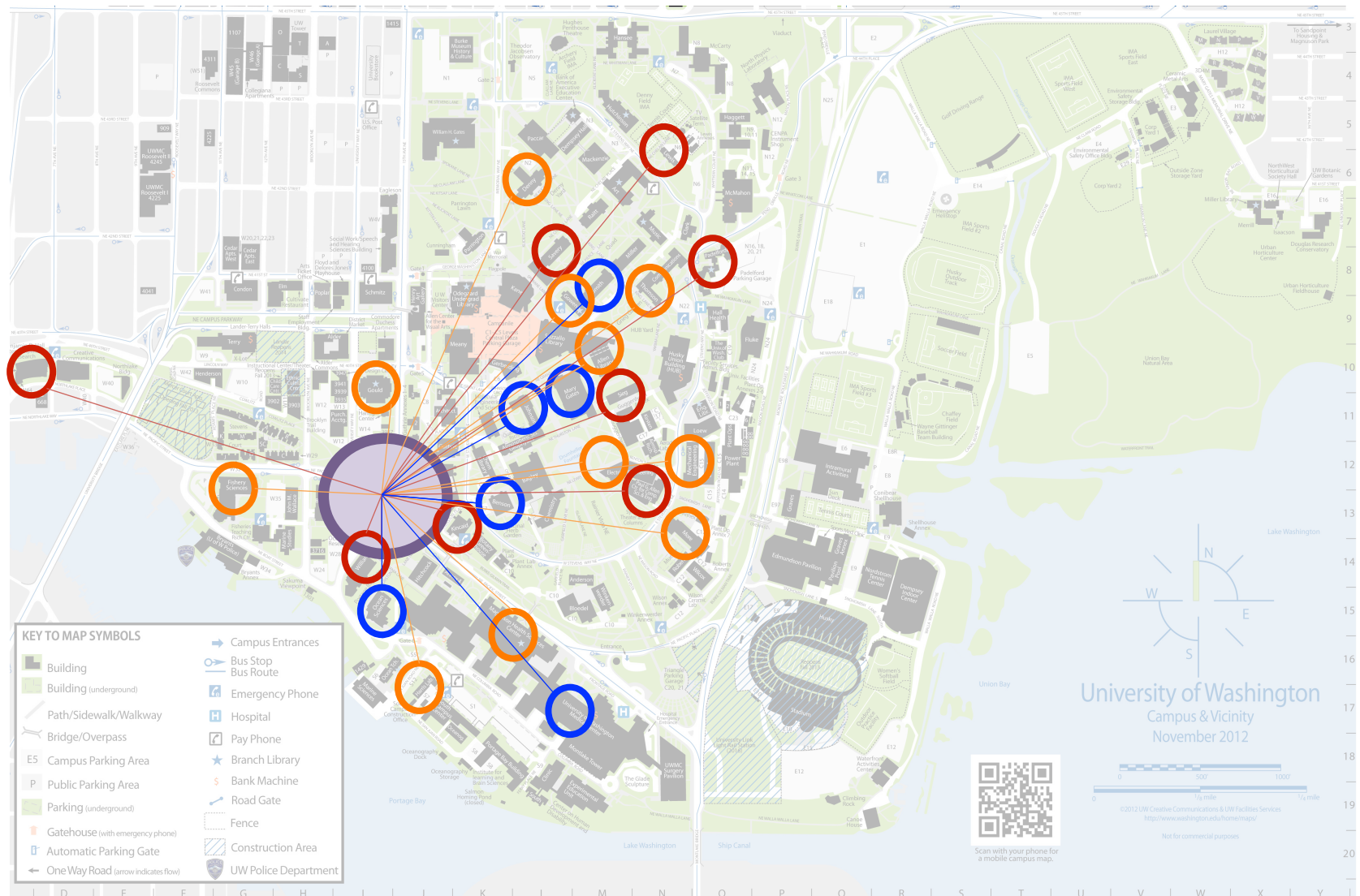
Library in the mix:

Office hours

Reproducibility and Open Science Group

Site visits

Data Science: the rising tide that lifts all boats



Problems - eScience vs. open, curated data

*How much time do you spend “handling data”
as opposed to “doing science”?*

Mode answer: 90%

(Bill Howe, 2015)

What qualifies as releasable data?

Open data constrained by evidential cultures -

Individualism vs. Collectivism (Collins, 1998)

Who takes responsibility for validity and meaning?

Why do we invest in data?

