

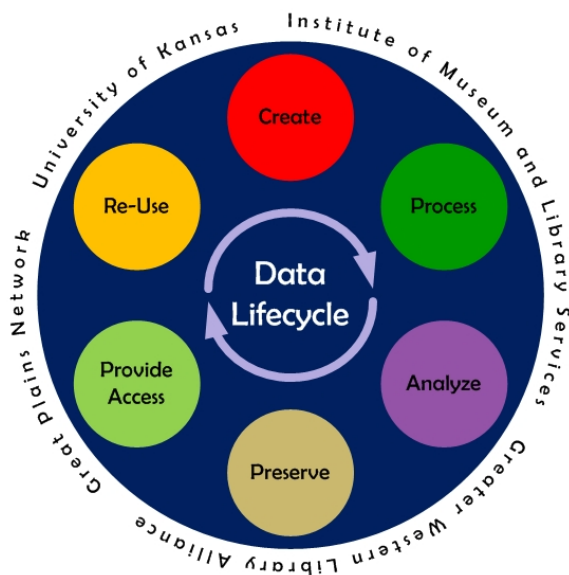
PRELIMINARY REPORT

PLANNING FOR THE LIFECYCLE MANAGEMENT AND LONG-TERM PRESERVATION OF RESEARCH DATA: A FEDERATED APPROACH

Phase 1, Environmental Scan

June 18, 2013

IMLS 51-12-0695



UNIVERSITY OF KANSAS LIBRARIES &
INFORMATION TECHNOLOGY

GREAT PLAINS NETWORK

GREATER WESTERN LIBRARY
ALLIANCE

INSTITUTE OF MUSEUM AND LIBRARY
SERVICES

*PRELIMINARY REPORT
PLANNING FOR THE LIFECYCLE MANAGEMENT AND LONG-TERM PRESERVATION OF
RESEARCH DATA: A FEDERATED APPROACH*

A report to the Advisory Council, IMLS Grant 51-12-0695

June 18, 2013 Version 1.1

Grant Staff

Deborah M. Ludwig, Principal Investigator
Scott R. McEathron, Investigator
Bob Lim, Investigator
Paul K. Farran
Joni M. Blake, Investigator
Gregory E. Monaco, Investigator
Nicole A. Potter, Project Coordinator

Acknowledgements: We wish to acknowledge the efforts and contributions to this report from many individuals, including the members of the committees and research teams listed in Appendix A. We also thank Lars Hagelin of the Greater Western Library Alliance for the development of graphics and technical advice.



This project is made possible by a grant from the U.S. Institute of Museum and Library Services

Views, analyses, and recommendations expressed by the authors of the environmental scan reports reflect the perceptions and opinions of the individual authors and may not necessarily reflect those of the University of Kansas, the Greater Western Library Alliance, the Great Plains Network, or the Institute of Museum and Library Services.

FORWARD

This report provides the preliminary findings of a collaborative planning project undertaken to identify and consider opportunities for a shared approach to management of research data by members of the Great Plains Network, the Greater Western Library Alliance, and the University of Kansas.

As such, our efforts build upon a number of current and recent preparatory activities among institutional members of GPN and GWLA connecting their distinct missions to better serve the research community. The final goal of our work is the development of a plan for a collaborative and federated approach to research data management services that seeks to leverage collective strengths of member institutions

This work is conducted under IMLS National Leadership Planning grant 51-12-0695 with three specific goals.

Goal #1: Undertake an in-depth environmental scan focused on current national and international data management initiatives and on the needs of our member universities for research data management services and infrastructure.

Goal #2: Bring together a GPN and GWLA member forum and two-day workshop for the university research, library, and technology communities focused on understanding challenges and solutions in managing, sharing, and preserving research data.

Goal #3: Create and disseminate a plan for a scalable multi-institutional approach to research data management to support the university members of GPN and GWLA and advance this plan for funding.

This report is a living report and will be updated as each goal is complete and shared with the advisory council and others connected to this work.

TABLE OF CONTENTS

FORWARD	3
EXECUTIVE SUMMARY (MAY 17, 2013)	6
<i>Overview of the Grant and Activity to Date</i>	6
<i>Advisory Council Meeting in Kansas City, May 29 and 30, 2013</i>	9
CHAPTER I: ENVIRONMENTAL SCAN	11
OVERVIEW OF THE ENVIRONMENTAL SCAN	11
SECTION A. REVIEW OF CURRENT LITERATURE	12
<i>Introduction to the Literature Review (Scott R. McEathron)</i>	12
1. <i>Introduction and general works (Scott R. McEathron)</i>	12
2. <i>Assessments of researcher behavior, attitudes and needs (Stephanie Wright)</i>	13
3. <i>Services, roles and responsibilities (Brian Westra)</i>	15
4. <i>Sharing, reuse, publication and citation (Scott R. McEathron)</i>	16
5. <i>Data management planning (Brian Westra)</i>	18
6. <i>Policies and standards [report pending]</i>	19
7. <i>Institutional repositories, approaches, and issues (Sarah Potvin)</i>	19
8. <i>Disciplinary or subject repositories, approaches, and issues [report pending]</i>	23
9. <i>Federated and collaborative approaches (Deborah M. Ludwig & Michael Bolton)</i>	23
10. <i>Economics/costs of data curation [report pending]</i>	34
11. <i>Archiving and preservation (Amalia Monroe-Gulick)</i>	34
12. <i>Metadata and description (Andrew Johnson)</i>	36
SECTION B. RESULTS OF SURVEY (PRELIMINARY AS OF 5/15/2013).....	40
SECTION C. REVIEW OF KEY PROJECTS AND TECHNOLOGIES.....	54
1. <i>Curation of Data</i>	55
2. <i>Preservation of Digital Materials</i>	77
3. <i>Archiving & Repository Services</i>	93
4. <i>Storage Systems</i>	108
5. <i>Enabling Technologies, Services, and Components for Data Management</i>	115
CHAPTER II: OUTCOMES OF DATA MEETING (NOT YET AVAILABLE)	150
CHAPTER III: PLAN FOR DATA MANAGEMENT SUPPORT (NOT YET AVAILABLE)	151
APPENDIX A: COMMITTEES AND STAFF	152
STEERING COMMITTEE	152
ADVISORY COMMITTEE	152
GRANT STAFF.....	153
RESEARCH TEAM A, LITERATURE REVIEW	153
RESEARCH TEAM B, KEY PROJECTS & TECHNOLOGIES REVIEW	153
APPENDIX B: GRANT ABSTRACT	155

APPENDIX C: MAP OF GPN AND GWLA MEMBERSHIP	157
APPENDIX D: BIBLIOGRAPHY (ALPHABETICAL)	158
APPENDIX E: BIBLIOGRAPHY (THEMATIC)	168
1. <i>Introduction and general works</i>	168
2. <i>Assessments of researcher behavior, attitudes and needs</i>	169
3. <i>Services, roles and responsibilities</i>	170
4. <i>Sharing, reuse, publication and citation</i>	172
5. <i>Data management planning</i>	173
6. <i>Policies and standards [pending]</i>	174
7. <i>Institutional repositories, approaches, and issues</i>	174
8. <i>Disciplinary or subject repositories, approaches, and issues</i>	175
9. <i>Federated approaches</i>	175
10. <i>Economics/costs of data curation</i>	177
11. <i>Archiving and preservation</i>	177
12. <i>Metadata and description</i>	178
APPENDIX F: FORM FOR REVIEW OF KEY PROJECTS & TECHNOLOGIES	180
APPENDIX G: SURVEY QUESTIONS	181
APPENDIX H: GLOSSARY OF TERMS AND CONCEPTS	192

EXECUTIVE SUMMARY (MAY 17, 2013)

The level of investment required to support computationally intensive research is large and growing. It makes no sense to replicate resources, skills, and services to all colleges and universities. Instead, institutions have an opportunity to establish explicit, long-term agreements to work with one another for mutual benefit. There will be serious challenges to overcome ... but these challenges must be met to sustain digital research efficiently and affordably.¹

Scholars, institutions, and libraries are significantly affected by the transformation of scholarly dissemination as noted in numerous works on the digital age of research. Over the last decade, it has become apparent that technology advances the opportunity to share research at earlier stages than previously possible. To some extent journal articles now represent the archived outcome of research. Research data, when shared, offers new opportunities to replicate studies and build upon existing knowledge, accelerating the pace of discovery and reflection – but only if the data can be discovered, interpreted, and re-used.

Libraries have a traditional role as the repository for mankind’s record of discovery and knowledge across all disciplines; but that role is now under stress as institutions embrace the rapid technological and ideological transformation of scholarship in order to be competitive in advancing new knowledge and demonstrating institutional value. Information technology is expected to provide more than the on-ramp to networked resources; the emphasis today is on service and technology integration. Researchers are challenged to supplement critical domain knowledge and research focus with data management planning and practice.

Overview of the Grant and Activity to Date

This grant, undertaken with support from the Institute for Museum and Library Services, seeks to discover and plan shared opportunities for the institutional members of the Greater Western Library Alliance (GWLA) and the Great Plains Network (GPN). Through partnership we will leverage our collective resources, skills, and services to meet the challenges of computationally enhanced research. These two membership

¹ Williford, Christa, and Charles Henry. 2012. “One Culture. Computationally Intensive Research in the Humanities and Social Sciences. A Report on the Experiences of First Respondents to the Digging Into Data Challenge.” *CLIR Publication*. No. 151 (June 2012): 4.

organizations have extensive experience and a number of projects underway that reflect their commitment to supporting the research and teaching missions of their member institutions. The map in Appendix C illustrates the large number of institutions potentially benefited by this partnership.

The Greater Western Library Alliance (GWLA) is a 33-member dynamic and project-oriented consortium of leading research libraries in the Central and Western United States. GWLA's strategic member-driven initiatives are collaborative and innovative, resulting in pragmatic outcomes that create user-focused services and programs that broadly support research, teaching and learning, and outreach while influencing the scholarly communications framework.

The Great Plains Network is a consortium of over 20 universities, primarily research intensive and extensive. Researchers at these institutions participate in projects that require advanced networking and are data- and computing- extensive.² GPN plays a role in connecting campuses and state networks to Internet2 and ESnet, facilitating research collaboration among its members, facilitating learning and implementation of emerging technologies through its professional development program and annual meeting (in partnership with GWLA), and engaging members in strategic partnerships that advance their missions.

Our work in the first phase of this grant, the environmental scan, is found in Chapter 1 and highlights a number of crucial areas to be considered and addressed in the second and third phases of the grant.³ Those areas include:

- emerging national policies and practices with the accompanying need to develop policy and practice at the institutional level;
- the complexity of research data and the equally complex challenges we seek to address;
- the heterogeneous approach to institutional and disciplinary repositories and associated approaches to data archiving and/or lifecycle management;
- researcher's attitudes toward stewardship and data sharing;
- and, ultimately, the need to examine new roles for our institutions and to build new communities of service and practice.

² <http://www.greatplains.net/display/Home/Data+Intensive+Projects+Across+the+GPN+Region>

³ In phase two, we bring together our advisory council and in phase three we complete our planning process.

National and international policies, and the practices driven by emerging policy, are being shaped through government, private and public interest in data. In the United States, these interests are reflected in the recent memorandum from the Office of Science and Technology Policy (OSTP) directing federal agencies “with over \$100 million in annual conduct of research and development expenditures to develop a plan to support increased public access to the results of research” (Office of Science and Technology Policy 2013). The literature review portion of the environmental scan takes a closer look at areas of policy as they relate to institutional services, roles, and responsibilities along with the sharing, reuse, publication, and citation of data. Our preliminary look at the survey data collected in the environmental scan indicates that 12 of 17 GWLA and GPN institutions that responded to questions about data policies have created general data “ownership” policies. Fewer respondents have policies related to data for externally funded research, for data that is not supported by grants and other external funding sources, or for managing data after grant funding has expired. Institutional policy development may represent a key area for further exploration.

The complexity of research data and the equally complex challenges we seek to address are also evident in our scan. Interest in the digital humanities and the complexity of texts, images, and media as unstructured data is noted in the literature reviews. A recent report by the Council on Library Resources (CLIR) suggests that we must “Expand our concept of research [and] our concept of research data” (Council on Library Resources, 2012).

The complexity of data and the equally complex approaches to data repositories and data storage is also evident throughout the environmental scan. The landscape includes a mix of disciplinary and institutional repositories, differing content policies for institutional repositories that favor either faculty scholarship or a wider body of materials, and high performance computing initiatives. These varying approaches and levels of stewardship for different types of research data are reflected throughout the literature review, the reports on key projects and technologies, and in the survey findings.

Researchers’ attitudes toward stewardship and sharing research data have been the subject of a number of recent studies also seen in the literature reviews. These studies provide somewhat inconsistent evidence about how data sharing is perceived by researchers. The University of North Carolina advances the belief that “disciplinary communities (including funding agencies for those communities) are best prepared to define what constitutes ‘research data’ for their respect fields” (Marchionini et al. 2012). *The Data Curation Profiles Directory* from Purdue University (Carlson and Brandt 2013) gives valuable insight into the sharing practices found within disciplines and in multi-disciplinary research projects.

Finally, the need is evident to examine new roles and new services for our institutions, our people, and our partnerships. The reviews of key projects and technologies offers several exemplars of curation and archiving that entail partnership, including *DataONE*, the *Data Conservancy*, and the *Texas Advanced Computing Center*. Several of those exemplars are also covered in more detail in the literature review.

We acknowledge that there are significant challenges as we seek to discover and develop a successful partnership to better manage the research data of our member institutions; but we also note with encouragement the numerous present-day success stories of institutions and partnerships stepping forward to meet these same challenges and to build rich communities around research data.

Advisory Council Meeting in Kansas City, May 29 and 30, 2013

We look forward to beginning our work with the advisory council to develop a common vision and plan for success. Please give some thought to the following questions as starting points for our discussion in May as we consider common ground for our planning efforts.

Question 1. How is your institution responding to existing funder, government or disciplinary mandates or initiatives to share and steward research data? What policies do you have in place or do you anticipate creating?

Question 2. To what extent is research data management an administrative priority for your institution? What are your most pressing needs: places to store and interact with data during active stages of research, places to archive and preserve data for future access, helping researchers write data management plans to satisfy grant requirements, or something else?

Question 3. Most policy efforts target those data and publications produced from scientific and engineering research. What are your institution's needs for solutions to discovery, access and preservation of non-scientific research data?

Question 4. Which divisions (IT, Office of Research, Libraries) does your institution expect to play what specific roles to ensure compliance with existing data policy or efforts to preserve data as an asset for scholarly communication? How are, or how will, multiple divisions work together, if required?

Question 5. With respect to discovery, access, and sharing of research data, what do you already have in place, and what are your institution's greatest needs for new research infrastructure, resource, and expertise?

Question 6. What solutions for research data management do you believe could be developed through a GWLA / GPN partnership that would meet some of the priorities for your institution?

Question 7. How could a multi-institutional model create value for your institution and what would the value be?

Question 8. What do you recommend as specific next directions in our planning as partners?



CHAPTER I: ENVIRONMENTAL SCAN

Overview of the Environmental Scan

This purpose of the environmental scan has been to prepare our project working groups with knowledge and to share that knowledge with members of our steering and advisory councils to inform the planning process.

The specific outcomes are:

1. *Comprehensive literature review* to identify existing studies, formally published literature, and reports developed by major agencies and organizations, (e.g., the December 2011 pre-publication report of the National Science Board, “Digital Research Data Sharing and Management”).
2. *Survey* designed and conducted of researchers at participating GWLA / GPN member institutions to validate current needs.
3. Identification of *current studies and projects* relevant to the conduct of faculty research, researcher attitudes, and lifecycle management of research data;
4. *Site visits* and/or phone visits with leading institutions that provide relevant existing and emerging management models. ⁴

⁴ Site visits will be conducted after the initial meetings with the advisory council.

Section A. Review of Current Literature

Introduction to the Literature Review (Scott R. McEathron)

In the process of reviewing the literature on the theme of lifecycle management and long-term preservation of research data, it became clear that many related sub-themes have emerged and continue to evolve. The literature related to this theme is like a growing tree with new branches forming and growing out from one or two main trunks. We have tried to identify the relevant branches and have grouped the literature into the following twelve topics:

1. Introduction and general works
2. Assessments of researcher behavior, attitudes and needs
3. Services, roles and responsibilities
4. Sharing, reuse, publication and citation
5. Data management planning
6. Policies and standards
7. Institutional repositories, approaches, and issues
8. Disciplinary or subject repositories, approaches, and issues
9. Federated approaches
10. Economics/costs of data curation
11. Archiving and preservation
12. Metadata and description

We have limited the literature review to works that have had a great impact (high citation rate), are excellent case studies, or show promise as examples of innovation. Thus the current literature review is selective and limited to English language works, most of which from the last 10 years. Some works could easily fit into more than one area. We have left it to the reviewer of any given section to decide which best fits that section. The selected works have been arranged thematically in APPENDIX E. An alphabetical bibliography of these works is located in APPENDIX D.

1. Introduction and general works (Scott R. McEathron)

As we began to review the literature relevant to lifecycle management and long-term preservation of research data, a number of bibliographies and guides surfaced. We are indebted to Charles Bailey's (2013) *Research Data Curation Bibliography*. With over 200 citations and growing, it is probably the most comprehensive bibliography on the subject. However, other distinctive bibliographies and guides are noteworthy. The Westra, et al. (2010) bibliography of *Selected Internet Resources on Digital Research Data Curation* presents a thematically organized bibliography of the more important internet

based resources. Witt and Giarlo (2012) provide a description of another unique guide, *Databib: An Online Bibliography of Research Data Repositories*. *Databib* currently provides records on over 500 repositories worldwide and is an example of the growth and geographical breath in digital data repository services.

Some of the more important early works may be described as “calls to action” in response to the growth in e-sciences (Gray et al. 2002; Hey and Trefethen 2003; Hey et al. 2009). Grey (2002; Hey et al. 2009) called for tools to support the whole research cycle--and specifically the curation, archiving, and publishing of digital data. Hey and Trefethen (2003) called for the creation of new types of digital libraries to archive and curate e-science data and provide other data-specific services.

Other articles have made similar “calls to action” beyond the e-sciences (Borgman 2009; Ogburn 2010). Borgman’s article is indicative of the growing interest in digital humanities. She makes comparisons of the data practices between the sciences and humanities, which she then uses to frame a series of lessons and questions. Ogburn (2009) makes a similar “call to action” in her article focused on the potential role of libraries in the area of data curation.

The digital or data curation theme has continued as a central focus for many writers (Higgins 2008; Ogburn 2010; Yakel 2007). Of note, Higgins describes the Digital Curation Centre’s Curation Lifecycle Model as a tool to help plan curation and preservation activities to different levels of granularity (135). Yakel (2007) explores the evolution of digital curation as becoming “an umbrella concept that includes digital preservation, data curation, electronic records management, and digital asset management” (335).

A number of compilations have also emerged. The most noteworthy example is Graham Pryor’s *Managing Research Data* (2012) with covering many of the same topics of this literature review--several of the chapters are specifically cited. What follows is a more focused literature review of each of the themes related to lifecycle management and long-term preservation of research data.