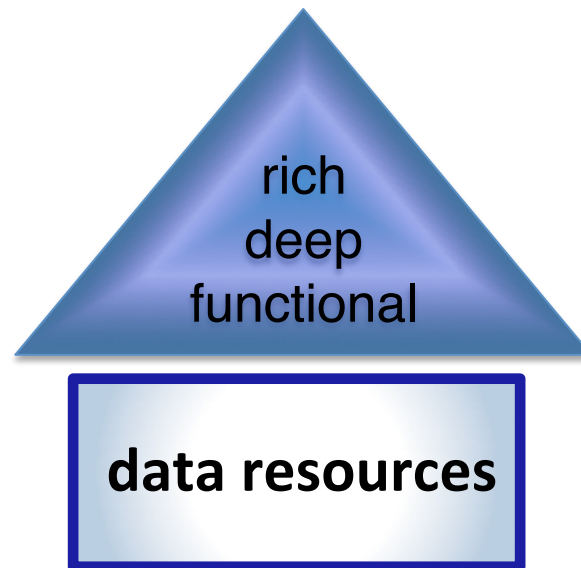
- Open data requirements and expectations
- Reproducibility, replication, and other “Rs”
- Stewards of the common good / scholarly record
- Competitive, innovative research
 - exemplars of “open” research
 - **centers of excellence, research prominence**

Optimizing data for reuse

Different objectives and expertise than:

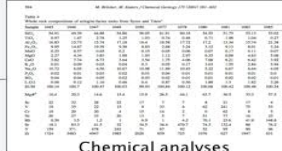
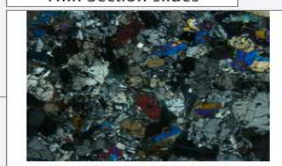
- preserving a record of research
- providing access and transparency

and much more resource intensive.

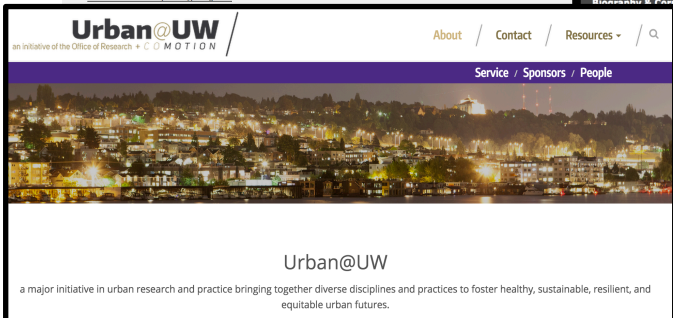
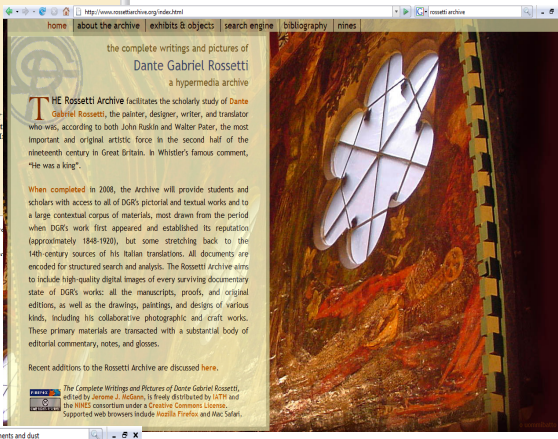
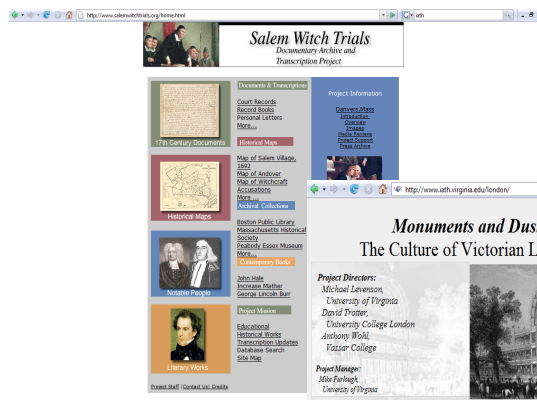


Promoting our deep, rich, functional data

Data Type



Metadata records



Vision

Using state-of-the-art research and technologies in a broad range of fields, from engineering to the humanities we will dramatically improve the efficiency, sustainability, resilience, social justice and equity of cities with a focus on Seattle but an impact on the world.