2. Assessments of researcher behavior, attitudes and needs (Stephanie Wright)

Publications assessing researcher needs, behaviors, and attitudes surrounding data management have proliferated over the last few years. While most assessments seem to use the same tools (surveys and/or interviews), they can vary widely in focus. Some focus on the disciplines of the researchers (Williams and Pryor 2009; life sciences), some on a particular phase of the data lifecycle (Feijin 2011; access and storage, Swan and Sheridan 2008; sharing, and Kuipers and van der Hoeven 2009, Sharpe 2006;

preservation), others on specific geographic area (Sharpe 2006, Swan and Sheridan 2008; UK) or particular institutions (Marchionini, et al. 2012; Scaramozzino, Ramirez and McGaughey 2012). Not surprisingly, some publications had more than one of these foci, and there are similar findings in more than one assessment. In addition, there are publications that analyze and synthesize data from multiple assessments - attempting to provide a broader picture of data management from the researcher environment (Feijin 2011) and some included perspectives from other research-data stakeholders (publishers and data managers).

Feijin (2011) falls into the latter category and is in itself a literature study. The report focuses on researcher needs in terms of data storage and access, but covers much more - with the primary conclusion that researchers do want and need support services for managing digital data. The authors provide a list of requirements necessary to make those support services successful, including making sure tools and services are easy to use and are “in tune with researchers’ workflows”. One of the publications reviewed by Feijin (2011) worth noting is the PARSE report focusing on digital preservation (Kuipers and van der Hoeven 2009). This assessment has broad geographic and disciplinary representation as well as responses from publishers and data managers, which provides a useful comparison of needs and motivations across stakeholders.

Swan and Sheridan (2008) interviewed over 100 researchers across multiple disciplines and assessed what researchers are actually doing in regards to data sharing, as well as uncovering their motivations and constraints with sharing data. In a similar vein, Scaramozzino, Ramirez and McGaughey (2012) looked at multiple data curation behaviors of California Polytechnic State University researchers and compared their actions to their expressed beliefs and attitudes.

Within the group of disciplinary studies, Williams and Pryor (2009) used information lab notebooks to supplement the more familiar assessment tools of interviews and focus groups to understand information exchange behaviors (including data sharing) of life sciences researchers within the context of their roles in a research group. This method allowed the authors to create diagrams showing the information flow within each research group and how they move through the stages of the data life cycle.

Finally, the *Data Curation Profiles Directory* out of Purdue University (Carlson and Brandt, 2013) is a publication worth exploring as it maintains a multidisciplinary collection of profiles of specific data set requirements as detailed by the researcher. In essence, the profiles are case studies identifying how researchers across institutions and disciplines deal with data management issues.

3. Services, roles and responsibilities (Brian Westra)

Research Data Services (RDS) and *Research Data Management* (RDM) services are two umbrella phrases authors have used to describe the suite of services related to data curation (Jones, Pryor, & Whyte 2013; Tenopir, Sandusky, Allard, & Birch 2013). These services may be viewed through the lens of organizational structures, degree of investment, competency requirements, or in relation to existing library services. Giarlo depicts the data curation services of academic libraries as “data quality hubs” (Giarlo 2013, 6), where curatorial practices address such factors as trust, authenticity, and usability. He then examines the implications for data curation practices. Others use the collection development metaphor to express data curation services (Choudhury 2010).

Lyon examines the implications of research data informatics on libraries and outlines how libraries can transform to meet these needs. She describes ten different data support services: surveys, planning, informatics, citation, training, licensing, appraisal, storage, access, and impact. She also outlines a framework of roles, responsibilities, requirements, and relationships for providing these services (Lyon 2012).

Lewis outlines nine areas where libraries can be active in relation to research data, ranging from developing library workforce data skills and confidence, to leading or partnering on the development of local data policies, and influencing national policy (Lewis 2010). Others have used a tier model for data management services activities that reflect increasing involvement or “embeddedness” with the researcher and the data, progressing from education to consultation to infrastructure (Reznik-Zellen, Adamick, & McGinty 2012). The data.bris project at the University of Bristol outlined four options in their business case for a pilot data management service: ‘do nothing, do little, preferred, and gold-plated’ (Whyte 2013).

Business cases, roadmaps, and strategic planning documents may be useful tools for defining the structures, partnerships, and organizational development required to provide new services. There are many examples including: (Beitz, Dharmawardena, & Searle 2012; Jones et al. 2013; Macdonald & Martinez-Urbe 2010; University of Edinburgh 2012; Cole 2013; Marchionini 2012; Whyte 2013; Witt 2012).

National and international initiatives such as DPN (Digital Preservation Network) and RDA (Research Data Alliance) aim to provide “an open global research infrastructure”, and avenues for influencing and developing services and collaborations that can have broad impact.

Pilot cases provide an opportunity to explore, develop and apply the infrastructure that is needed to support the full lifecycle of research data. Examples include the University of California San Diego's work with five projects, in neuropsychology, archaeology, oceanography, earth science/topography, and astrophysics (Moore 2013), and the University of Manchester's work with biomedical data (Poschen et al. 2012).

The preceding examples and others exemplify successful collaborations and shared oversight between libraries and other institutional partners. For example, Purdue University Libraries partnered with the research office and information technology to develop and fund the Purdue University Research Repository. At the University of California San Diego, the Research Cyberinfrastructure Oversight Committee and Implementation Team have representatives from academic research departments, computing, the supercomputing center, the office of research, and the libraries (University of California San Diego).

Other kinds of infrastructure collaborations may include the evaluation and implementation of electronic lab notebooks (University of Wisconsin - Madison 2012), resource navigation systems (Haendel, Vasilevsky, & Wirz 2012), and image management systems (Linkert et al. 2010).

Policy development is at the higher level of Lewis' hierarchy (Lewis 2010). Some institutions are developing guiding principles that not only state policies for researchers, but institutional responsibilities as well (Flach & Price 2012). While researchers may define responsibilities for data stewardship based on a data "ownership" conceptual framework (Marchionini 2012), institutional responsibility statements can signal the commitments that the larger organization is making to data management.

4. Sharing, reuse, publication and citation (Scott R. McEathron)

From the White House in Washington, D.C. to an unassuming pub called the Panton Arms in Cambridge, UK, scientists, scholars, attorneys, and others are now talking about the importance of sharing data in research and more specifically, the need for open data. In a recent memorandum from The Office of Science and Technology Policy (OSTP) (2013), John Holden has directed "each Federal agency with over \$100 million in annual conduct of research and development expenditures to develop a plan to support increased public access to the results of research funded by the Federal Government" (2). Similar activity can be found at the grass-roots level. In the UK, a group of scientists and others have developed a set of guidelines called the Panton Principles that aim to make open data "freely available on the public internet permitting any user to download, copy, analyze, re-process, pass them to software or use them for any other purpose

without financial, legal, or technical barriers other than those inseparable from gaining access to the internet itself” (Murray-Rust, et al. 2010). The growing interests of scientists and the current and proposed future requirements of national research funders will continue to make data sharing a central theme within the literature.

Many articles used the theme of lifecycle management and long-term preservation of research data center as one of the primary reasons to preserve research data: **so it may be shared and reused**. The articles vary in their scope--from broad overviews, framing research agendas to multi-scaled or disciplinary specific data sharing and specific publisher policies.

A broad overview of the reasons for sharing (and not sharing) research data and an agenda for future research is provided by Christine Borgman (2012). A similar opinion, from the point-of-view of a scientist, on the need to share data is offered by Vision (2010). Vision concludes that instead of individual publishers or journals as the repository/archive of data, a “superior approach is a disciplinary repository that has data as its primary focus and is shared by a scientific community larger than a single journal or publisher” (330). The model he gives as an example is Dryad⁵. Vision also forwards the opinion that “journals are in the best position to promote the practice of archiving ‘small-science’ data upon publication” (330).

Cragin, et al. (2010) studied the data characteristics and sharing practices of “small-science” research areas and the implications for institutional repositories. They concluded that because of the high level of variation and complexity in data forms and sharing practices,...resource demands for curation services...will be high” (4036). Faniel and Zimmerman (2011) present a very thorough review of current research on data sharing and reuse and offer a research agenda in three areas: 1) how researchers manage their data; 2) questions around the re-users of data; and 3) the influences on how researchers make their data available (65-66).

An example of disciplinary concerns for data publication can be found in the field of Bioinformatics as described by Chavan and Ingwersen (2009). These authors detail specific concerns of data within the discipline such as data not being easily discoverable or accessible and the lack of recognition for publishing it. Further, they offer a “Data Publishing Framework” to incentivize the goal of sharing data.

Another way of conceptualizing the issue of data sharing is from the publishing framework. What does it mean to publish data? Lawrence et al. (2011) provide an

⁵ www.dryad.org

overview of the structures that they feel are needed in order to improve a more formal system of data peer review, publication, and citation. They and others (Simons 2012) also suggest Digital Object Identifiers (DOIs) to provide a permanent identifier and locator for datasets (7). Mooney and Newton (2012) also found that the majority of articles they reviewed lacked adequate citation of data used. They concluded that full citation of data is not a normative behavior in scholarly writing (15).

Piwowar and Chapman (2008) explored the correlation between a journal's data-sharing policies and the likelihood of having accessible datasets. They reviewed the policies for data sharing of journals that publish results of research utilizing gene expression microarray data. They found that "high impact journals tended to have strong data sharing policies" (11) and also "articles published in journals with a strong data-sharing policy are more likely to have publicly available datasets" (15).

5. Data management planning (Brian Westra)

Although the NIH had required since 2003 that researchers address data sharing for grant awards over \$500,000, the 2011 NSF requirement for data management plans (DMPs) had a more significant impact on proposal-writers and raised the importance of the services provided academic institutions. The recent Office of Science and Technology Policy (OSTP) mandate expands the basis for data management plan requirements to all federal agencies providing over \$100 million in research grants (Holdren 2013).

A 2011 survey of ACRL libraries in the US and Canada (Tenopir, Birch, and Allard 2012), highlights the range of services provided by academic libraries. At the time of the ACRL survey, only 20% of library respondents were providing consultations for faculty and graduate students on data management plans (DMPs), but another 22% planned to do so within the next 2 years.

Services specific to data management plans for grant-funded research may include consultations with grant writers, DMP training and workshops, and form-based tools for creating a DMP. Some libraries have begun to review larger sets of DMPs (Parham and Doty 2012).

Libraries are also working alone and in collaboration with other campus partners to develop service models in support of the constituent elements of a plan, such as storage and backup, electronic lab notebook systems for description (University of Wisconsin-Madison 2012), and data preservation and sharing. Some aspects of these services are considered elsewhere in this document.

Understanding researcher needs, and presenting services with measurable positive impacts on those needs are critical to the success of DMP services. The data curation profile provides a framework for determining data management practices and needs of researchers (Carlson 2012). Establishing trust relationships with researchers is important, and embedded librarianship may provide one of several paths toward this end (Carlson and Kneale 2011).

The DMP Online tool developed in the United Kingdom, and its relative, the DMPTool developed by the California Digital Library and partners in the U.S., are employed by some libraries to walk grant-writers through the process of developing a data management plan for submission with a grant proposal. Sallans and Donnelly (2012) compare and contrast these two form-based web resources. The DMPTool links to data plan requirements published by the funding agency units, and local guidance materials can also be incorporated into the web pages.

Cornell and Syracuse librarians investigated funder data policies and assessed their comprehensiveness and level of detail against a rubric across a range of data policy elements (Dietrich et al. 2012). Funder data policies are usually general in scope, though a few programs have more explicit guidance and form-based resources for developing plans and reporting on data activities, such as IEDA (Integrated Earth Data Applications).

Libraries can provide training specifically about data management plans for researchers (Johnston, Lafferty, and Petsan 2012), but the data management plan structure can also be used as a rubric for broader data management training.

6. Policies and standards [report pending]

7. Institutional repositories, approaches, and issues (Sarah Potvin)

What is an institutional repository? Clifford Lynch's early, formative definition urged: "An institutional repository is not simply a fixed set of software and hardware" (2003, 2). Rather, it is "a set of services that a university offers to the members of its community for the management and dissemination of digital materials created by the institution and its community members. It is most essentially an organizational commitment to the stewardship of these digital materials..."(2).

Were these institutional repository services intended to manage data over their lifecycle, and how have they evolved to do so? Lynch's broad definition suggests the possibility; he predicted that "a mature and fully realized institutional repository ... will also house experimental and observational data" from community members (2). Subsequent research and surveying point to institutional willingness to engage with data, though financial, staffing, and technology constraints have intervened. Lynch and Lippincott (2005) report on results of a Coalition for Networked Information survey of relevant individual and consortial member institutions, aimed at assessing the "current state of institutional repositories (IRs) in the US." In preparing the survey, they observed "two views" of IRs: "One characterizes an institutional repository as primarily addressing dissemination of various forms of e-prints for faculty work..."; "The second approach conceives of an institutional repository as broadly housing the documentation of the intellectual work—both research and teaching—of the institution, records of its intellectual and cultural life, and supporting evidence for present and future scholarship." This second variety "will include e-prints, certainly, but also datasets, video, learning objects, software, and other materials."

The 2005 CNI survey results indicated that "a significant number of institutions are committed to institutional repositories that go far beyond e-prints" (Lynch and Lippincott 2005).⁶ While only four responding institutions indicated that their IRs currently held data sets, twenty-six indicated plans to do so over the next 1-3 years—the largest number of responses garnered for any type of planned content. Lynch and Lippincott conclude that institutional repositories in the US "are being positioned decisively as general-purpose infrastructure within the context of changing scholarly practice, within e-research and cyberinfrastructure, and in visions of the university in the digital age" (2005). A 2009 ARL report on repository services echoed the observation that "repositories are developing rather than developed" (Moore, et al. 2009, 8): "Just a few years ago, many libraries were acting on a vision of repositories that focused on preprints and postprints of faculty publications and theses and dissertations. ... We now understand better that institutions produce large and ever-growing quantities of data, images, multimedia works, learning objects, and digital records..." (8).

How have particular institutional repositories engaged with data? Many institutions have formulated data management plans, policies, and resources. Some universities have extended their "set of services" around data management while depending on repository systems beyond the university to disseminate the data.

⁶ Italics removed.

The University of North Carolina's 2012 "Research Data Stewardship at UNC: Recommendations for Scholarly Practice and Leadership," a working report produced in response to a charge from the UNC Provost, prioritizes the placement of research data into subject/disciplinary repositories, rather than a UNC IR. It notes: "Individual researchers are best positioned to identify which repositories are most appropriate to their data and are encouraged to take advantage of public repositories whenever possible" (Marchionini, et al., 18). A central campus data registry is recommended to "provide a single place where anyone with access could get an overview of the data preservation efforts at UNC and at a minimum, a list of individual data sets that have been stored" (21). There is recognition, however, that public disciplinary repositories might not be ideal for all data. The recommendations thus extend to the development of "A repository for UNC research data where public repositories do not exist and/or data must be locally managed by contractual or sensitivity requirements" (24). This approach loosens the university's role in publishing its researchers' data, while ensuring that the data is accessible and published in a repository preferred by the researcher (or funder).

Other institutions have taken a more active role in overseeing the full lifecycle management of institutional research data. Notably, the Purdue University Research Repository (PURR)⁷, built on a HUBzero platform, aims to support researchers as they develop data management plans, collaborate on research, and, ultimately, publish and archive their datasets, which are issued with DOIs (PURR website, 2013). Still others, such as Indiana University Bloomington, have integrated specialized ingest processes, which prompt the collection of data-specific metadata for those contributing datasets to the institutional repository (Konkiel, 2013).

How have other institutions approached the decision by researchers to deposit data in either disciplinary or institutional repositories? Universities retain lists of potential repositories for deposit, often relying on other institutions to develop and supplement these lists (University of Oregon; University of Idaho Library).⁸ They have also

⁷ PURR website. (Accessed April 26, 2013) <https://purr.purdue.edu/>

⁸The University of Oregon site lists repositories by discipline and includes the note: "You may want to use the University of Oregon Libraries' institutional repository, called Scholars' Bank, for your data." University of Oregon (Accessed April 26, 2013) "Data Repositories," in "Research Data Management," <http://library.uoregon.edu/datamanagement/repositories.html>. University of Idaho Library (Accessed April 26, 2013) "Data Management: External Links." http://www.lib.uidaho.edu/services/data/data_management/links.html. The UNC Data Stewardship report recommends that UNC "design a process whereby researchers are both informed of existing resources and encouraged to use existing resources as a first choice. ... it is unlikely the university will have the wherewithal or inclination to produce a complete listing of all data repository resources along with a full enough specification of their intended use-cases.

developed tools to facilitate data curation and, ultimately, repository deposit. Some tools are designed to interoperate with both institutional and disciplinary repositories for deposit. Cornell University's DataStaR (Data Staging Repository), "a platform and a set of services meant to facilitate data sharing" takes a destination-repository-neutral approach (Steinhart 2011, 16). The intermediate repository allows Cornell researchers to "store and share data with selected colleagues, select a repository for data publications, create high quality metadata in the formats required by external repositories and Cornell's institutional repository, and obtain help from data librarians with any of these tasks" (Steinhart, 16). Monash University had previously explored the role of intermediate or "collaboration" repositories that facilitate researchers actively working on their data prior to their appearing in "publication domain" repositories (Treloar, Groenewegen, and Harboe-Ree 2007). This approach, too, strengthens the university's set of offered services while retaining neutrality on publication site.

One issue regarding including data in institutional repository services is the question of whether data can simply be incorporated into the existing infrastructure, or whether it requires separate handling and applications. Some authors have enforced the distinction suggested by Lynch and Lippincott (2005), pointing to "publication" or "data" repositories as sites worthy of differentiation.

Early efforts involved developing open-source digital repository applications such as Fedora (developed at Cornell in 1997) and DSpace (developed at MIT, with support from Hewlett-Packard, in 2002 and now managed, developed, coordinated, and supported under the umbrella of DuraSpace). Together with ePrints, these applications represent many of the institutional repository platforms currently in use in the US and form the basis of institutional repository services. Universities customize their applications of these platforms (sometimes extensively). They depend on (and ideally participate in) networks of developers and committers who improve and update the open source products. Universities might choose one application, or maintain multiple repository platforms, perhaps specifying separate instances for publication or data repositories.

Each of the three institutions participating in the DISC-UK DataShare Project (2007-2009), a JISC-funded project to investigate and "contribute to new models, workflows and tools for academic data sharing ... " in institutional repositories, relied on its own distinct repository platforms—DSpace, ePrints, Fedora—for institutional datasets. Two

Many universities are developing such repository lists... and UNC should collaborate with others whenever possible to ensure deep coverage." (Marchionini et al.,18).

of the institutions involved incorporated data into their existing repository instance, while one launched a dedicated repository for datasets, though using the same application as the parallel institutional publication repository. A key conclusion of the project was that “IRs can improve impact of sharing data over the internet” (Rice 2009, 5).

As Macdonald and Martinez-Uribe (2010) summarize the evolution of institutional repositories toward data:

Integral to the whole research base are research outputs such as publications and digital data as both evidence and the means to verify intellectual endeavor. University strategies to harvest these products have developed around the concept of digital repositories developed by academic libraries. The first realization of such information systems were publication repositories built to manage and disseminate research articles and aimed to provide open access to a significant proportion of newly published academic papers. The development of research data repositories has been seen as the next coherent step in the growth of repositories. (6)

However, this coherent step toward research data repositories has not been collectively trod. Rather, as these examples demonstrate, they have been undertaken in different ways. Currently, US research universities offer distinct and different repository services and tools. While recognizing the importance of and the need to intervene in the data lifecycle management of their researchers, universities have developed and expanded their repository services along particular lines, according to institutional priority. These distinct approaches continue to evolve.

8. Disciplinary or subject repositories, approaches, and issues [report pending]

9. Federated and collaborative approaches (Deborah M. Ludwig & Michael Bolton)

*Databib*⁹ currently lists 573 repositories for research data, many of which involve partnerships or collaborative efforts. The literature about federated and collaborative approaches to managing and sharing research data looks at participative communities that build shared services and infrastructure in a variety of disciplinary and institutional contexts. This review covers literature pertaining to a few federated or collaborative examples.

⁹ Databib is a tool for discovering sources of online research data.

Collaborative Communities. Beyond a place to store and access data, successful long-term curation of research data involves communities of practice, working alongside disciplinary specialists to develop tools, standards, best practices, and programs for education and training. The National Science Foundation (NSF) *DataNet*¹⁰ program envisions, “the creation of new (virtual) organizations [integrating] library and archival sciences, cyberinfrastructure, computer and information sciences, and domain science expertise....” (Hanisch and Choudhury 2009, 2) and Treloar (2009), highlight the NSF *DataNet* call

...to build sustainable infrastructure by creating a new type of organization that the NSF does not believe exists today. It is looking for librarians, archivists, and computer/computational/ information scientists who will work together to build excellent infrastructure for science and/or engineering, while engaging closely with intended users; domain scientists will be full partners in the process” (134-135).

Michener et al. (2011) describes the *Data Observation Network for Earth (DataONE)* initiative as a participatory endeavor, engaging scientists as the primary stakeholders in the network with numerous secondary stakeholder communities. Libraries have been prioritized as the most important secondary community network “because integrative science is data-driven and information-reliant and because libraries provide support services in each of the five [*DataONE*] science research environments” (7-8).

Many authors of articles and papers about research data emphasize the community aspects of data sharing and stewardship (Allard 2012, Michener et al. 2011, Hanisch and Choudhury 2009, Schaeffer et al. 2011, Treloar 2009, Williford and Henry 2012). Collaborative communities allow people with different skills and knowledge to work together and to contribute to successful and complex solutions. Partnerships create a foundation for success by building on the knowledge and skills of a rich community of people. Without tools and standards to make data management and data sharing practical and without the requisite best practices, education and training necessary to adhere to the established standards and best practices, data stewardship projects and programs may not be sustainable (Williford and Henry 2012). Many of the projects noted in this review of the literature offer perspectives on the essential human connections.

What do we mean by federated approaches? Heimbigner and McLeod (1985) defined a federated database architecture as one that

¹⁰ http://www.nsf.gov/funding/pgm_summ.jsp?pims_id=503141

...allows a collection of database systems (components) to unite into a loosely coupled federation in order to share and exchange information. The term federation refers to the collection of constituent databases participating in the federated database" (254).

In a data repository context, we can conceptualize federations as a distributed approach to storage or repository services that manage, share and perhaps curate data from multiple sites. The point of federation is often at a registry level, which houses metadata and pointers to the actual location of the data within a larger connected network (Allard 2012; Michener et al. 2011; Treloar 2009; Warner et al. 2007).

Williford and Henry (2012), in discussing the *Digging Into Data Challenge* program¹¹ sponsored by ten international research funders, refer to "a digital ecology of data, algorithms, metadata, analytical and visualization tools, and new forms of scholarly expression"(2). The authors note, "the digital raw materials upon which today's humanists and social scientists rely are heterogeneous, complex, and as massive as 'big data' in the sciences"(2-3). They recommend the adoption of more models for sharing, noting that it "makes no sense to replicate resources, skills, and services at all colleges and universities"(4) and that agreement to work together offers mutual benefit.

Many articles note the challenges of cross-repository interoperability based on differences between disciplinary communities and the need for an interoperability framework to connect many heterogeneous systems housing data. Warner et al. (2007) discusses the *Pathways* project, a partnership of Cornell and Los Alamos National Laboratory to develop a lightweight interoperable data model and workflow to better enable interoperability across complex metadata and the formats represented.

Federated Approaches to Data Management in the Sciences. Two exemplars of federated approaches to managing scientific research data are found in *DataNet* projects the National Science Foundation funded in 2009.

- 1) The *Data Observation Network for Earth (DataONE)* with headquarters at the University of New Mexico, a member institution of the Greater Western Library Alliance.
- 2) The *Data Conservancy*¹² with headquarters at the Sheridan Libraries at Johns Hopkins University.

¹¹ <http://www.diggingintodata.org/>

¹² <http://dataconservancy.org/>

DataONE

Allard (2012) and Michener et al. (2012) discuss *DataONE* as a multi-institutional, multi-national, and interdisciplinary collaboration to support the full information lifecycle of biological, ecological, and environmental data and to provide tools for researchers, educators, and the public. Michener et al. (2012) explains that environmental sciences represent a challenge for discovery because of their extremely heterogeneous data. *DataONE* reaches across different scientific disciplines and institutions for sharing of data and findings, expertise and tools, and for the development and utilization of compatible data management strategies and best practices.

Allard (2012) and Michener et al. (2012) describe *DataONE's* distributed technical architecture of member nodes and coordinating nodes or repositories with accompanying infrastructure and services. The overall architecture includes repository services and data replication as well as the *Mercury* metadata catalog for data discovery across geographically distributed member node repositories. The project addresses the need for a consistent researcher interface with analysis and visualization tools spanning member repositories. The services include quality metadata, shared identity management and access control policies.

Michener, et al. (2012) and Allard (2012) both describe *DataONE's* virtual community as one that creates a strong network of people and working groups comprised of scientists, academic researchers, educators, government and industry representatives, and leading computer, information, and library scientists. To encourage a cross-disciplinary approach, the network's scientific communities are not categorized by domain. *DataONE* facilitates data preservation and re-use, interoperability solutions, and best practices for data management across its lifecycle. Michener et. al. (2012) describes scientists as *DataONE's* primary stakeholder with libraries as the most important secondary stakeholder community. Library partners include the Libraries at the University of New Mexico; the College of Communication and Information, University of Tennessee; and the UC Curation Center, California Digital Library, University of California. Allard (2012) concludes that *DataONE* is a successful partnering of librarians, information science researchers, and scientists and notes that community engagement of scientists is an important user-centric focus of the *DataONE* project.

The Data Conservancy

Hanisch and Choudhury (2009) describe the *DataNet Data Conservancy* project as a partnership between John Hopkins University's Sheridan Library, the US National Virtual Observatory, and a wider list of partners found on the web site¹³ that includes

¹³ <http://dataconservancy.org/community/partners/>

several national research centers along with Cornell, the University of Illinois, and the University of California, Los Angeles. Use cases in astronomy, seismology, and international land use policy have informed initial development. “The Data Conservancy team is interdisciplinary and multifaceted, and rooted in the university research library at the John Hopkins University” (6). The authors explain three key terms: stewardship, sustainability, and multiple-scales.

As with *DataONE*, Hanisch and Choudhury (2009) describe the technical elements of the federation as a common user access interface layer that includes a registry of aggregated metadata, database query and data access, and distributed storage for the actual data objects. The project envisions a black box repository environment that could be widely deployed by various organizations from which metadata is replicated out to *The Data Conservancy*.

Both *DataONE* and *The Data Conservancy* initiatives represent national and international approaches to federated management of scientific research data with metadata for discovery and tools or services providing points of federation for discovering and using distributed collections of data.

Federated Approaches to Data Management in the Social Sciences. Two social science data initiatives, *ICPSR* and the *IQSS Dataverse* illustrate federated approaches to social science research data management.

The Inter-university Consortium for Political and Social Research (ICPSR)

*ICPSR*¹⁴ is a non-profit, membership-based, centralized data archive located at the University of Michigan and focused on social science data. Green and Gutmann (2007) from *ICPSR* discuss two primary repository types: institutional and discipline/domain. *ICPSR* represents the latter type of repository. “The current digital repository landscape is made up of a blend of repository types. Repositories can be grouped into two broad categories: Institutional digital repositories with no specific discipline focus and discipline or domain-specific data archives” (38). The authors note that many institutional repositories (e.g., *asudigitalrepository*¹⁵, *KU ScholarWorks*¹⁶, *MOSpace*¹⁷, *Scholars Bank*¹⁸) may allow deposit of data sets, “but support services for data processing, metadata production, or analysis are not usually offered as part of the repository service....[T]hese repositories position themselves at or near the end of the

¹⁴ <http://www.icpsr.umich.edu/>

¹⁵ <http://repository.asu.edu/>

¹⁶ <http://kuscholarworks.ku.edu>

¹⁷ <https://mospace.umsystem.edu/xmlui/>

¹⁸ http://library.uoregon.edu/diglib/irg/SB_Role.html

scientific research life cycle. Their goal is less to partner with researchers or with domain-specific repositories throughout the research life cycle than it is to garner the value of the institution's productivity, to gather this productivity, and possibly to lower the local or community-wide cost of scholarly publications" (38-39).

The authors remind us that the social science domain reflected in discipline-specific repositories have been in existence for decades. "Rather than focusing on publication-related materials from multiple subjects areas within a single organization, domain-specific digital repositories hold collections of materials grouped by type, subject, or purpose and intrinsically support domain- or discipline-oriented research needs. Domain-specific digital repositories in the social sciences have a history of providing infrastructure for data sharing and strive to provide support throughout the data life cycle. These data archives hold the raw materials that faculty and students can reuse, repurpose, analyze, and recompile in teaching, learning, and research environments" (39).

Both disciplinary and institutional repositories have shared goals as well as distinctive differences. The authors note the similarities of shared metadata and common discovery platforms such as Google Scholar. "Resource discovery in the social sciences now extends far beyond consulting a stand-alone research aid or search tool. Alliances have been developed across repositories and a new set of tools allowing researchers to do complex and innovative searches to locate and explore data is emerging" (41).

Dataverse, the Institute for Quantitative Social Science, Harvard University

The Harvard *Dataverse* Network¹⁹ is described by Crosas (2012) as a collection of social science research data contained in virtual data archives called "dataverses" open to all researchers worldwide to share, cite, reuse and archive research data. Institutions can also use the open source *Dataverse* software to develop their own local institutional data repositories, which through metadata harvesting and sharing can then become federated partners within the larger *Dataverse* network. Repositories based on other software can also share metadata with the *Dataverse* network using the OAI-PMH (Open Archives Initiative Protocol for Metadata Harvesting) protocol.

According to Crosas (2012), a key focus for *Dataverse* is maintaining the researchers connection to the data.

Most professional archives, although often considered the most reliable solution, do not usually facilitate control and ownership of the data by the author. Once the author submits the data, the archive becomes fully

¹⁹ <http://dvn.iq.harvard.edu/dvn/>

responsible for the data management, cataloging and future updates. While this can be advantageous for some researchers, many prefer to maintain control of their data and to receive increased recognition.

The author explains that *Dataverse* allows researchers to create “virtual collections” called *Dataverses*, which include custom branding for the researcher or research team – helping meet the researcher’s needs for recognition, visibility and ownership.

Federated Approaches to Data Management in the Humanities and Arts. Researchers in the humanities and in the arts create data that may be both big and complex.

Humanities Projects at Scale – The Digging Into Data Challenge

Williford and Henry (2012) cover eight humanities collaborations sponsored by the previously mentioned *Digging Into Data Challenge*. All projects engage with large data corpora, apply computational analysis, require collaboration from a variety of professionals, and conduct a research process. All projects cross disciplinary and international boundaries. The authors note some of the challenges these multi-institutional, multi-disciplinary, multi-national partnerships have encountered including working as virtual teams, and the lack of effective training for students and junior scholars pursuing computationally intensive research. Long-term sustainability of the projects was an additional concern of participants in the *Challenge*. The authors make nine recommendations and expand on these in their report:

1. Expand our concept of research.
2. Expand our concept of research data and accept the challenges that digital research data represent.
3. Embrace interdisciplinary.
4. Take a more inclusive approach to collaboration.
5. Address major gaps in training.
6. Adopt models for sharing credit among collaborators.
7. Adopt models for sharing resources among institutions.
8. Re-envision scholarly publication.
9. Make greater sustained institutional investments in human infrastructure and cyberinfrastructure (2-3).

Digging Into Image Data... (DID-ARQ)

DID-ACQ is one of the Simeone, et al. (2011) speaks of lessons learned from the international *Digging Into Image Data to Answer Authorship-Related Questions (DID-ARQ)*²⁰. This project looked for ways to examine an archive of images too large to

²⁰ <http://isda.ncsa.illinois.edu/DID/>

examine manually. It provides a template for future collaborations involving multiple datasets in geographically distributed locations. The project includes researchers from the University of Illinois, the National Center for Supercomputing Applications, Michigan State University, and the University of Sheffield.

TextGrid

Another example of a federated humanities data management is the *TextGrid*²¹ project involving ten institutions in Germany, which was described by Neuroth, et al. (2011). The authors describe *TextGrid* as the first large, multi-year project in Germany dealing with developing a research infrastructure and virtual research environment for the arts and humanities. The environment has two parts: an entry point to a virtual research environment and a data archive (or set of archives). Disciplinary communities served include textual philology, linguistics, art history, classical philology, and musicology. *TextGrid* is part of the scientific *D-Grid* program affiliated with various e-Science endeavors. Within the *TextGrid* context, federation refers to the shared development of tools for the virtual research environment, the shared cost model, and a future federated repository infrastructure.

Closer to Home: Collaborative Research Data Projects of GWLA / GPN Members. In addition to *DataONE*, several projects have been developed at Greater Western Library Alliance and/or Great Plains Network institutions that exemplify possibilities for managing research data. These projects include *tDAR* and *MaizeGDB* described below.

tDAR

Spielmann and Kintigh (2010) discuss development of *the Digital Archaeological Record (tDAR)*²², an international digital repository for data access and preservation from archaeological investigations. The *tDAR* repository encompasses datasets, documents, and images from current archaeological research and legacy data. *tDAR* is part of the Digital Antiquity partnership, a multi-institutional group that ensures sustainability.

MaizeGDB

The *MaizeGDB* project, described by Schaeffer, et al. (2011) and Lawrence, et al. (2004), is a public database that serves the community of maize researchers by storing and curating data related to the genetics and genomics of maize (corn). The “curation focus is to facilitate data integration of very large data sets and to provide insight into development of easy-to-use interfaces and data displays. The efforts toward data integration involve gene nomenclature considerations as well as ontology development

²¹ <http://www.textgrid.de/en/community/>

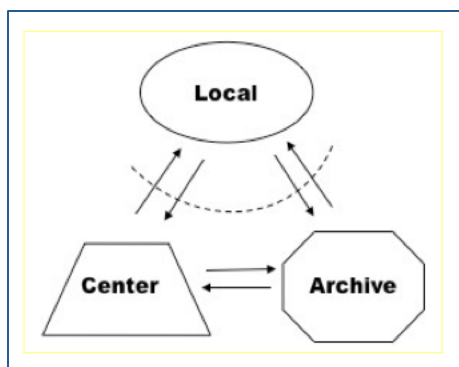
²² <http://www.tdar.org/>

and implementation”(1). This data repository is over 20 years old and has matured from a focus on comprehensive curation of the literature, genetic maps and stocks to a current approach including the recent release of a reference maize genome sequence, multiple diverse maize genomes and sequence-based gene expression data. Originally housed at the University of Missouri, the Maize DB is the work product of a research community that provide data, establish nomenclature standards and recommends directions, priorities and strategies.

Federation Between Different Types of Repositories. The repository is a key element in federating long-term stewardship of data, reflecting the often-rich research partnerships built around research data. Baker and Yarmey (2009) note distinctive types of repositories that have different goals, participants, and purpose; but these repositories are also similar in terms of roles, activities and responsibilities. The authors consider repository relationships from the vantage point of connectedness to, or distance from, the researcher. The repository first created by, or for, the researcher is more connected to the active research project while a repository that provides “final” archiving of a dataset is further away from the researcher.

The authors illustrate the relationship of different types of repositories and the data flow between three different categories of repositories:

- local repositories where data lives during the active research process
- center repositories where data may be accessed during or after a project
- archival repositories concerned with long term data preservation



Data may begin the lifecycle in a repository largely controlled by the researcher(s) involved in the acquisition of research to answer a particular question. As the data matures and is ready for sharing, it may move to a “center” repository with the function of making data available or shared with a broader audience, and finally, if the data has future value in disciplinary and interdisciplinary research, the data may move to an archival

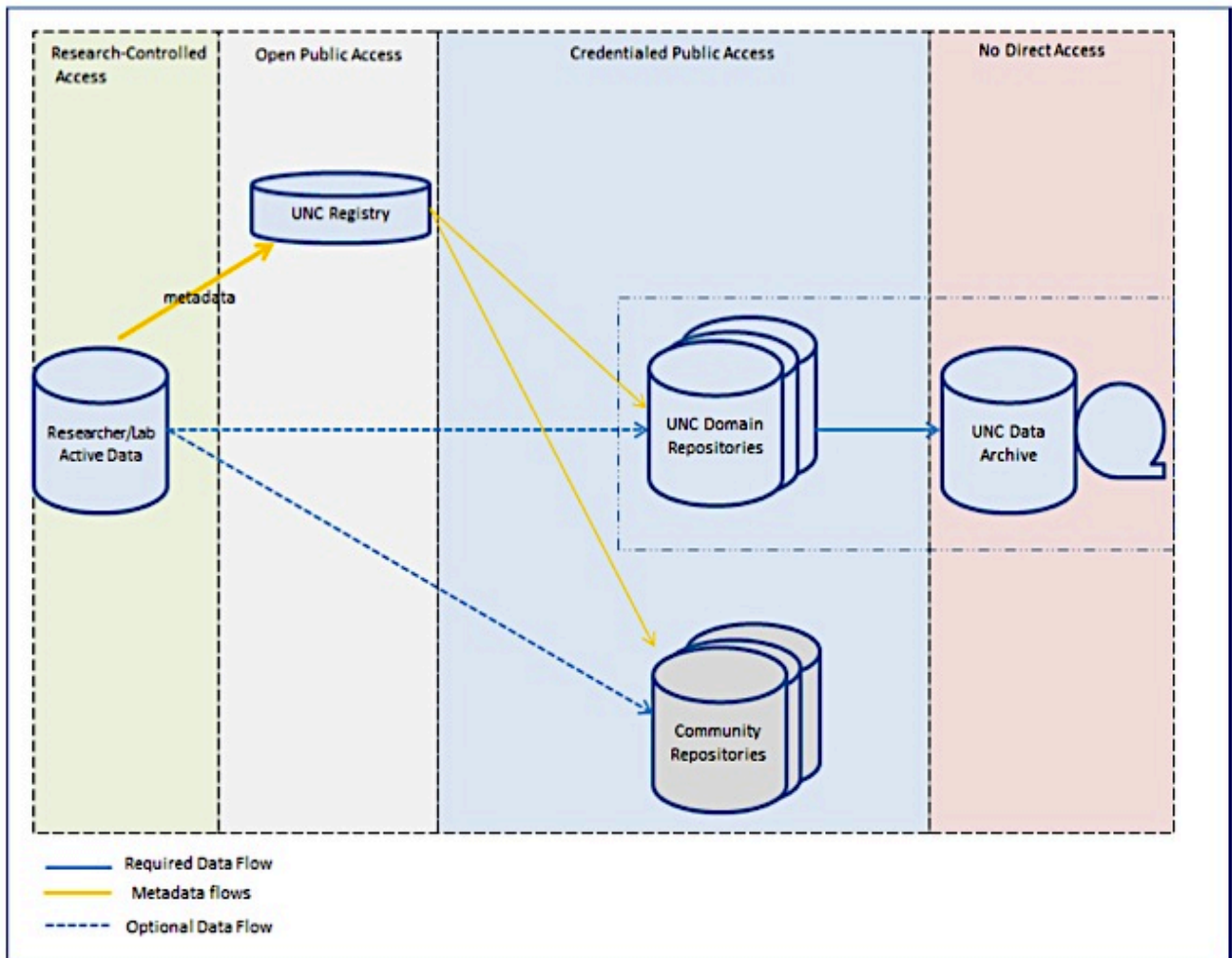
repository for the purpose of ensuring future re-use.