Empirically derived reuse principles

Data Curation
Profiles Project

Data Conservancy

Site-Based Data
Curation at YNP

- *Releasable \neq reusable*
- *Producer sets /
consumer subsets*
- *Indicators of reuse value*
- *Primacy of method*



Information School
UNIVERSITY of WASHINGTON

True reusability for site-based data



Retain value and promote reuse of data from scientifically significant sites.



Geobiology data from Yellowstone National Park

Reuse dependent on
Sampling procedures

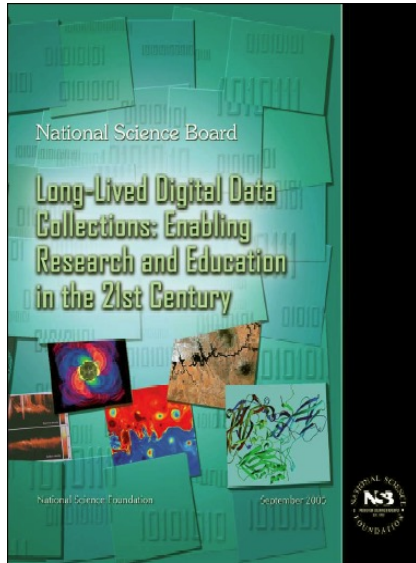
within

Field campaign context:
geological feature
***** new measurements**
vent location, etc.

Used with permission from B. Fouke



Crisis in resource collections

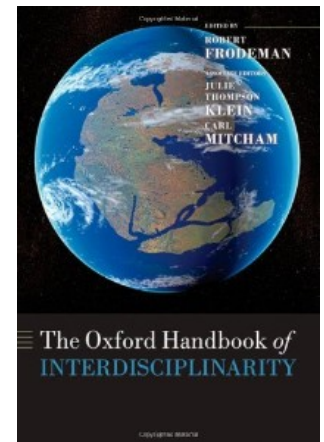


NSB 2005:...ever increasing investment in creating and maintaining collections, and the rapid multiplication of collections, with a potential for decades of curation.

Atkins and Unsworth: Value-added ... widely shared ... collections...enabling ...interdisciplinary research ...

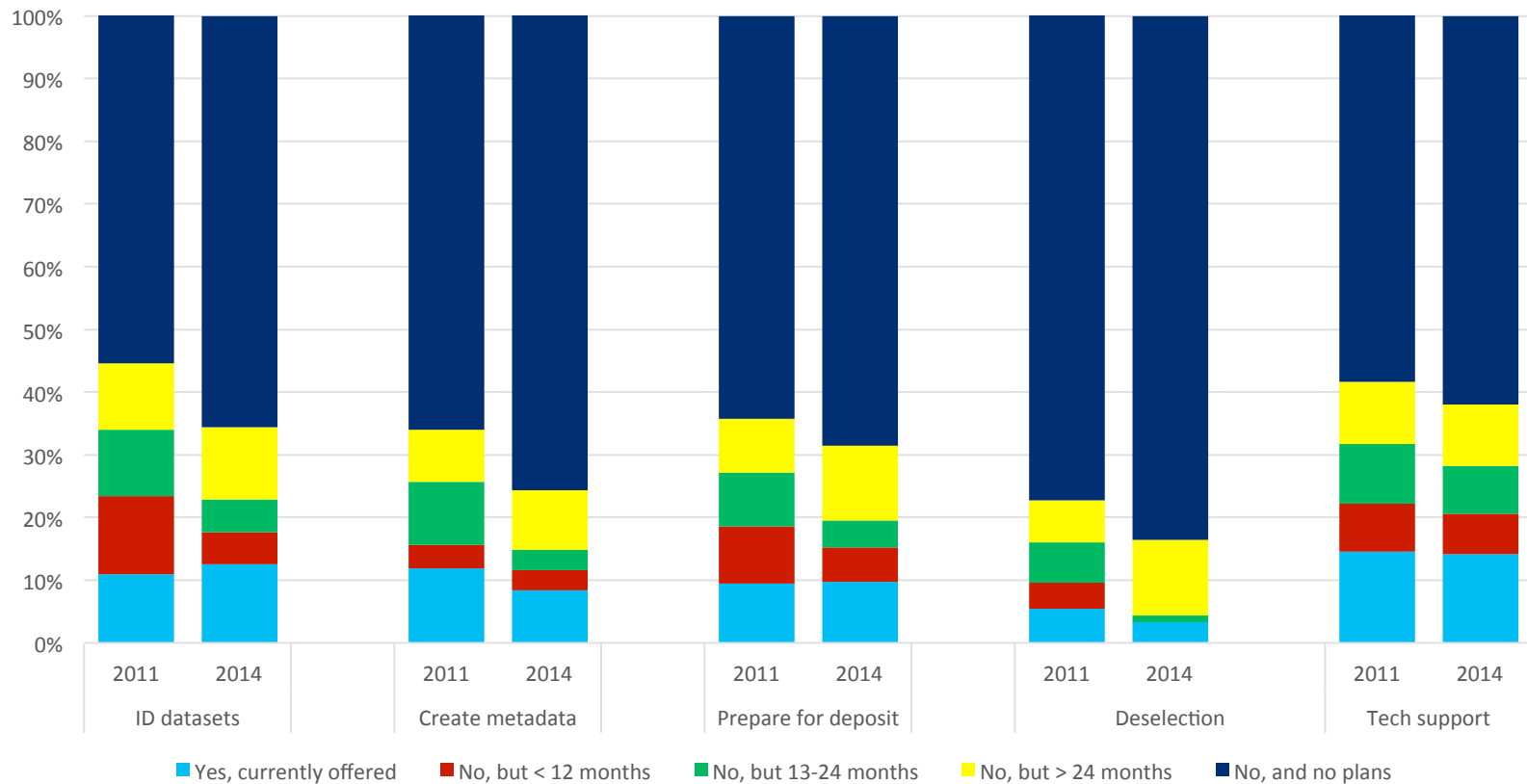
Greatest challenges not ability to move across disciplinary boundaries but in maintaining the increasingly long and mutable intellectual paths to our disciplinary past.

(Palmer, 2010)



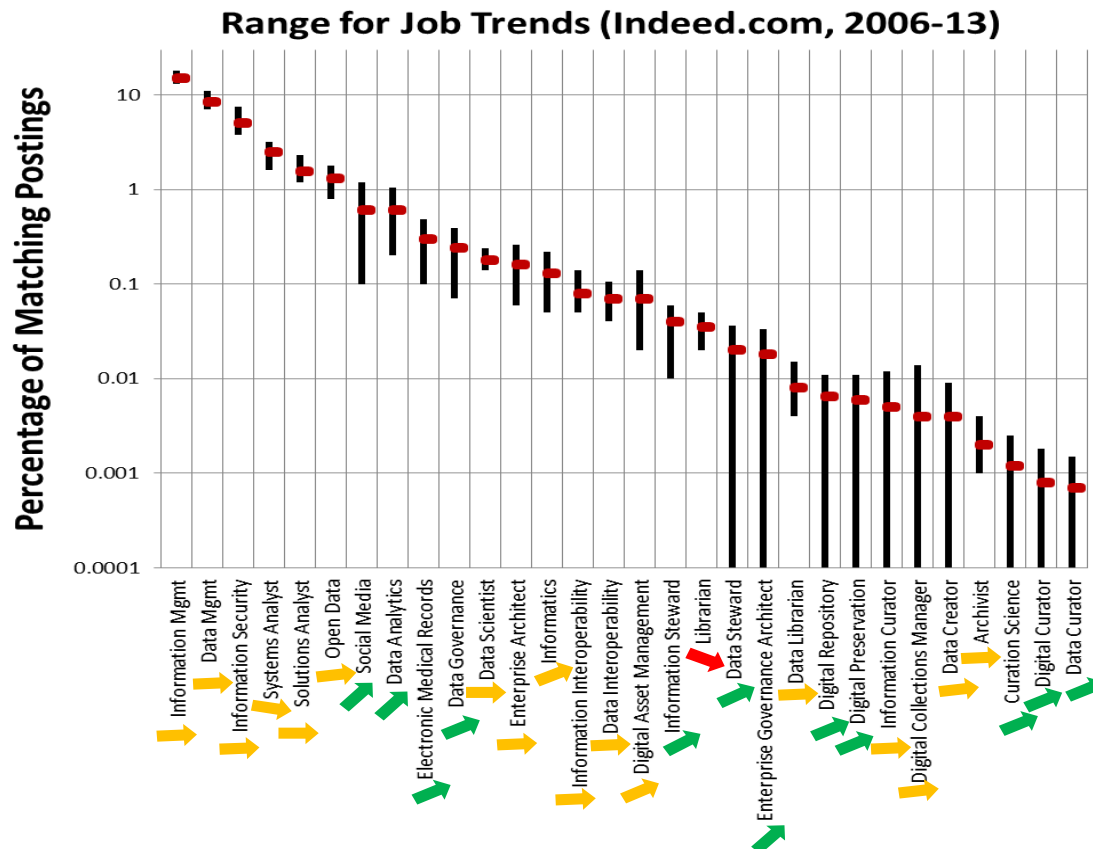
Where are we with the workforce?