Baker and Yarmey (2009) explain that the three types of repositories form a collective approach to research data management with each type of repository filling a different role within the federation and serving a different audience. Data can flow in both directions within the federation. For example, the archive may be used to restore data to

the local researcher's repository or, in some cases, data might move directly from the local repository to an archive. Data might be synchronized between an archive and center (accessible) archive so changes and new data versions are automatically sent to the archive. An archive might periodically check to see if its version of the data matches that of a center repository and initiate appropriate events when the two versions differ. The *Digital Preservation Network (DPN)*²³, for example, represents a specialized type of federation to reserve research data and other digital artifacts of interest, including cultural heritage data. Baker and Yarmey (2009) would consider this a federation of dark archives used for preservation and restoration of data access when lost by a “center” archive.

Connecting Local and Remote Repositories – the University of North Carolina View.

The University of North Carolina (Marchionini, et al. 2012) undertook a campus study of digital research data stewardship. Working under a charge from the Provost, the report



²³ <http://www.dpn.org/>

describes how a task force conducted an environmental scan of research data stewardship policies and trends, discussed issues, and using interviews and a survey collected data on campus. The task force developed a set of principles and associated courses of action for the campus to consider. The UNC data flow diagram (above) shows expected interaction of UNC's repositories and with possible community/disciplinary repositories. This diagram shows metadata flowing from a researcher or lab to a university data registry for public discovery with metadata and optional data flows to university repositories or community repositories as appropriate. Data flows to university domain repositories would be coupled with deposit into a "dark" archive for presentation that excludes public access and could provide future restoration of the data to the access repository if needed.

In Summary: Federated and Collaborative Data Management in the Literature.

Authors cited in this literature review point to the need for data to remain close to the researchers and for the researchers to play a major role in the development of any plan for the management, stewardship and curation of research data. The literature suggests developing a community to tackle complex issues of data management in a user-centric approach, broadly inclusive of partners with different skills and knowledge.

Various authors (particularly Green and Gutmann 2007) also point out that both institutions and disciplinary communities are in the business of building and maintaining repositories for access to data. Universities need to consider and articulate how they wish to balance expectations for data sharing between local campus and non-local repositories. Disciplinary communities may also need to consider how to work with other disciplinary and with local institutional repositories to federate discovery of data for researchers and perhaps to provide common sets of tools for interacting with the data.

Data is complex, originating from a variety of domains. Scientists, social scientists, and humanists use data that extends beyond numeric and encompasses a wide variety of formats. Data management planning needs to take into account the various types of data to be curated. The University of North Carolina captures data complexity in the following assumption.

The term 'research data' is considered in its broadest form and includes the digital traces of processing and documentation associated with the research enterprise. We consider only electronic (digital) research data, yet, what this means varies enormously across disciplines and scholars. We believe that disciplinary communities (including funding agencies for those communities)

are best prepared to define what constitutes ‘research data’ for their respective fields and that university policies and implementations should be guided by what UNC scholars and their communities define as research data. We also note that more than 40% of UNC researchers in our survey reported regular use of non-digital text, hand-written notes, and other forms of non-digital data as part of their research. Another one-third of our respondents report using biological, organic, or inorganic samples and/or specimens in their research. As UNC develops policies regarding research data – digital or otherwise – this enormous diversity must be accommodated and embraced. UNC simply cannot seek to develop a “one-size-fits-all” policy in this environment (Marchionini, et al. 2012, 5).

10. Economics/costs of data curation [report pending]

11. Archiving and preservation (Amalia Monroe-Gulick)

The preservation of digital data and objects is a major concern in both the research and archivist communities. In a 2002 report from the Research Libraries Group (RLG), digital preservation was defined as “the managed activities necessary for ensuring both the long-term maintenance of a byte stream and continued accessibility of its contents” (Beagrie, et al. 2012, 3). The archiving and preservation of digital data has been a significant challenge to researchers and information professionals since the beginning of the “data deluge.” Berman (2008) argues that it is assumed that digital data will always be available; however, it is actually fragile and can easily be lost due to many reasons including technology failure, as well as software and storage media becoming outdated.

According to Cragin, et al., “sharing is at the heart of success, as collecting, storing and making use of data can only [occur] after the means for sharing are in place” (2012, 1023). Another major concern are the lack of incentives and even disincentives for scientists to make data available for sharing. Borgman (2007) identifies disincentives for sharing which include the lack of rewards for effective data management, the amount of effort to preserve data for later use, and scientists desire to share data only after publication. In addition, Tenopir, et al. (2011) indicate scientists identify problems with funding and time as the main obstacles to sharing data. However, many funding agencies, such as the NSF, are now requiring data to be made available. Beyond funding requirements, one author argues that the scientific community needs to be fully aware of

the importance of the continuum of data that is necessary for the advancement of future research (Durr 2008).

Jantz and Giarlo (2005) argue that no matter the scale of the data preservation purpose (i.e. “big data” or “small data”), methods, policies, standards, and technologies all need to be considered. In an article addressing data preservation for astronomy, they found that while many research facilities provide the archival space, a complete cycle - which includes capture, curation, preservation, and long-term access - is missing (Choudhury, et al. 2007).

One of the first steps in establishing a complete cycle or “end to end process” is risk assessment and management. Risk management of data preservation is another important issue, and challenge, when deciding how to preserve digital objects. Different types of preservation actions in terms of risk management for archivist and scientists can be clarified through the framework a Preservation Network Model (PNM) Conway, et al. (2012). The goal of PNM is to capture the interrelationships of the data, software, and other digital objects in order to accurately assess the acceptable risks. The preservation strategies identified are:

- Risk acceptance and monitoring
- Capture of software and extension through the stack
- Description
- Migration

These different strategies require that the risks with preserving each type of digital object first be identified. In addition, different categories of data have special considerations that need to be addressed during these initial assessment and planning phases. For example, McGarva (2009) identifies the specific actions with Geospatial data, including: format, systems, legal, and community actions.

Akmon (2011) indicates that another essential element for successful data preservation is developing of a preservation plan early in the research lifecycle. Because of the unique lifecycle of different types of data, this is an issue that presents challenges and cannot be completely addressed by standardized policies. Wallis, et al. (2012) discuss the importance of understanding different types of data lifecycles (including e-science and government documents) as essential, especially in interdisciplinary research. In an article discussing the preservation of remotely sensed data, Faundeen reports on the data migration systems used by the U.S. Geological Survey which address the major issue of technological obsolescence, and argues that “planning for these preservation activities is extensive and must be done before data are threatened” (2003, 162).

Akmon (2011) argues that libraries and archivists can play an important role in the preserving and archiving of data, even without specific domain knowledge. Tenopir, et al. (2012) found that most academic libraries do not develop research data policies for universities. However, because of the expertise they can offer in archiving and preservation policies, they could serve as clearinghouse for potential policies.

Jantz and Giarlo (2005) found that there are many different types of systems and methods for long-term access to research data. Institutional repositories managed by libraries, with established guidelines and procedures are one option that has been greatly explored. Another potential method for dealing with the challenge is utilizing cloud technology. Mattmann et al (2010) report on storing and preserving NASA data with cloud technology. However, regardless of the type of system, the same fundamental questions of long-term needs must be addressed.

This review of the literature addresses some of the issues and potential solutions to the many challenges facing those attempting to preserve data. This is an evolving topic that will require ongoing exploration and discussion as the most effective, scalable, and discipline-specific processes to ensure long-term access and simplified sharing of data develop.

12. Metadata and description (Andrew Johnson)

Metadata is central to the lifecycle management and long-term preservation of research data. A number of metadata models and schemes have been developed to describe data from particular domains or disciplines. These include the Data Documentation Initiative for social science data and the Ecology Metadata Language for ecology data. There have also been attempts to describe scientific data more broadly. Matthews et al. (2010) developed a general model, called the Core Scientific Metadata Model (CSMD), to represent scientific study metadata for data generated in large-scale scientific facilities. The CSMD was designed to provide a core model that can be extended when necessary. Thus, the CSMD can be more specialized than general metadata models, yet broader than those that only describe a particular scientific domain (Matthews et al. 2010, 114).

In addition to describing research data in either general or domain-specific terms, metadata is also valuable for discovering and citing research data sets. The DataCite Metadata Scheme attempts to address the issues of data discovery and citation at a global level (Starr and Gastl, 2011). The DataCite schema features a small set of mandatory elements with a larger optional set that can be used for more detailed description. The core elements include general descriptive information as well as a

Digital Object Identifier (DOI) to describe relationships between the cited data set and other objects.

Metadata plays a key role in the long-term preservation of research data. In order for research data to remain authentic and useable over time, it is necessary to capture contextual and administrative metadata as well as technical metadata about the digital object(s). As Wilson (2010) notes, this latter type of technical metadata, exemplified by the PREMIS standard, is not solely sufficient for long-term preservation of research data, and he argues that archivists and scientific researchers should work together to ensure that all of the necessary “recordkeeping” metadata, including contextual and administrative metadata as well as technical metadata, is captured in order to maintain data integrity, authenticity, reliability, and usability over time (215). Shaon and Woolf (2012) provide an example of one such attempt to incorporate contextual, domain-specific metadata with technical preservation metadata. In their work with the Infrastructure for Spatial Information in the European Community (INSPIRE) Spatial Data Infrastructure, they convey that INSPIRE’s application of the ISO 19115 metadata model for geospatial data captured important contextual information necessary for discovery and understanding of data; however, it failed to meet preservation metadata requirements like those described in the Open Archival Information System (OAIS) Reference Model. To solve this problem, they developed a preservation profile of ISO 19115 that includes both preservation metadata adhering to OAIS and PREMIS specifications as well as core ISO 19115 metadata that allows for accurate description and contextualization of geospatial data.

Despite the emphasis in the literature on capturing a variety of types of metadata throughout the research data lifecycle, there are numerous barriers and challenges to achieving this goal in practice. Mayernik (2010) discusses significant metadata challenges for curating and sharing data. He notes that responsibility for metadata creation is ambiguous and could include information professionals, scientists, and hardware/software tools (Mayernik 2010, 2). Another challenge lies in the discrepancy between the information professionals’ formal approach to metadata creation and the informal practices of scientists. This discrepancy is further complicated by the distributed nature of metadata creation in research settings as well as the changing role of metadata through different stages of the data lifecycle (Mayernik 2010, 2). In additional work on the metadata practices of scientists, Mayernik (2011) describes how researchers who rarely create documentation beyond what is required for their own use also rarely share data with users outside of their immediate projects (248). When asked to create metadata for a widely available metadata registry, the researchers Mayernik

(2011) observed still described data in ways targeted mostly to researchers with similar expertise (11). Mayernik (2011) argues that the metadata necessary for widespread data sharing is not in line with the current metadata practices of many researchers (280).

Ferguson (2012) describes additional metadata challenges for biomedical researchers attempting to understand and repurpose publicly available microarray data sets, and she notes that many of these data sets lack sufficient contextual information and metadata (51). In an effort to address this issue, a team of data curators at the Bioinformatics Core Group in the Harvard School of Public Health sought to annotate and contextualize some of this publicly available data (Ferguson 2012, 52). According to Ferguson (2012), the data curators used a suite of open source software tools to provide metadata about the experiments and data sets. However, even with tools and standards available to assist with metadata creation, she describes this process as time consuming and requiring significant subject matter expertise (55). Ferguson (2012) also notes that the data curators expressed uncertainty about how much metadata was necessary to allow for discovery and reuse (55).

Greenberg et al. (2009) also address the question of what constitutes sufficient metadata in their work developing metadata best practices for the Dryad data repository. They determined that the amount and type of metadata for data ingested into digital repositories should support both the immediate operational needs and long-term project goals of the repository (208). For example, Dryad's metadata approach addresses the immediate need of making content available in a DSpace repository via an XML schema while responding to the long-term goal of alignment with the Semantic Web via a metadata application profile (Greenberg et al. 2009, 209).

Brownlee (2009) focused on metadata challenges particular to research data in general purpose institutional repositories (IRs), and she concludes that metadata used to describe and manage data can be dissimilar to the metadata required for submission to IRs (2). At the University of Sydney, metadata describing research data was found to be stored primarily in databases or Excel spreadsheets with variability in use of and conformance to standards (Brownlee 2009, 3). The University of Sydney Library chose to address this variability and lack of standardization by ingesting the native metadata records as digital objects along with Dublin Core (DC) metadata describing these objects (Brownlee 2009, 6). Other options the library considered were mapping native metadata to DC without ingesting the original files and creating a custom metadata schema in the IR for each new native metadata element set (Brownlee 2009, 4-6).

Variability among numerous discipline-specific metadata schemes can create barriers to interdisciplinary research. Willis et al. (2012) examined nine metadata schemes for physical, life, and social science data to identify common objectives across schemes. Their analysis indicated that despite being constrained by historical disciplinary practices and workflows many metadata-driven goals are largely independent of scientific discipline or data type (32). Willis et al. (2012) identified eleven fundamental metadata goals for documenting scientific data in order to enable sharing across disciplines, including: abstraction, extensibility, flexibility, modularity, comprehensiveness, sufficiency, simplicity, data interchange, retrieval, archiving, and publication (27). Willis et al. (2012) suggest that further research on these common goals could help to break down the metadata barriers for interdisciplinary data sharing and research (33).

Section B. Results of Survey (Preliminary as of 5/15/2013)

Amalia Monroe-Gulick and Deborah Ludwig

The survey questions are found in Appendix G and were developed by Research Team B members, named in Appendix A. The survey was sent to contacts representing the highest level of administration for libraries, information technology, and offices of research.

18 unique institutions out of 47 surveyed provided responses from at least one contact. The response rate was highest from libraries, with 13 of the 23 responses. Responding institutions all offers some level of data management services.

Given a fairly low response rate with responses skewed toward respondents from libraries, it is not possible to gain a comprehensive picture of research data management services in GWLA and GPN institutions at this time; however there are some generalizations we can make from these results and additional comparative information may be gleaned from looking at other national studies such as the recently-conducted SPEC Kit Survey by IMLS.

In general, all institutions responded that they provide some level of data management services. Information Technology organizations seems more likely to provide short term or long term storage for research data and to be involved in high performance computing initiatives. Libraries seem more likely to be involved in helping researchers develop data management plans, locate repositories for data deposit, or help with deposit data in an institutional repository. Research centers are most likely to provide data analysis and visualization services. Formal education and training services in research data management are not common.

With respect to data services for storage, archiving, preservation, and sharing, the most commonly offered services were storage and help identifying repositories where researchers could deposit data. The services most frequently not offered are local institutional repositories. When repositories are provided, the libraries are the most common provider of that service.

Detailed Review of Data

Response Rate

Total Number of Surveys Sent Out	134
Total Number of Universities Contacted	47
Total Responses	23
# of Total Universities Responding	18
University-Level Response Rate	38%
Total Survey Response Rate	18%

- Out of the 47 individual universities contacted, there was a least one respondent from 18 (38%) of the universities.
- Three individuals were contacted at each university (libraries, IT, research administration), with a total of 134 surveys distributed. The resulting overall individual response rate was 18%.
- One university returned surveys from all of the departments/units contacted and two universities returned surveys from two of three departments/units contacted.
- Respondents also indicated if they were completing the survey on behalf of their organizational unit or on behalf of their institution. The majority of respondents indicated they completed the survey for their organizational unit (74%).

Organizational Unit/Department Responding

Library/Library Unit	13	57%
Information Technology/IT Department	2	9%
Research Studies	6	26%
Provost	1	4%
Digital Services	1	4%

- Over half of the surveys were completed by those affiliated with a university library or a library unit. The second largest group represented was research studies.
- The overall individual response rate was not high. When this is considered with the majority of respondents reporting information on their organizational units, this will limit the overall analysis of understanding data services at a university-wide level because not all organizational units may be aware of services, policies, and challenges associated with other university units. In addition, with over half of the respondents representing libraries, the results may also be skewed toward knowledge of library-based data services rather than campus-wide services.

Services

In the service provision section of the survey, five main categories were addressed:

1. General Support and Services for Research Data
2. Storage, Archiving, Preservation, and Sharing of Data
3. Accessing and Using Research Data
4. High Performance Computing

5. Support services management, preservation, and access to digital and non-digital research data for long-term access to campuses

General Support and Services for Research Data

In this section, questions were asked about support and service the following areas:

1. Development of data management plans
2. Consultation on data management
3. Consultation on options for data licensing agreements for open or restricted access
4. Dedicating funding resources that support long-term management of research data
5. Formal training on for researchers on data management planning
6. Ensuring university compliance for research data in accordance with commercial licenses, government regulations, and funding agency mandates

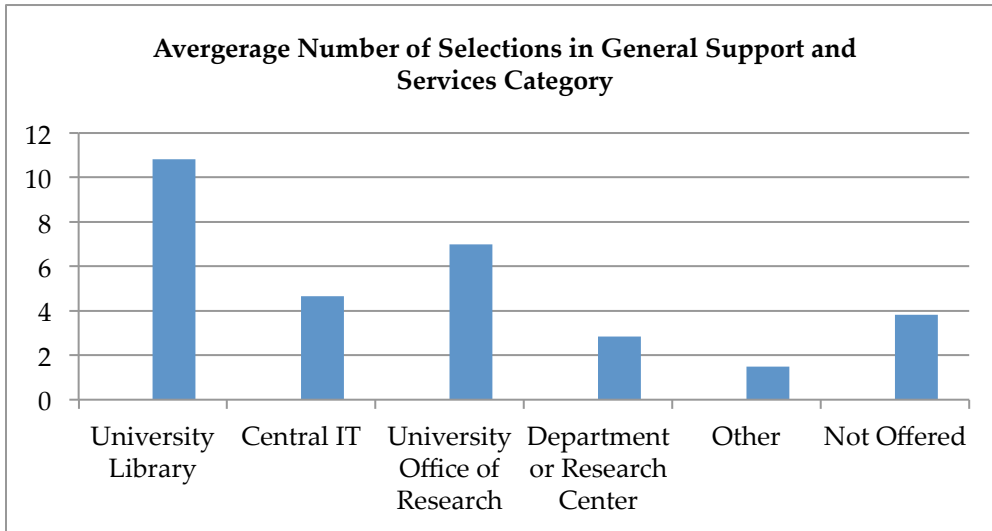
Overall, general support and services for data management is present at the responding universities.²⁴ The average number of positive responses for the entire category of multiple questions was 27.