Research Data Services Offered or Planned in ACRL Libraries



(Tenopir, ASIST, 2014)

Forthcoming - Preparing the Workforce for Digital Curation



Trending up:

- Information Steward
- Data Steward
- Digital repository
- Digital preservation
- Curation Science
- Digital Curator
- Data Curator

Trending down:

Librarian

(R. Larsen, IDCC, 2014)

Illinois data curation placements

Academic

- 40% of placements,
 $\frac{1}{4}$ of those outside library
- Many focused on metadata and technology

Positions that (probably) didn't exist 5 years ago

- Research Data Management Service Design Analyst
- Data Management Consultant
- Data Science & Informatics Librarian
- Data Curator
- Assistant Dean, Digital Humanities Research

Non-academic positions

- Data Steward Consultant
- Solutions Analyst
- Senior General Engineer
- GIS Specialist
- Director of Archive Technology
- Digital Asset Manager
- Information Architect
- Information Systems Associate
- Digital Project Coordinator
- Media Content Specialist

Classroom experiences with multiple experts

- Earth science data center services
- Cyberinfrastructure R & D
- International data sharing coordination
- Funding & policy perspectives



Field experiences with multiple mentors

Data / Science / Peer mentors



NCAR internships

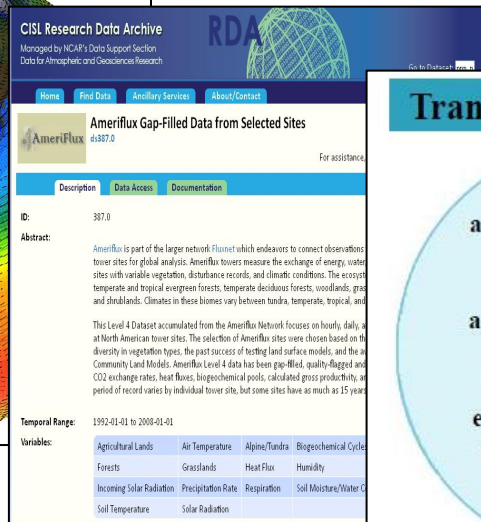
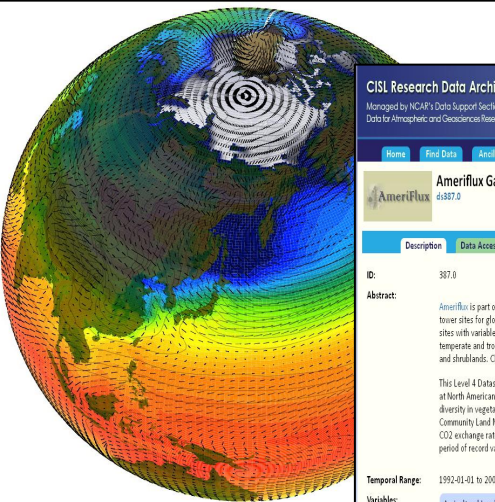
Climate model
metadata

Sensor data
archiving

Social science
data
organization

Time-series
temporal spatial

Analog data for
digital access



Translator and Facilitator

Understands and articulates the needs and goals of scientist

Understands and articulates data manager needs for curation

Creates guidelines to enhance communication and efficiency between scientists and data managers