Support for the development of data management plans received the largest number of responses, (38). Formal training for data management planning received the fewest responses (19). This result is interesting because of the federally mandated inclusion of data management or sharing plans for NSF and NIH grants. It appears that the responding universities are offering assistance with developing data management plans, but do not have formal programs in place. There is potential for further research in this area to identify the nature of data management consultation services and if they are decentralized and on an “ad hoc” basis.

The libraries had the highest frequency of responses regarding general data services, with an average of 11. The make-up of the survey respondents may have influenced this potential skewing towards the high rate for the library selection. However, the question on assistance with university compliance with commercial licenses, government regulations, and funding agency mandates received 20 overall responses, with offices of research receiving 13 of those 20 (54%). Offices of research were the second largest respondent group.

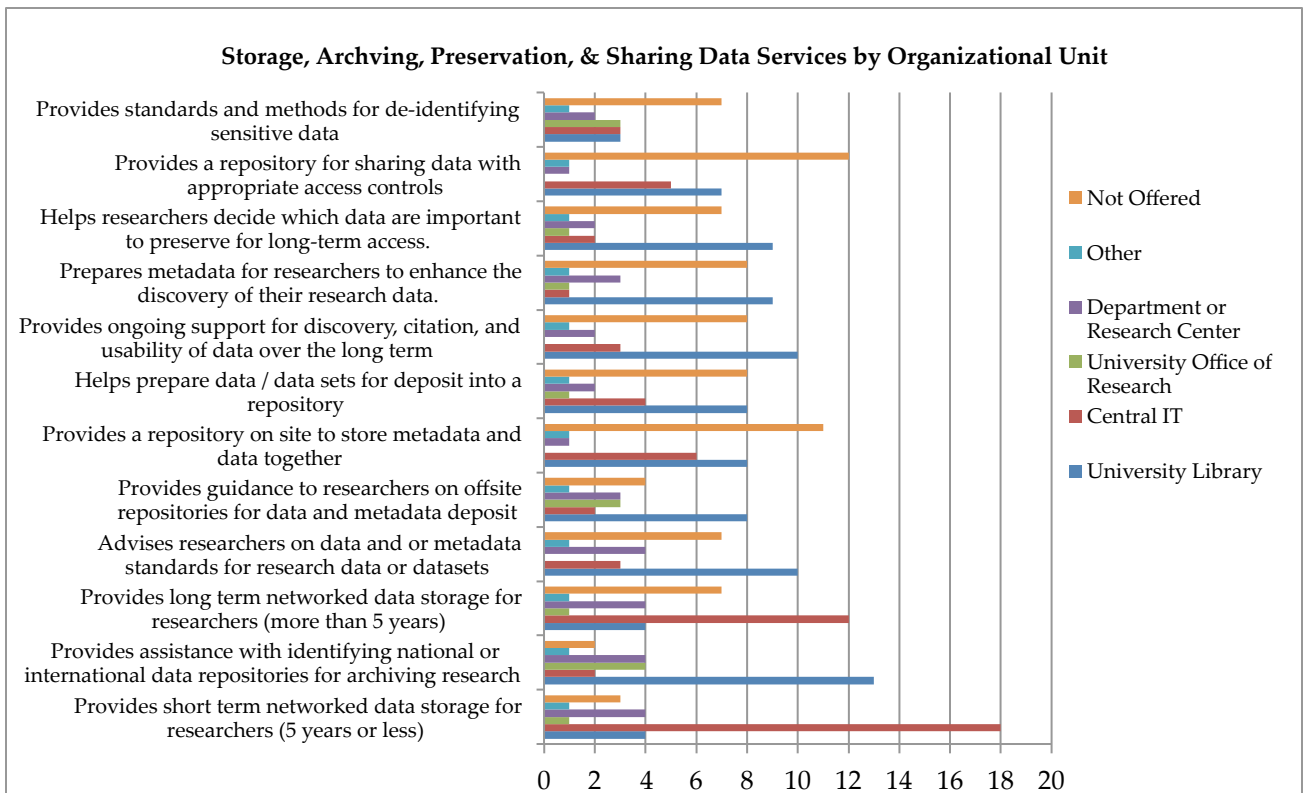
²⁴ The following discussion uses data that excluded the “not offered” choice.

All questions had responses that indicated a service was not offered, with the two highest being data management best practices (6) and dedicated funding resources (6).



Storage, Archiving, Preservation, and Sharing of Data

In this section of the survey, respondents were asked about services related to storing, archiving, preserving and sharing data. The average positive response rate was lower than the general services section, 18 compared with 27, indicating that this is an area that the responding universities have not addressed as much as general services, which focus more on issues related to data management.



The library and central IT options were most frequently selected in this category of questions. The library received the largest number of responses with the questions related to identification of repositories, metadata standards, and long-term discoverability and usability of data.

According to the results, central IT provides short-term networked data storage (5 years or less) with 18 responses out of a total 29 positive responses (62%). Respondents also indicated that central IT provides long-term data storage, but a smaller frequency (55%). This result is interesting, in part, because there were only two respondents from IT organizational units, but respondents representing different organizational units recognized the role of central IT in both long and short-term data storage.

The options of offices of research or research centers were rarely selected in questions related to storage, archiving, preservation, and storage. The office of research option was not selected for metadata standards, long-term support for discovery, citation, and usability of data, or providing a repository for sharing data.

However, in all questions, respondents indicated that the specific service was not offered through their unit or institution. Out of 26 total responses to the question “provides a repository for sharing data with appropriate access controls,” there were 12 (46%) responses of “not offered.” Also of note is the high rate of “not offered” responses to the question “provides a repository on site to store metadata and data together.” Out of the 27 total responses to this question, there were 11 (41%) “not offered” responses. Low rates of indicating these services are present in a unit or institution to these questions could indicate that universities are not yet offering centralized data repositories, for either short or long term. This is not an unexpected result because of the complicated nature of the issue.

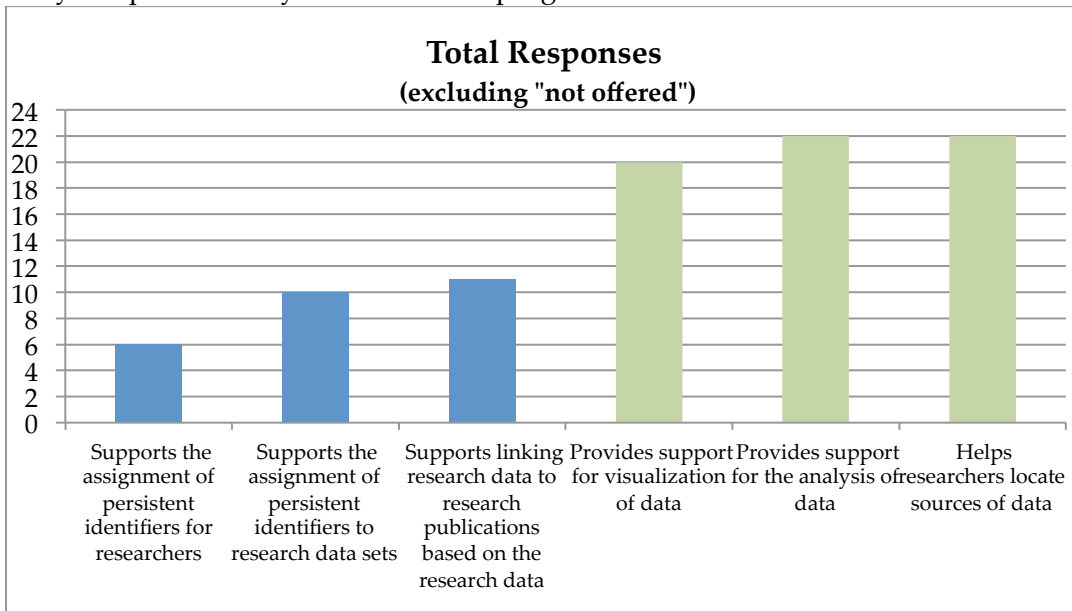
Accessing and Using Research Data

The section of the survey addressing accessing and using research data can be divided into two sections:

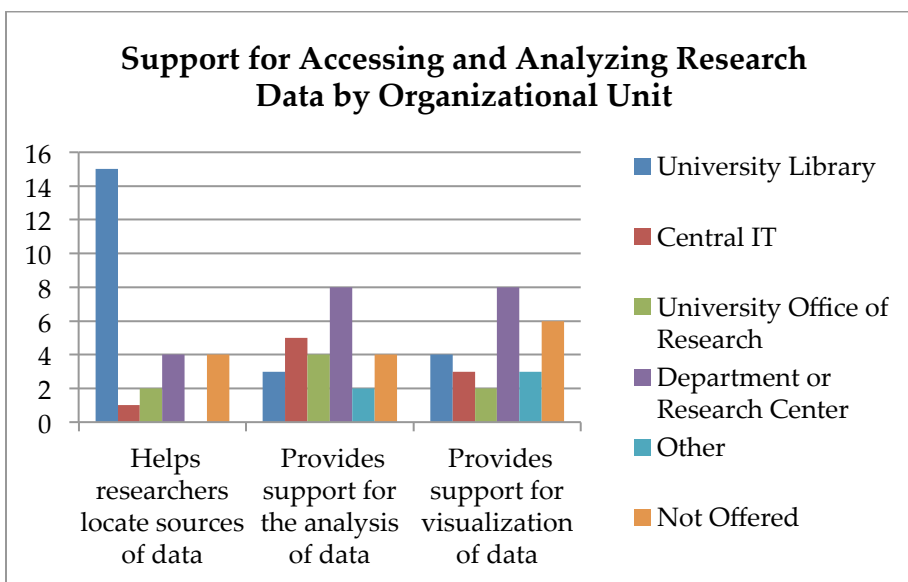
1. Support for locating and using/analyzing data
2. Support for linking and assigning unique identifiers to researchers and research data.

There was an evident gap between positive response rates for the question in the two sub-sections of the category. The average response rate (i.e. excluding “not offered” choice) was 21 for locating/analyzing data and 9 for linking/assigning unique identifiers. This is not a surprising result because supporting researchers with finding and using data has historically been the focus of data service programs at university libraries. However, linking research and research data to publications and assigning unique

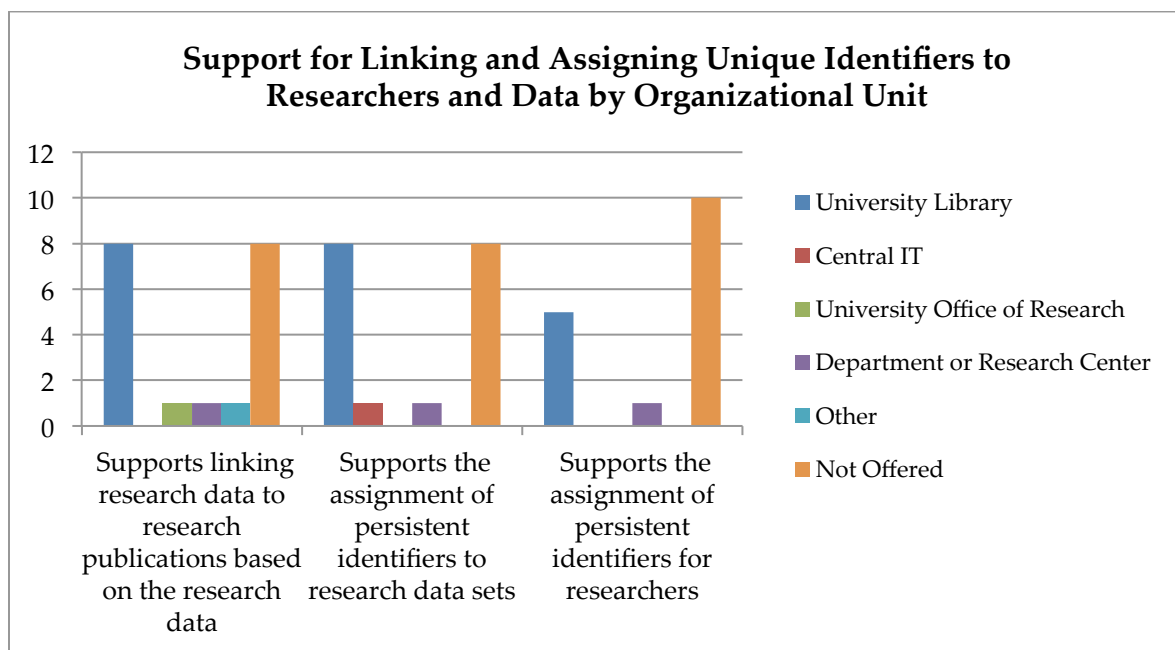
identifiers is an emerging area across universities, and university department, therefore, not yet represented by formal service programs.



The responding universities indicated that their libraries most frequently assist researchers with finding data, but campus and research centers most frequently offer data analysis and visualization support. However, “not offered” is the second most frequent response to data visualization support, demonstrating while services are being offered, this is still an area of potential growth for data support services at universities. Offices of research and central IT do not play a large role in this category, which is not necessarily surprising, and potentially indicates a trend of decentralized data services across campuses.



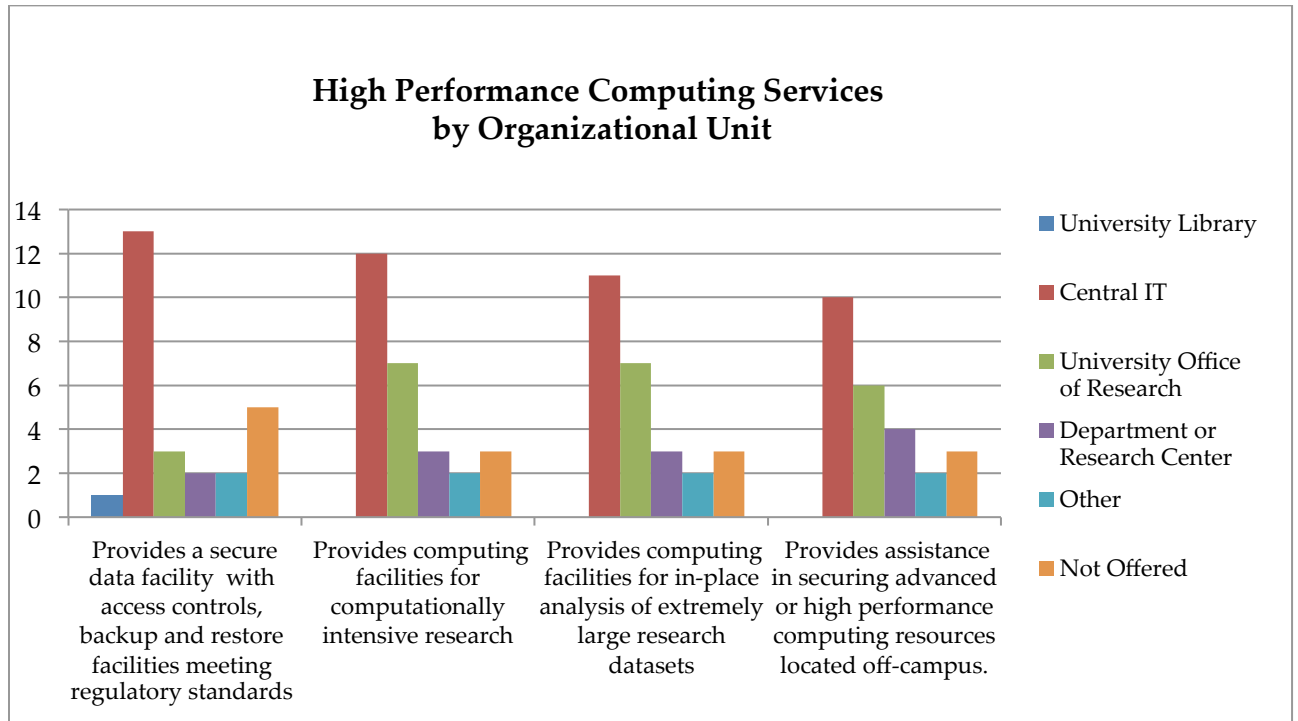
Within the sub-section of the category of accessing and using data, the question related to linking and assignment of unique identifiers to researchers and research data sets, there was low response rate from the survey respondents. Respondents did indicate that 8 libraries offer assistance with linking research data to publications and assigning persistent identifiers to research data sets. However, this was matched by the same number of respondents selecting “not offered,” indicating that this is a service not yet frequently offered. The trend within this survey could be indicating that libraries may be the emerging leaders with these types of support services, but significant conclusions cannot be reached from this sample.



High Performance Computing

Central IT was the dominant organizational unit for providing services related to high performance computing, which is not a surprising result because of the inherent role of IT in computing. A potential area for further research is the role of the university offices of research in high-performance computing. Except for providing secure data facility with access controls, backup and restore facilities meeting regulatory standards, offices of research had an average selection frequency of 7, compared with IT’s average frequency of 11. There were six open-ended responses to elaborate on an “other”

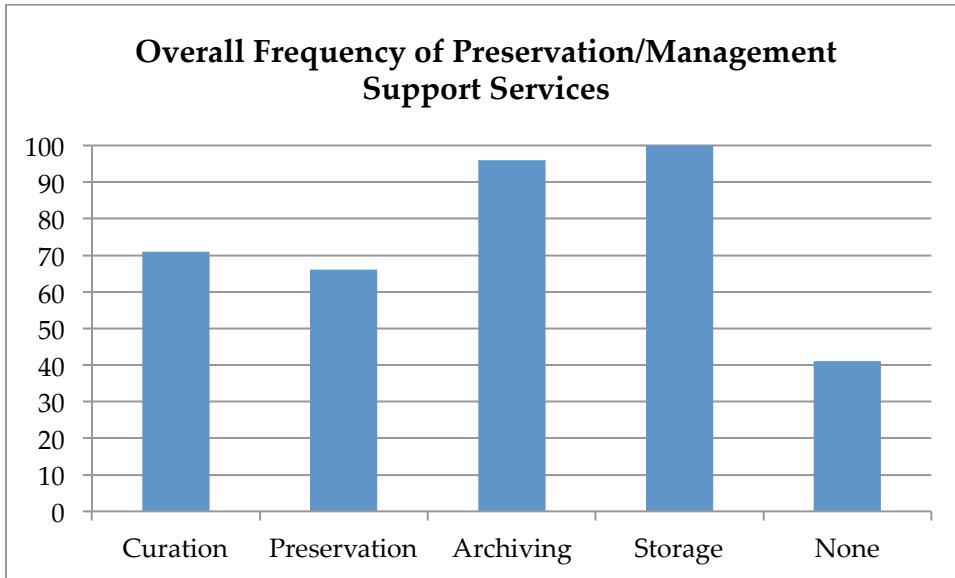
selection. These indicated that some institutions use off-campus services, including other universities high performance computing facilities and other external vendors or partners.



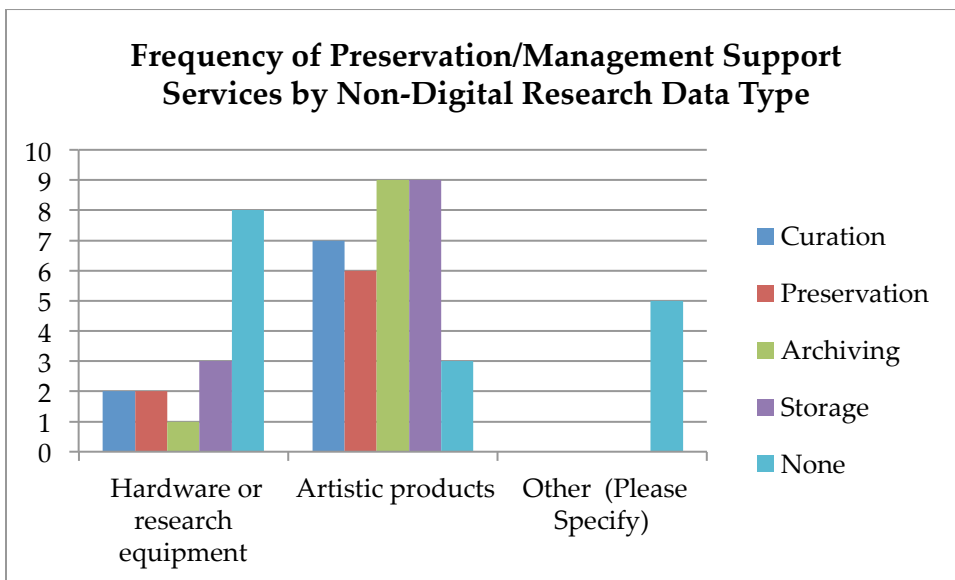
Types of Digital and Non-digital Research Data that are Supported to Manage and Preserve Long Term Institutional Access

Overall, there was high indication of services present at the responding institutions in the category of supporting long term institutional access to both digital and non-digital data. Storage was the most common support service, which is not unexpected since it is considered the most basic approach to providing access. However, archiving was the second most frequent response indicating that the responding universities are acknowledging and taking action on the necessary service of archiving data, which goes

beyond preservation. However, curation and preservation, which take many more resources, were less frequently selected.



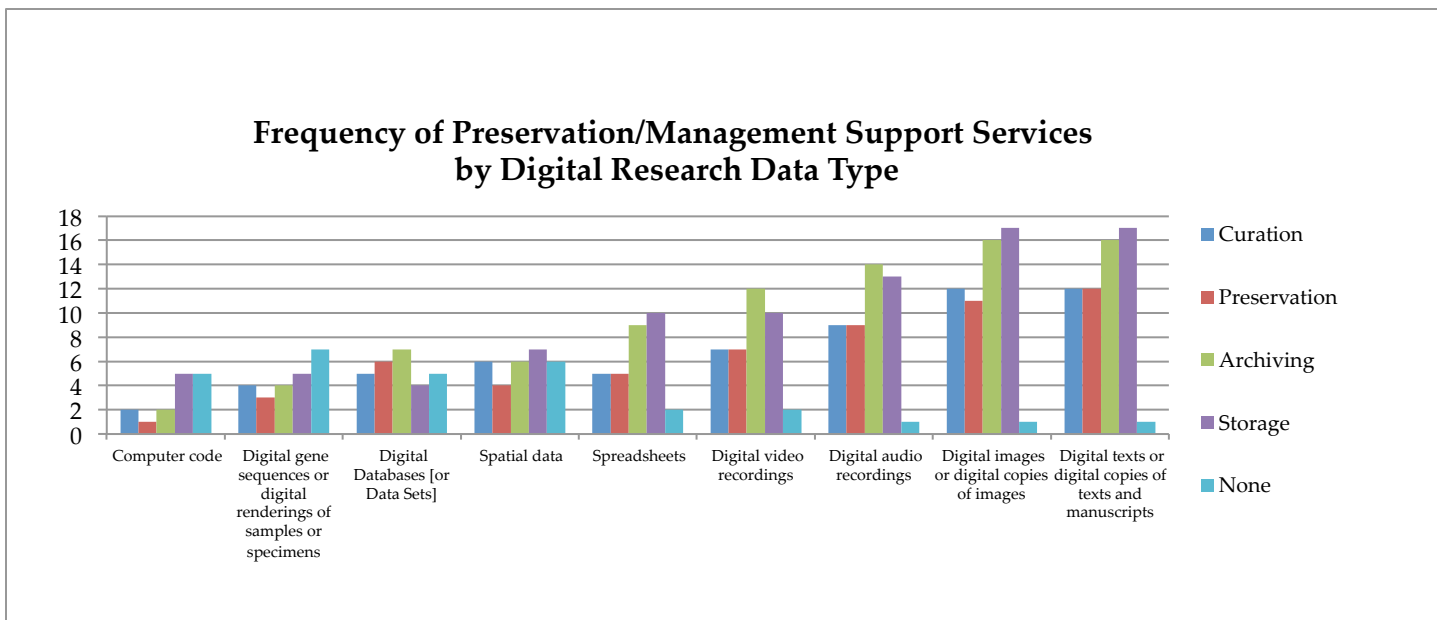
The long-term access to non-digital data had a much lower response rate than the digital data options. The area of most concern is the lack of support services for maintaining long term access to hardware and research equipment from the responding universities. With technology quickly changing and becoming obsolete, maintaining hardware may become an issue because of data that might potentially be unusable if it is not (or cannot) be put in a format for newer technologies. However, it is possible that the high response rate from library units potentially skewed this section because libraries have not traditionally played a role in hardware management and preservation.



Support services for long-term access to digital data are more prevalent, at least among the responding universities, than support services for non-digital data. The average number of total responses indicating support among all digital types (excluding “other”) was 30. The most commonly supported digital data formats are:

1. Digital texts (57)²⁵
2. Digital images (56)
3. Digital audio recordings (46)

All three of these digital formats only received one “not offered” option. Since access to the above data types is commonly supported by libraries, this is not an unexpected result. As indicated by the survey respondents, the two least supported digital data formats for long term access are computer code (10) and digital sequences gene sequences or similar digital renderings of biological/organic/inorganic samples or specimens (16). One explanation could be that access to these data types are not generally supported by the units surveyed, and therefore, support services are unknown by respondents.



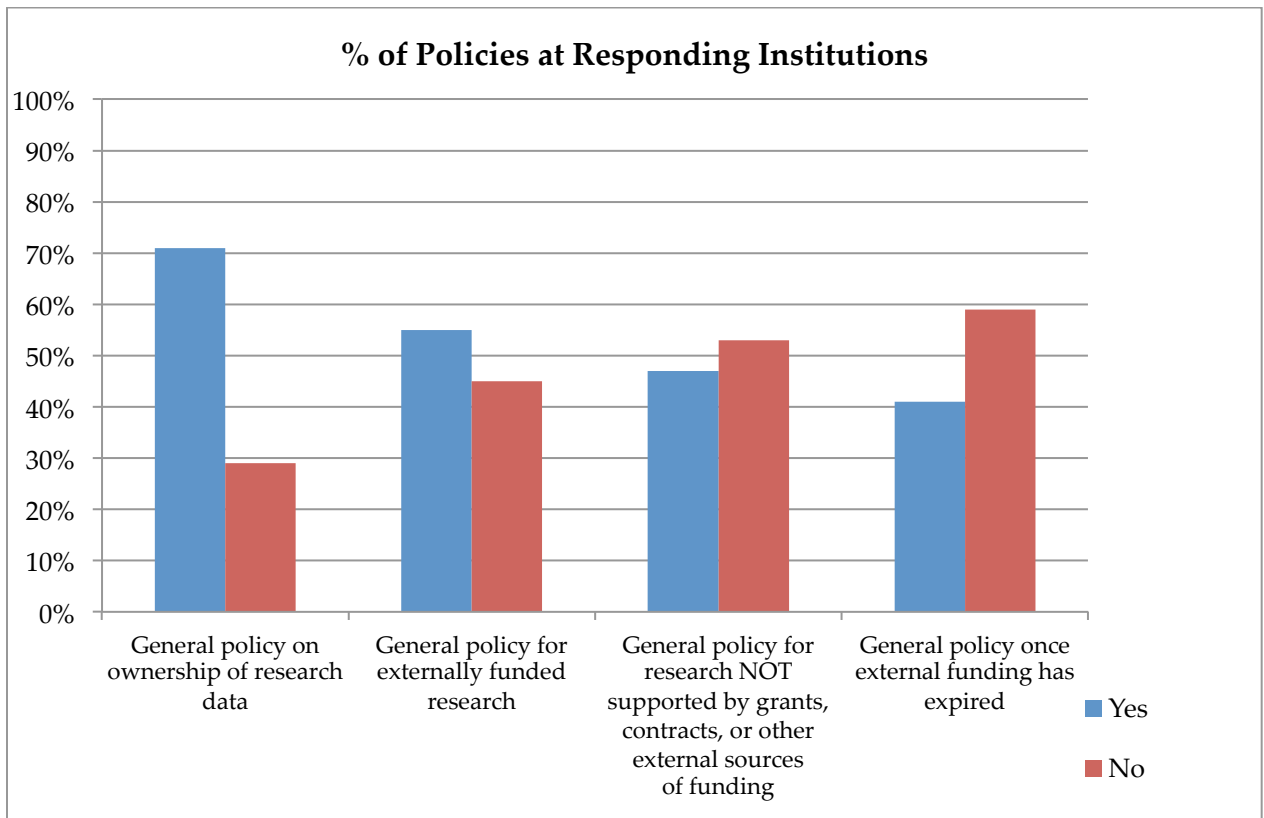
²⁵ Total responses w/out “not offered” option

Research Data Policies

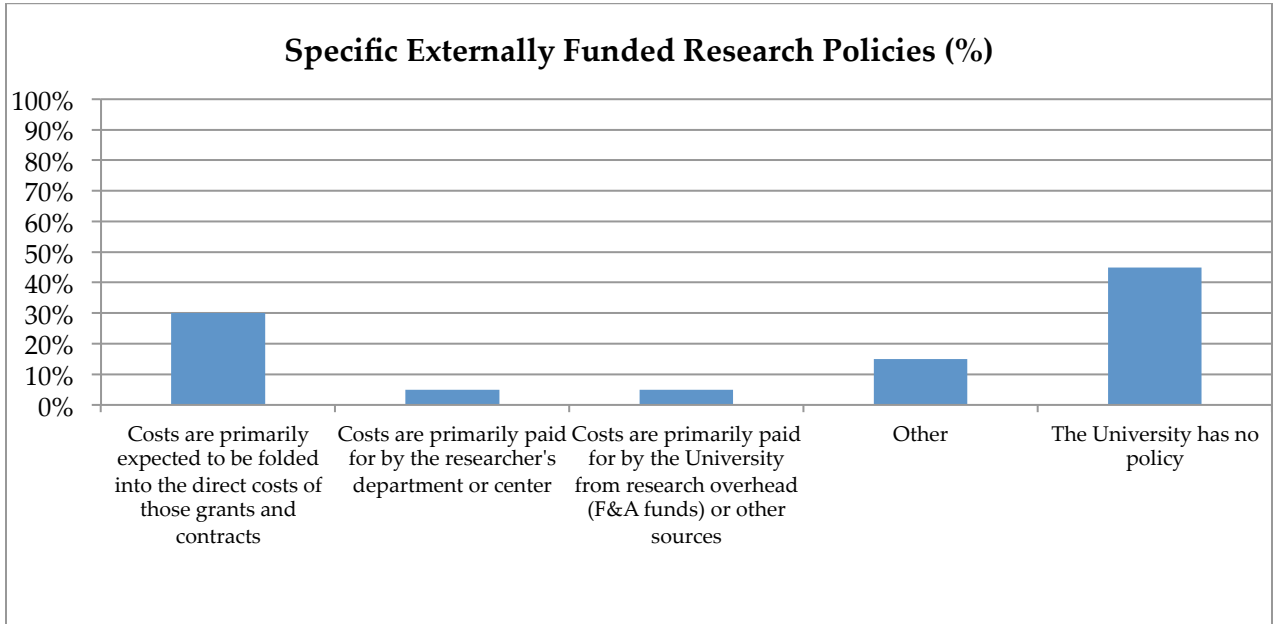
The research data policy section of the survey investigated four types of potential institutional policies:

1. General research data ownership policy
2. General policy for externally funded research
3. General policy for research NOT supported by grants, contracts, or other external sources of funding
4. General policy once external funding has expired

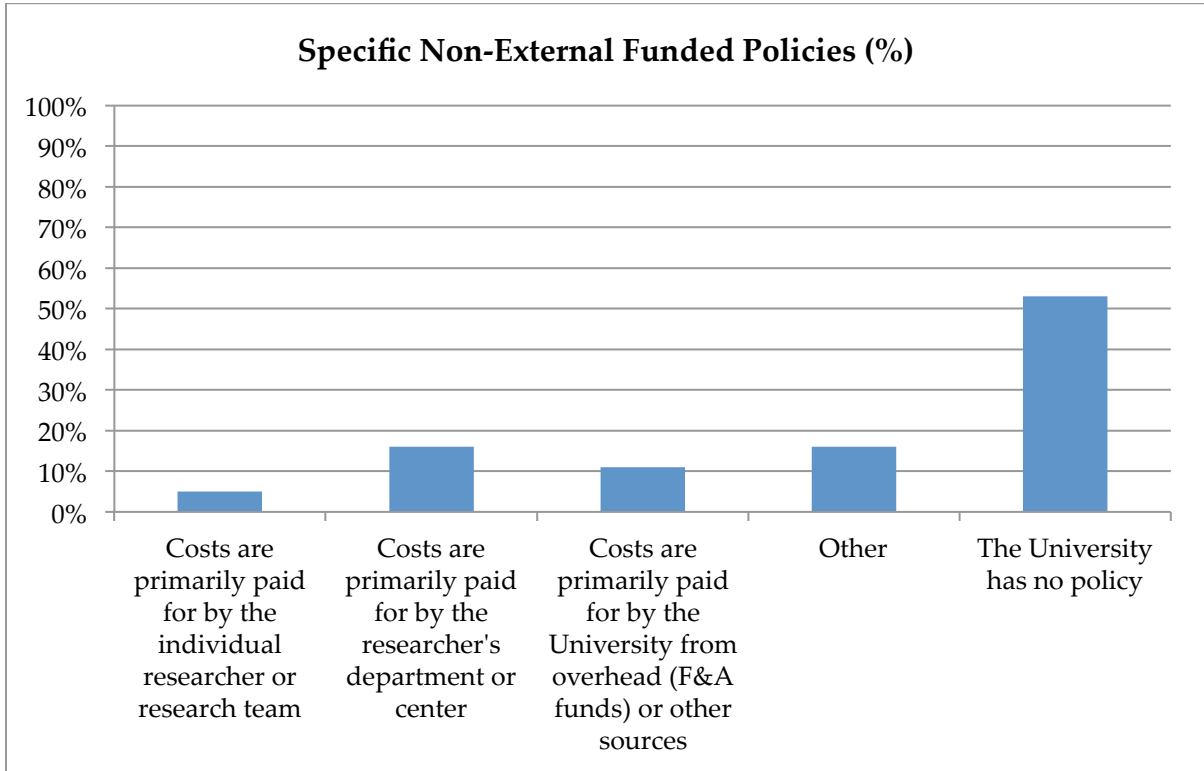
Out of 17 responses, 12 (71%) respondents indicated that their institutions have a general data ownership policy. The remaining policy areas had much closer results. A slight majority indicated that their institutions have policies on externally funded research. The remaining policy areas had slightly more negative responses, indicating areas for further research at the institutional level. However, the results of this section could once again be influenced by the internal-institutional affiliation of the respondents. Since libraries were highly represented, knowledge of specific types of funded research policies may not be known.



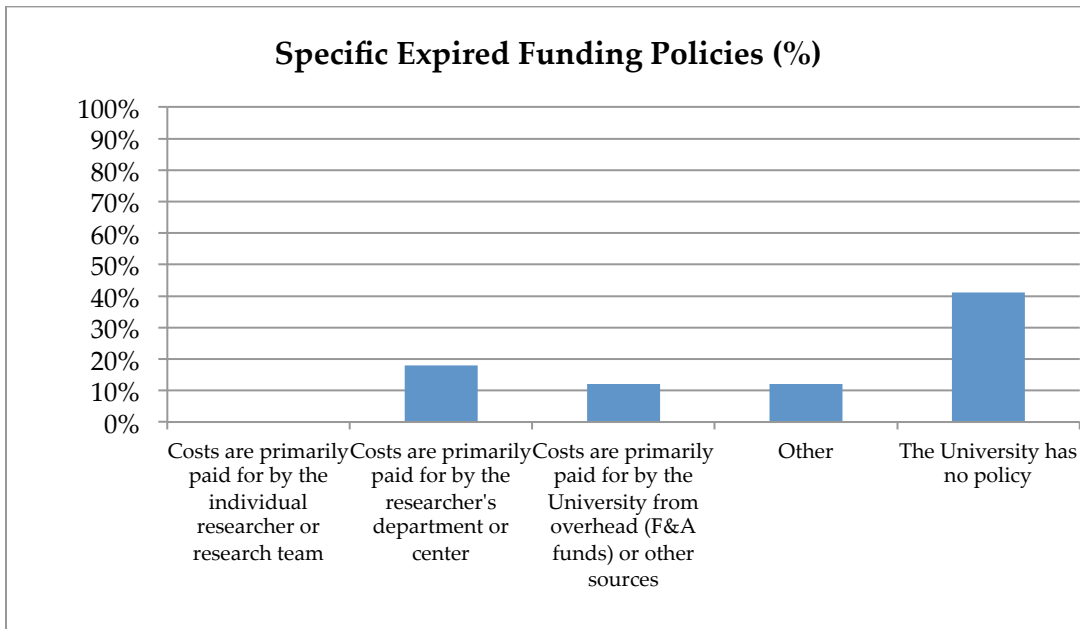
For policies related to externally funded research projects, 30% of respondents indicated that costs are primarily folded into the direct costs of the grants. According the survey results, research departments or university-level funding are rarely responsible for financial contributions. However, 45% indicated that their institutions do not have specific policies on externally funded research.



Over half of the respondents indicated that their institutions do not have a policy specifically related to non-externally funded projects (53%). There were three open-ended responses to the “other” choice. These indicated that at one institution a policy is being drafted, a second institution stated all of the choices were a part of the institution’s policy, and a third institution wrote that they will store research data from this type of project for free on a short-term basis.



For research projects that have expired, 41% of the respondents indicated that their institution has no specific policy. The next highest response was that costs were paid by the researcher's department or center. No respondent indicated that an individual researcher or research team took on the costs. There were no open-ended answers to the two "other" selections.



Section C. Review of Key Projects and Technologies

This section provides a useful reference with guide to various projects and technologies as exemplars representing five categories of data management services.