metadata
harvesting,
standards
compliance,
quality

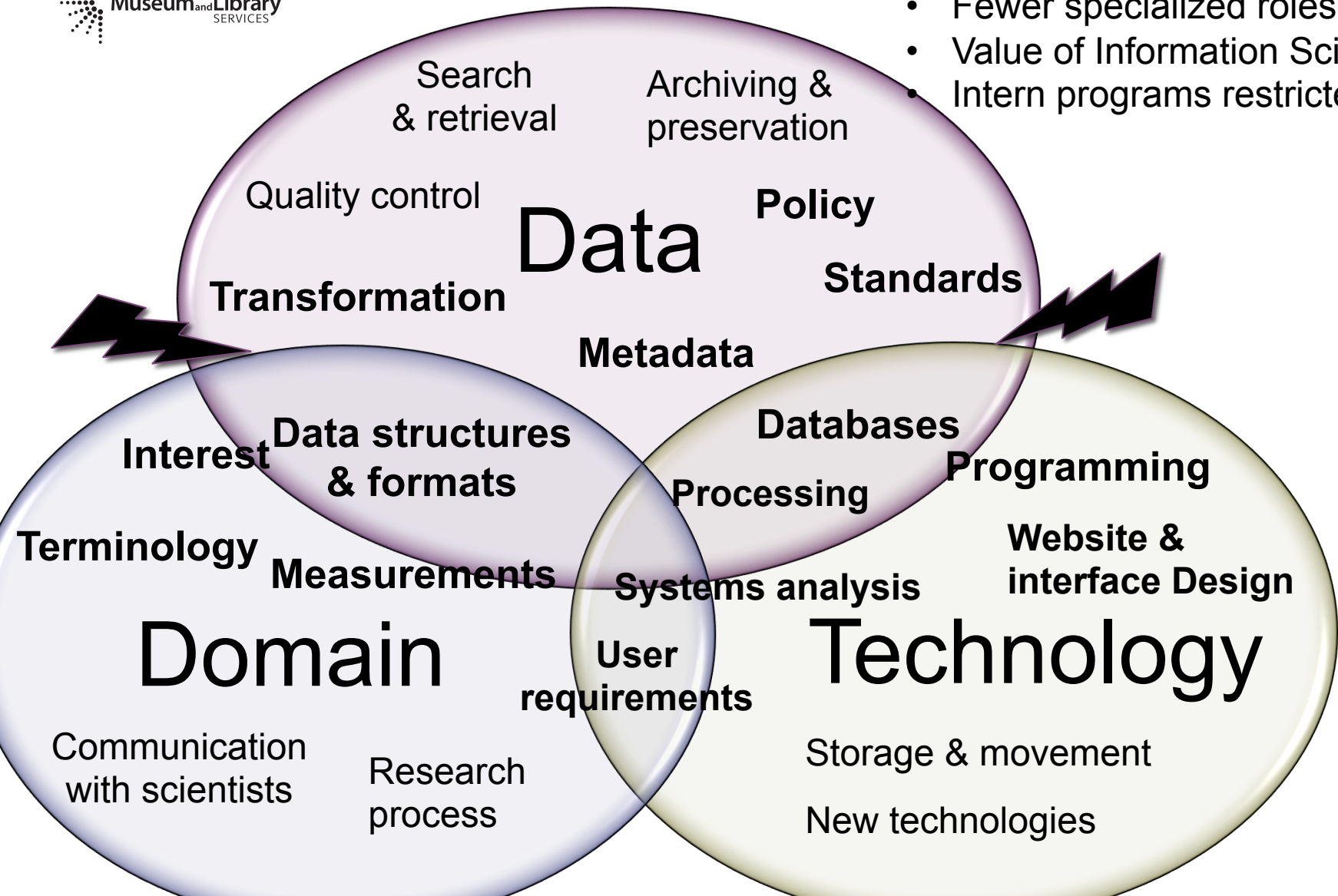
processing &
file migration

cross-disciplinary
data curation;
subsetting

high resolution,
provenance,
NetCDF

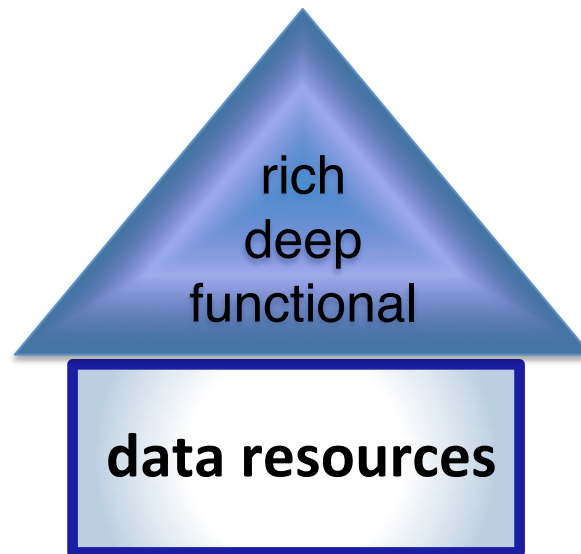
50 international
collections,
OAIS, DOIs

- Some growth in positions
- Fewer specialized roles
- Value of Information Science
- Intern programs restricted



Too much to lose, if we don't get it right.

“Your analytics are only as good as your curation.”



- marshal our strengths in LIS
- leverage progress across disciplines
- build a new LIS foundation in the science of data

Thank you for your attention.