1. Curation of Data
2. Preservation of Digital Materials
3. Archiving & Repository Services
4. Storage Systems
5. Enabling Technologies

1. Curation of Data

Data curation is the active and ongoing management of data through its lifecycle of interest and usefulness to scholarship, science, and education. Data curation enables data discovery and retrieval, maintains data quality, adds value, and provides for re-use over time through activities including authentication, archiving, management, preservation, and representation. ²⁶

In this section, we share collaborative or federated efforts to curate data including:

1. Australian National Data Service (ANDS)
2. Data Conservancy
3. DataONE
4. Dataverse
5. Digital Curation Centre
6. Dryad
7. The Interuniversity Consortium for Political and Social Research (ICPSR)
8. Texas Advanced Computing Center TACC
9. The Digital Archaeological Record (tDAR)
10. UK Data Archive

²⁶ Graduate School of Library and Information Science, The iSchool at Illinois.
http://www.lis.illinois.edu/academics/programs/specializations/data_curation

Deborah Ludwig

<http://www.ands.org.au/>

Brief Description of the Project

The Australian National Data Service (ANDS) is a project arising out of the need for platforms for collaboration under a plan by the Australian National Collaborative Research Infrastructure Strategy (NCRIS). It is one of several discrete services that addresses collaboration.²⁷ It is lead by Monash University.

According to the website:

ANDS aspires to build Australia's Research Commons, as a "cohesive collection of research resources from all research institutions, to make better use of Australia's research data outputs." The data commons component registers data, which resides in repositories across a network of institutions.

ANDS is transforming Australia's research data environment to:

- *make Australian research data collections more valuable by managing, connecting, enabling discovery and supporting the reuse of this data*
- *enable richer research, more accountable research; more efficient use of research data;*
- *and improved provision of data to support policy development*

Monash University²⁸ leads the ANDS partnership with the Australian National University²⁹ (ANU) and the Commonwealth Scientific and Industrial Research Organisation (CSIRO)³⁰.

The scope of ANDS as a collaborative effort is extremely broad. ANDS includes a rich mix of public sector partners such as the Australian Bureau of Statistics, the National Archives of Australia, the Australian Institute of Health and Welfare working alongside university partners and NCRIS organizations. There are several communities of practice for metadata, software infrastructure and tools, data managers, public sector, and specialist communities. The lengthy list of projects is on the ANDS website.³¹

²⁷ Treloar, Andrew. "Design and Implementation of the Australian National Data Service." *International Journal of Digital Curation*. 4, no. 1 (2009): 125-137.

²⁸ <http://www.monash.edu.au/>

²⁹ <http://www.anu.edu.au/>

³⁰ <http://www.csiro.au/>

³¹ <https://projects.ands.org.au/getAllProjects.php?start=all>

Specific services of ANDS include:

- *Cite My Data* for assignment of persistent Digital Object Identifiers
- *Publish My Data* for registering descriptions of data collections. ANDS does not store the data itself.
- *Register My Data* is a more complex metadata registry for data collections which also allows metadata to be exposed and included in other discovery services
- *Identify My Data* is a service to attach Handles (Corporation for National Research Initiatives or CRNI service) to data.

Reviewer's Analysis

ANDS is a project to create a data commons by registering metadata for data held in various repositories at a wide group of institutions. It seems similar to DataONE in terms of providing a metadata registry with links to distributed nodes. A difference is in scope with DataONE being focused on earth science data and ANDS being focused on many disparate types of research data. More would need to be discovered about local or institutional data sources and expectations for the treatment of those data sources for permanence.

Considerations and Recommendations

If GWLA and GPN member institutions are interested in creating a research commons approach to create more knowledge about available data, this would be an interesting model to further digest. Note that it was built initially on three years of external funding.

Key Contacts: contact@ands.org.au

Sponsors: Australian Commonwealth Department of Education, Science, and Training (DEST).

Funding: 3 years of initial funding to develop the project. Ongoing is not clear from sources consulted.

Inception: 2007

Geographic Location: Australia

Data Conservancy

CURATION

<http://dataconservancy.org/>

Philip Konomos, Arizona State University

Brief Description of the Project

The Data Conservancy (DC) is an initiative to support the preservation, management, and re-use of scientific data, particularly data that results from grant funded research projects. The project started as a collaboration between the Sheridan Libraries at Johns Hopkins University and the Sloan Digital Sky Survey (SDSS). Sheridan was tasked with curating the resulting astronomical data from the SDSS and developed a data curation framework to complete the work.

Over time, using funding from the 2007 NSF DataNet solicitation as well as the Institute for Museum and Library Services, the DC has expanded its focus to four areas of science information research:

- conducting ongoing needs assessment among science researchers to develop new digital data curation tools and services;
- researching and developing the cyberinfrastructure needed to curate, preserve and make science data accessible;
- working with Library and Information Schools on data curation education and professional development; and,
- examining strategies for the long-term sustainability of repositories and data curation centers (business model).

Reviewer's Analysis

The Data Conservancy is still in development, but aims to include an open, digital repository of science data as well as tools and services to enhance the re-use of data. One example of a research tool is the Feature Extraction Tool - the ability to find and integrate disparate data sets in the network of DC repositories, by querying using taxonomic, spatial or temporal terms. There is no public interface allowing querying or access to data currently, however. Instead, staff and affiliated scientists are working on analyzing the needs of the science communities and building appropriate curation and analysis tools.

Considerations and Recommendations

DC is actively seeking partner institutions to download and install an "instance" of their software. For GWLA members who have strengths in astronomy, earth science, life

sciences, or social sciences, this might be a worthwhile partnership. However, institutions would have to be willing to commit significant technical staff time to installing, configuring and developing the DC software.

Key Contacts: Sayeed Choudhury, Associate Dean for Library Digital Programs and Hodson Director of the Digital Research and Curation Center Sheridan Libraries, Johns Hopkins University

Sponsors: The Data Conservancy began at the Sheridan Libraries at Johns Hopkins University, but the current model is to create a community of partners (universities, research centers, institutions) with instances of the cyberinfrastructure and data, who help with continued research and development of new tools and services.

Funding: In 2007 NSF issued a call for proposals to create cyberinfrastructure that would curate and make accessible the increasing amounts of grant funded science research data. Two proposals were funded starting in 2008: DataONE, a repository covering ecology, evolutionary, and earth science data; and the Data Conservancy, a repository network focused on astronomy, earth science, life sciences, and social science data. Additional funding for DC has also come from the Andrew W. Mellon Foundation and the Institute of Museum and Library Services (IMLS).

Inception: 2007-2008

Geographic Location: Located at Johns Hopkins University Library, but including a network of other universities with instances of the infrastructure.

DataONE (Data Observation Network for Earth)

CURATION

<http://www.dataone.org/>

Michael Bolton, Texas A&M University, Sterling C. Evans Library

Brief Description of the Project

From the web site, DataONE is one of the original DataNets supported by the U. S. National Science Foundation grant #OCI-0830944. It provides a distributed framework and sustainable cyber infrastructure to provide open, persistent, robust and secure access to well-described and easily discovered Earth observational data. In short, DataONE provides a comprehensive set of tools and a repository for collecting metadata on datasets dealing with Earth, environmental, atmospheric and ecological sciences. The infrastructure is composed of THREE (3) coordinating nodes to provide the network-wide services. A growing number of member nodes provide data to the coordinating nodes where it is indexed and replicated across all 3 nodes. In addition to the repository service, DataONE also provides a rather long list of tools referred to as the Investigator Toolkit. These tools cover the full spectrum of data management, from preparing the data management plan to archiving metadata in the ONEShare repository. DataONE also has a strong outreach component that includes educational services as well as a User's Group that helps promote the service. DataONE is a fairly mature service, well documented and accessible. Its governance model is equally mature with an Executive Team, a Leadership Team and a fairly large number of working groups.

Reviewer's Analysis

DataONE provides a very good model for DataNet, either regional or focused on a particular discipline, such as Earth Sciences. The architecture lends itself to a consortia approach in that a series of super-nodes, Coordinating Nodes for DataONE, aggregate data for a larger number of member nodes. Multiple, geographically dispersed coordinating nodes are key to the preservation model used in DataONE.

It should be noted DataONE does not store data, just the metadata. The actual datasets reside at the originating member node. DataONE facilitates an integrated search and discovery service and also provides replication services to the member nodes. So I do not see DataONE as a real preservation network at this time, which puts a premium on deploying something like DPN [the Digital Preservation Network] for deep archiving of data.

It includes a number of nice tools, which we can elect to add to our offering. Some of these tools have been in existence for some time while others were developed as part of the DataONE grant. The use of Mercury as the search engine, along with the harvesting into a centralized index makes for a very strong search interface.

Considerations and Recommendations

This would be a good place to start in considering a model. In looking at the Digital Preservation Network model, the idea of coordinating and contributing nodes seems to be a common model. I think we should investigate how we could use the DataONE architecture for our project. The DataONE Architecture documents are online at <http://mule1.dataone.org/ArchitectureDocs-current/index.html>

Key Contacts: PI – William Michener, Executive Director – Rebecca Koskela

Sponsors: University of New Mexico

Funding: DataONE is supported by NSF grant #OCI-0830944 – part of the initial DataNets project. The grant is estimated to expire July 31, 2014.

Inception: From the NSF grant, August 1, 2009 is the start date.

Geographic Location:

Current configuration has THREE (3) Coordinating Nodes located at;

- The University of New Mexico
- The University of California Santa Barbara
- The University of Tennessee (collaboration with Oak Ridge National Laboratory)

Member nodes, of which there are 10 at this time, are all over the country.

Brief Description of the Project

The Dataverse Network is an open source web application for hosting research data and provides flexible tools for the management of institutional branding, user accounts and use/access requirements. Additional features include dynamic analysis and visualization tools, versioning, reformatting of some statistical file types, and use tracking. Depending on the original file type, some data can be converted to preservation format on upload, with DDI and Dublin Core metadata stored in XML format. Additional metadata profiles can be specified through templates. Each Dataverse Network is composed of one or more individual Dataverse, and may be distributed across multiple departments or organizations. Individual Dataverse can be defined at various levels – per institution, per researcher, per project, etc.

Reviewer's Analysis

The Dataverse Network is a feature rich application with regard to publication and post-publication data management activities. In particular, the facets of data sharing and reuse are well supported through granular access management and user comment features. Many frequently cited concerns regarding data publication are addressed - user permissions can be specified at the file level, allowing for multiple access options even within a single project. Access terms can be customized, and the guestbook features allow researchers and administrators to require information from users prior to authorizing data download. Organizations running a Dataverse Network can enable online analysis and visualization by installing an optional R server package, allowing users to interact with the data for validation or selection purposes. Formatted data citations and handles are provided. Finally, archival and preservation processes are supported through scheduled XML metadata exports as well as LOCKSS and OAI harvesting utilities. (Metadata profiles for the physical sciences and humanities can be specified in addition to the default DDI.)

Other phases dealing with data creation and analysis workflow are less well represented. User comments and versioning features apply to static, posted data sets, but robust collaboration and version control tools would have to be managed by other means.

With regard to implications for a collaborative or consortial approach to research data management, the flexible networking and granular access controls are significant. Making use of these features, a lead institution could centrally manage the Dataverse Network application while contributing institutions remotely manage their individual Dataverses. Use and access policies can be configured locally, so a one-size-fits-all approach would not be required unless agreed upon by consortia members.

Considerations and Recommendations

The project does merit further review and analysis, as it is under active development and has already achieved a sizable user base within the social sciences. Additionally, a Dataverse Network plugin for the Open Journal System is in development, which may broaden the adoption of this resource.

Key Contacts: Sonia Barbosa, Eleni Castro

Sponsors: Institute for Qualitative Social Science at Harvard University, Massachusetts Institute of Technology

Inception Date: Development began in 2006

Brief Description of the Project

From the web page:

The Digital Curation Centre (DCC) is a world-leading centre of expertise in digital information curation with a focus on building capacity, capability and skills for research data management across the UK's higher education research community.

DCC provides expert advice, training, support and consultancies for the UK higher education's research community.

The U.S. works with DCC in international organizations, for example in the Committee on Data for Science and Technology CODATA (<http://www.codata.org>). The U.S. National Research Council's Board on Research Data and Information represents US in CODATA. Francine Berman and Clifford Lynch among members of US/CODATA.

Reviewer's Analysis

DCC is a comprehensive reference source on research data curation. The Center website offers well-structured information, resources, and tools to everyone interested in data curation: briefing papers, how-to guides, curation reference manual, curation lifecycle model, policy and legal resources, data management plans, sets of tools including Collaborative Assessment of Research Data Infrastructure and Objectives CARDIO (<http://cardio.dcc.ac.uk/>), Data Asset Framework DAF (<http://www.data-audit.eu/index.html>), Digital Repository Audit Method Based on Risk Assessment DRAMBORA Interactive toolkit (<http://www.repositoryaudit.eu/>); case studies, standards, training courses, and research and development resources.

DCC is one of major reference sites, so there is always something useful, including training materials; standards watch; tools, programs, and reports:

- Collaborative Assessment of Research Data Infrastructure and Objectives (CARDIO) tool can be useful. <http://www.dcc.ac.uk/projects/cardio>
- Neil Beagrie, JISC Benefits from the Infrastructure Projects in the JISC Managing Research Data Programme. Final Report, version 5.0, Sept. 2011

- http://www.jisc.ac.uk/media/documents/programmes/mrd/RDM_Benefits_Final_Report-Sept.pdf
- Disciplinary metadata standards: <http://www.dcc.ac.uk/resources/metadata-standards>
 - Catalog of resources for curators and researchers: <http://www.dcc.ac.uk/resources/external/tools-services>

Considerations and Recommendations

The DCC Web Site can be included to the list of useful resources. DDC training materials and tools can be used especially when similar U.S. materials have not yet developed. In more direct approach, specific training and assessment materials, tools, and recommendations may be selected and adapted as needed during planning and implementation of the GWLA/GPN initiative.

Key Contacts: Kevin Ashley, Director kevin.ashley@ed.ac.uk +44 131 651 3823; Sarah Jones, Senior Institutional Support Officer sarah.jones@glasgow.ac.uk;
Twitter: sjDCC Phone: +44 141 330 3549

Sponsors: HATII, UKOLN and STFC, with the University of Edinburgh (from March 2004 through February 2010). From March 2010 the DCC has reorganized into a three-cornered consortium, led from Edinburgh, with the following Principal Partners: Humanities Advanced Technology and Information Institute (HATII) at the University of Glasgow and UKOLN, a center of expertise in digital information management, based at the University of Bath.

Funding: JISC

Inception: 2004-03-01

Geographic Location: Edinburgh

Dryad

<http://datadryad.org/>

Susan Matveyeva, Wichita State University

CURATION

Brief Description of the Project

From the Dryad's web page at NCSU:

Dryad is a joint project of NESCent and the UNC Metadata Research Center, with North Carolina State University participating as a development partner, along with the University of New Mexico, and Yale University, focusing on creation of a repository for data underlying scientific publications, with an initial focus on evolution, ecology, and related fields. Dryad allows investigators to validate published findings, explore new analysis methodologies, repurpose the data for research questions unanticipated by the original authors, and perform synthetic studies such as formal meta-analyses."

Reviewer's Analysis

Dryad addresses all key facets of data lifecycle management. It is focused on sustainability, preservation, reuse of data, and automation of data flow and exchange. Developers continue adding new functionality. The Dryad UK Project (2010-2011) created a mirror of the Dryad repository and developed a sustainability plan to help Dryad become established as an international not-for-profit organization empowered to ensure long-term preservation and accessibility of its data holdings.

Software: DSpace Manakin XMLUI customized by @mire (<http://atmire.com/website/>) with the following features: dataset embargo; dataset security; Discovery SOLR Search System customization; enhanced submission and workflow; configurable workflow; item versioning; integration with EZID, DOI, DataCite and PubMed. Some customizations are now standard features of DSpace.

Search: SOLR Discovery faceted browse and search. Includes indexes:

- authority -- terms for auto-completion using controlled vocabularies, including HIVE
- dataoneMNlog -- log of accesses through the DataONE API
- dryad -- local storage of DOIs
- search -- primary search index
- statistics -- log of accesses to item pages and bitstream downloads

Access to Data: Authors release their data in public domain under the terms of a Creative Commons Zero (CC0) waiver to emphasized data sharing and increase data reuse. Embargo is an available option.

Content submission: Data is submitted as part of the publication process. “Journal integration with Dryad is available at no cost for any journal that wishes to implement low-burden data archiving and enhance their published articles with links to data.” No restrictions on file format. Data packages include data files and “ReadMe” file. Each data package is assigned DOI and linked to the journal article. Submission is simple: (see 2-min. video: <http://www.youtube.com/watch?v=RP33cl8tL28&feature=youtu.be>)

Content sharing: Dryad uses EZID service from the California Digital Library that manages DOI registration with DataCite. Dryad is a node of DataONE.

Reusability: LinkOut functionality. PMID in Dryad metadata. Citation sharing technology (Cite and Share in RIS and BibTex). Handshaking (the process of coordinating submission between Dryad and specialized repositories in order to (a) lower user burden by streamlining the submission workflow and (b) allow Dryad and specialized repositories to exchange identifiers and other metadata in order to enable cross-referencing of the different data products associated with a given publication (integration with GenBank and TreeBASE).

Integration: Integration with ORCID/DataCite is in a planning stage.

Metadata: Metadata profile: Dublin Core with addition of few fields from DDI, Darwin Core and PRISM (see http://www.cendi.gov/presentations/11-17-09_cendi_nfais_Greenberg_UNC.pdf). Metadata generation is automatic or semi-automatic.

Dryad implemented an automated integrated submission from multiple publishers’ sites: http://wiki.datadryad.org/Submission_Integration. The permanent link between the article and data packages is created during submission. A similar approach can be used to link records and/or full text/ or data in the local repository/storage and the consortium’s central hub.

Workflow: Dryad’s workflow practices of assigning DOI to data packages is another process to look at during the planning step of the regional consortium development.

Considerations and Recommendations

I like Dryad's careful approach to the enhancement of qualified Dublin Core records with fields from other metadata. These enhancements do not destroy the integrity of a Dublin Core record; inclusion of the elements of other schemas is as minimal as possible. Only few necessary elements are included: scientific name of a plant (Darwin Core dwc:ScientificName); journal name (Prism prism:publicationName). I would recommend similar approach for multidisciplinary repository: have major schema (e.g. Dublin Core) and few fields from the metadata schema of the appropriate discipline (if it was developed).

I recommend looking more closely at Dryad's best practices in metadata automation, use of identifiers, integration with multiple publisher's platform during submission, data packages sharing with other repositories (<http://wiki.datadryad.org/BagIt> Handshaking), versioning, curation practices and reports, experience of partnership with DataONE, DataCite, CLOCKSS, ORCID. Dryad's site has good documentation: http://wiki.datadryad.org/Main_Page See also Curation Manual: http://wiki.datadryad.org/wg/dryad/images/8/85/Curation_man_2012-12-21.pdf The Dryad practices analysis and usage would be especially useful if the committee choose DSpace as a software platform.

Key Contacts: Jane Greenberg, metadata management (Metadata Research Center; University of North Carolina at Chapel Hill) Phone: 919-962-8066 ; Fax: 919-962-8071; Email: janeg@email.unc.edu

Hilmar Lapp, technical management (NESCent); Todd Vision, project director (NESCent/University of North Carolina at Chapel Hill)

Sponsors: Dryad developed by the [National Evolutionary Synthesis Center](#) and the University of North Carolina [Metadata Research Center](#), in collaboration with several [Partner Organizations](#)..

U.K. development partners: JISC (the Joint Information Systems Committee), Oxford University, and the British Library (see the DDC project DryadUK:

<http://www.dcc.ac.uk/projects/dryaduk>)

Dryad is governed by a twelve member Board of Directors elected by members, representing publishers, societies, research and educational institutions (

http://wiki.datadryad.org/wg/dryad/images/9/94/Dryad_ByLaws_April2012.pdf

Funding: NSF grant DBI-0743720 (2008-2012);
NSF grant 2012-2016 (Abstract [DBI-1147166](#));
NESCent, the NSF-funded DataONE;
IMLS grant (LG-07-08-0120-08)

Business plan and Sustainability: Dryad is currently applying for status as a 501(c)3 not-for-profit to be incorporated in North Carolina. It also plans to charge membership and submission fees

(http://wiki.datadryad.org/Business_Plan_and_Sustainability). Staff will continue applying for R&D grants.

Inception: 2008

Geographic Location: hosted by the North Carolina State University

**Inter-university Consortium for Political and Social Research
(ICPSR)**

CURATION

<http://www.icpsr.umich.edu>

Deborah Ludwig

Brief Description

The Inter-university Consortium for Political and Social Research is a consortium of over 700 institutions. Membership includes research institutions, colleges and universities. ICPSR's mission, according to the website is to provide "leadership and training in data access, curation, and methods of analysis for a diverse and expanding social science research community. "

The website notes that "ICPSR maintains a **data archive** of more than 500,000 files of research in the social sciences. It hosts 16 specialized collections of data in education, aging, criminal justice, substance abuse, terrorism, and other fields." In addition to acting as a data archive, ICPSR includes a focus on training and education for data professionals and researchers.

Unlike some models that seek to unify data collections held in disparate locations, ICPSR is a centralized data archive. Data are deposited in the ICPSR archive and held there with an advanced emphasis on long-term preservation.

Reviewer's Analysis

ICSPR celebrated its 50th anniversary this year, so they have delved into issues such as metadata and data preservation at a level that not many other data archives have approached. They utilize a subscription based membership model. They focus on the social science domain and are a key disciplinary repository.

Considerations and Recommendations

If a centralized repository approach were envisioned for GWLA and GPN, a closer look at their operational and business model could inform our work and help us develop our success criteria. Noting their strategies for preservation could also inform our work.

Key Contacts: <http://www.icpsr.umich.edu/icpsrweb/content/membership/contact-people.html>

Sponsors: University of Michigan

Funding: Membership-based model

Inception: 1962

Brief Description

The Texas Advanced Computing Center is supported by the National Science Foundation, the University of Texas, Austin, and through other grants. It is one of 11 national centers for advanced computing to support computational research. Project partners include the Department of Energy, the National Oceanic and Atmospheric Administration (NOAA), and the National Archives and Records Administration (NARA).

The website states:

Computational science has become the third pillar of scientific discovery, complementing theory and physical experimentation, allowing scientists to explore phenomena that are too big, small, fast, or dangerous to investigate in the laboratory. Thousands of researchers each year use the computing resources available at TACC to forecast weather and environmental disasters such as the BP oil spill, produce whole-Earth simulations of plate tectonics, and perform other important research.

Within TACC, there is a Data Management and Collections Group (DMC)³² that was established in 2008. The DMC group “builds and maintains large data-management and storage resources and consults with collections’ creators in all aspects of the data lifecycle, from creation to long-term preservation and access. The DMC group actively seeks out research and grant proposal collaborations with researchers and institutions with collections of interest.” Data management is built on the iRODS platform. A recent presentation at the Preservation and Archiving Special Interest Group³³ by TACC’s Dan Stanzone states that there are 20 true data “collections” mixed in with “user storage” for data and projects.

TACC has two funding models: pay annually or pay once, store forever (which works well for researchers on grant funding.) Forever is about 5 years for most projects.

Data collections are hosted in a system called Corral, which includes storage for structure and unstructured data, lots of different interfaces to the data and iRODS-based data management services underlying the collections. Corral had over 700 TB of data as of 2012. The Corral services include hierarchical storage with multiple disk speeds,

³² <http://www.tacc.utexas.edu/tacc-projects/dmc>

³³ PASIG 2012, Austin <http://lib.stanford.edu/preservation-archiving-special-interest-group/presentations-pasig-meetings-austin-texas-january-2012>

multiple encryption and data mechanisms for HIPPA and other secure or confidential data. Preservation of collections is also under development. TACC Deputy Director Dan Stanzione is part of The Digital Preservation Network (DPN) leadership group.

Reviewer's Analysis

Corral data collections as implemented by TACC offers a grand view of what can be done at the intersection of high performance computing, research, and building collections of reusable research data.

Consideration and Recommendations

While building the level of infrastructure housed at TACC might be out of range for GWLA and GPN, perhaps there are opportunities to work with one or two high performance computing centers on a regional basis as part of a next-phase grant project. This should be explored.

Key Contacts: Dan Stanzione, dan@tacc.utexas.edu

<http://www.tacc.utexas.edu/staff/dan-stanzione>

Sponsors & Funding: NSF & University of Texas

Inception: 2001

Geographic Location: Austin Texas

tDAR – (the Digital Archaeological Record)

CURATION

<http://tdar.org>

Philip Konomos, Arizona State University.

Brief Description of the Project

tDAR, the Digital Archaeological Record, is an international repository for digital archaeological data, reports, and other files. Developed by faculty and staff at Arizona State University, tDAR's mission is to make accessible, and preserve in perpetuity, the digital files (data as well as text) that are generated as a result of archaeological investigations. Physical materials recovered during archaeological projects (artifacts, organic material, paper forms, etc.) are typically curated in museums but digital records from these projects are vulnerable to loss or neglect. Often the CD or DVD media that contain the digital files are treated as additional "artifacts" from the project and boxed with the pottery and stone tools. tDAR provides a stable preservation environment that makes these digital files discoverable and accessible.

Reviewer's Analysis

This is one of several archaeological disciplinary repositories that have been developed over the past decade. One strength of tDAR is that it accepts a wide variety of file types (GIS data, tabular data sets, spreadsheets, text files, images, 3D scans, etc.) from any user (university faculty member, Federal agency archaeologist, state historic preservation officer, etc.).

Users can browse and search tDAR for free, without creating an account. Downloading files requires an account, which is free and simply involves providing a name, email address, user ID and password. Fees are charged to upload files to the repository.

tDAR is primarily an open repository but a variety of accommodations can be made for different user needs. Sensitive information in tDAR (such as exact archaeological site locations) can be redacted from reports and specific files can be marked confidential



Note: 1 Geographic Search Screen

with access limited to specified individuals only. Metadata fields are extensive and customized to support the professional vocabulary common to archaeology. Users can search by drawing a box on a map as well as by using key words related to culture area, time period, material type, or archaeological site type :

Investigation Type(s)

Inherit values from parent project

<input type="checkbox"/> Archaeological Overview	<input type="checkbox"/> Architectural Documentation
<input type="checkbox"/> Architectural Survey	<input type="checkbox"/> Bioarchaeological Research
<input type="checkbox"/> Collections Research	<input type="checkbox"/> Consultation
<input type="checkbox"/> Data Recovery / Excavation	<input type="checkbox"/> Environment Research
<input type="checkbox"/> Ethnographic Research	<input type="checkbox"/> Ethnohistoric Research
<input type="checkbox"/> Geophysical Survey	<input type="checkbox"/> Ground Disturbance Monitoring
<input type="checkbox"/> Heritage Management	<input type="checkbox"/> Historic Background Research
<input type="checkbox"/> Methodology, Theory, or Synthesis	<input type="checkbox"/> Reconnaissance / Survey
<input type="checkbox"/> Records Search / Inventory Checking	<input type="checkbox"/> Remote Sensing
<input type="checkbox"/> Research Design / Data Recovery Plan	<input type="checkbox"/> Site Evaluation / Testing
<input type="checkbox"/> Site Stabilization	<input type="checkbox"/> Site Stewardship Monitoring
<input type="checkbox"/> Systematic Survey	

Material Type(s)

Inherit values from parent project

<input type="checkbox"/> Basketry	<input type="checkbox"/> Building Materials	<input type="checkbox"/> Ceramic
<input type="checkbox"/> Chipped Stone	<input type="checkbox"/> Dating Sample	<input type="checkbox"/> Fauna
<input type="checkbox"/> Fire Cracked Rock	<input type="checkbox"/> Glass	<input type="checkbox"/> Ground Stone
<input type="checkbox"/> Hide	<input type="checkbox"/> Human Remains	<input type="checkbox"/> Macrobotanical
<input type="checkbox"/> Metal	<input type="checkbox"/> Mineral	<input type="checkbox"/> Pollen
<input type="checkbox"/> Shell	<input type="checkbox"/> Textile	<input type="checkbox"/> Wood

Cultural Term(s)

Inherit values from parent project

Culture

- Pre-Clovis
- Paleolndian
- Archaic
- Hopewell
- Woodland
- Plains Village
- Mississippian
- Ancestral Puebloan
- Hohokam
- Mogollon
- Patayan
- Fremont
- Historic

Other

[+ add another cultural term](#)

Note: 2 Other metadata fields.

tDAR also provides a database integration tool for advanced research. Users looking for specific information (show me all of the database tables that contain information on

fish remains) can query the database files in tDAR and get results that only include the rows and columns specified in the query.

Considerations and Recommendations

tDAR contains a large, varied collection of text and data files useful for teaching and research. Staff are developing example curriculum modules for use in university classes. tDAR provides a helpful research tool for undergraduate and graduate students beginning research on a particular archaeological topic or geographic area. Librarians, particularly anthropology subject specialists, will find this a useful resource.

Key Contacts: The tDAR email is comments@tdar.org . The Digital Antiquity email is info@digitalantiquity.org . Francis Pierce-McManamon, Executive Director, Center for Digital Antiquity
Arizona State University
PO Box 872402
Tempe, AZ, 85287-2402

Adam Brin, Director of Technology, Center for Digital Antiquity
Arizona State University
PO Box 872402
Tempe, AZ, 85287-2402

Sponsors: [tDAR](#) is managed by the non-profit organization Digital Antiquity, at Arizona State University. Digital Antiquity is [run by a staff of five](#), and is governed by a [Board of Directors](#) and a [Professional Advisory Panel](#).

Funding: tDAR, through Arizona State University faculty and staff, has received over 4.3 million dollars in support from the National Science Foundation, the Andrew W. Mellon Foundation, the National Endowment for the Humanities, and the UK-based Joint Information Systems Committee (JISC).

Inception: Software development began in 2004 and the public website (supporting ingest, search, browsing and download) was launched in 2009.

Geographic Location: Offices are located at Arizona State University, but the repository contains international archaeological data.

UK Data Archive

CURATION

<http://data-archive.ac.uk>

Greg Monaco, Great Plains Network (GPN)

Brief Description of the Project

The UK Data Archive is an all-in-one portal to create and store/deposit, manage, and find shared data. This project *curates the UK's largest collection of digital social and economic research data, including data from government departments, researchers and research institutions, public organisations and companies.* The archive may be searched using the *Economic and Social Data Service* website at <http://www.esds.ac.uk/Lucene/Search.aspx>. The steps of this project's data curation process are located at <http://data-archive.ac.uk/curate/process>.

Reviewer's Analysis

The project provides a one-stop portal for all aspects of the lifecycle management process for a slice of research data: "Our collection encompasses a significant range of data relating to society, both historical and contemporary, covering the social sciences, economics and humanities, as well as the societal aspects of environmental and medical data." Rather than distributed data, it appears that all data is centrally located and centrally curated.

This is an example of a complete approach to lifecycle management, and I suggest that we explore what works and what the limits to this approach may be. (One question is whether this is scalable to other domains of knowledge/data?)

Considerations and Recommendations

This appears to be a high visibility project. I recommend that we do a further review of this project in order to address similarities and differences in approach between this project and our ultimate proposed project.

This project, in association with JISC, uses federated access management, likely Shibboleth (see Shibboleth review). We should probably determine this, for sure. (<http://data-archive.ac.uk/about/projects/identity-management>)

Key Contacts:

Sponsors: University of Essex, JISC, European Commission

Funding: <http://data-archive.ac.uk/about/projects/past>

Inception: Website states that it was established over 40 years ago.

Geographic Location: United Kingdom

2. Preservation of Digital Materials

For the purpose of this report , we consider the practice of preserving digital research data to be aligned with the preservation of other types of digital information and refer to “digital preservation” rather than “data preservation.”

*Digital preservation combines policies, strategies and actions to ensure access to reformatted and born digital content regardless of the challenges of media failure and technological change. The goal of digital preservation is the accurate rendering of authenticated content over time.*³⁴

EXAMPLES OF DIGITAL PRESERVATION SERVICES AND FRAMEWORKS

1. Archivemata
2. Dark Archive in the Sunshine State(DAITSS)
3. DuraCloud
4. Digital Preservation Network (DPN)
5. Ex Libris’ Rosetta
6. OAIS Reference Model for Preservation Services
7. TRAC Certification Process

³⁴ Definition prepared by the ALCTS Preservation and Reformatting Section, Working Group on Defining Digital Preservation. ALA Annual Conference, Washington, D.C., June 24, 2007.
<http://www.ala.org/alcts/resources/preserv/defdigpres0408>

Archivematica

<https://www.archivematica.org>

Nicole Potter and Deborah Ludwig

PRESERVATION

Brief Description of the Project

Archivematica is free, open-source Linux/MySQL software for digital preservation system. It is installed locally on an institution's servers. Archivematica provides a web-accessible dashboard that allows users to ingest digital objects. It also allows for administrative operations so that the user can modify preservation plans, storage locations, configuration of micro-services and user access levels. Archivematica also permits the user to enter metadata using various standards including METS, PREMIS, Dublin Core and other metadata standards.

Archivematica is compliant with the Open Archives Information System model, which is an ISO standard that defines the functions of digital preservation and uses a micro-services approach, which relies on an integrated suite of software applications to handle various tasks.

Archivematica addresses three preservation strategies: emulation, migration, and normalization, to help create a standard technology platform for institution to ensure compatibility now and in the future.

Archivematica was originally developed for the city of Vancouver, BC to store records after the Olympics. The municipal archive of the City of Vancouver has a blog containing more information about the tool.³⁵ The Library of Congress has also mentioned Archivematica in it's a blog post.³⁶

Reviewer's Analysis

Archivematica is essentially a wrapper around a set of software tools that create a pipeline to preservation for digital objects. It is a really interesting platform that continues to gain ground because it is reasonably lightweight in terms of installation and start up.

Considerations and Recommendations

If GWLA / GPN institutions are interested in a preservation service associated with any type of digital content, including research data perhaps held in a accessible open repository, Archivematica merits a closer look. It is still very early in release cycles. Alpha 0.9 was just released this fall.

³⁵ <http://opensourcearchiving.org/content/archivematica-city-vancouver-archives>

³⁶ <http://blogs.loc.gov/digitalpreservation/2012/10/archivematica-and-the-open-source-mindset-for-digital-preservation-systems/>

Key Contacts: Peter Van Garderen, President, Artefactual Systems, email: info@artefactual.com

Sponsors: The UNESCO Memory of the World's Subcommittee on Technology, the City of Vancouver Archives, the University of Alberta Libraries, the University of British Columbia Library, the Rockefeller Archive Center, Simon Fraser University Archives and Records Management, Yale University Library

Funding: NA

Inception:

Geographic Location: NA

Brief Description of the Project

DAITSS is a dark archive focused is on preserving “digital masters” and reconstituting those masters for access upon request. Access is handled either outside of the archive through other delivery mechanisms or by creation and delivery of a “dissemination information package” from the “archival information packet.” (Refer to the review of the OAIS reference model if these terms are not familiar.)

DAITSS was developed by the Florida Center for Library Automation (FCLA) for use by the Florida Digital Archive (FDA) which is a digital repository shared by the eleven universities in the Florida public university system. DAITSS is strictly modeled on the Reference Model for an Open Archival Information System (OAIS). DAITSS can accept a Submission Information Package (SIP), transform the SIP into a stored Archival Information Package (AIP), and transform the AIP into a Dissemination Information Package (DIP) on request. To do so, DAITSS directly implements four of the six OAIS functional entities: Ingest, Data Management, Archival Storage, and Access. FDA staff performs functions of the remaining two OAIS entities, Administration and Preservation Planning, with support from DAITSS reporting and data management functions. DAITSS is unique among repository applications in that it was designed to ensure the long-term render-ability of authentic digital materials. DAITSS maintains standardized preservation metadata including digital provenance, and performs continuous fixity checking on multiple stored copies. The preservation protocol implemented by DAITSS combines bit-level preservation, format normalization, and forward format migration.³⁷

Considerations and Recommendations:

How does the project address, or potentially address, key facets of lifecycle management? DAITSS is available for use through a GPLv3 license. The DAITTS website provides links to access to a fully configured VM version of DAITSS that can be downloaded to run under any VM manager, along with sample SIPs (submission packages) and documentation.

DAITSS was written for a multi-user environment and supports consortial as well as institutional preservation repositories. If a solution of interest was to provide a centralized point for multiple institutional deposit of research data coupled with a catalog for access and a

³⁷ Caplan, Priscilla. “DAITSS, an OAIS-based preservation repository” in Proceedings of the 2010 Roadmap for Digital Preservation Interoperability Framework Workshop. Article #17.

service to request research data sets and to deliver those “on demands”, DAITSS might provide a viable solution.

DAITSS could possibly be a component in a curation approach; however there may be other emerging choices to consider as well. DAITSS is certainly a quality open software solution available at no cost and built by a team that has been doing this work for a number of years for a consortium. There are several published articles that have been written about DAITSS as well.

Key contacts: Pricilla Caplin

Sponsors: Florida Center for Library Automation, Florida Digital Archive (11 universities)

Funding: Florida Institutions and IMLS

Inception: Spring 2007

Geographic Location: Florida

Digital Preservation Network (DPN)

www.dpn.org

Deborah Ludwig

PRESERVATION

Brief Description

Verbatim from the website:

The Digital Preservation Network (DPN) was formed to ensure that the complete scholarly record is preserved for future generations. DPN uses a federated approach to preservation. The higher education community has created many digital repositories to provide long-term preservation and access. By replicating multiple dark copies of these collections in diverse nodes, DPN protects against the risk of catastrophic loss due to technology, organizational or natural disasters.

Reviewer's Analysis

DPN is an effort to develop a group of preservation nodes that are not accessible except to carry out preservation functions. The idea is that content can be restored to an access repository if lost and that data replication coupled with messaging between nodes will underpin preservation efforts. Members have paid into the development of the network ahead of its actual design and implementation to further these efforts.

Recommendations

Key Contacts: Email: inquiry@dpn.org / Phone: (434) 286-3436. Steven Morales is the director of DPN.

Steven Morales, steven.morales@dpn.org

DPN Program Director

434-286-3436

Brief Description of Project

DuraCloud is a storage and preservation solution that is hosted or “in the cloud.” DuraCloud allows storage of redundant copies at multiple storage provider sites based on a simple dashboard approach. The website advertises:

...Replication and backup activities, preservation and archiving, repository backup, and multimedia access. DuraCloud also acts as a mediation layer between you and cloud storage providers, therefore eliminating the risk of vendor lock-in. ... [DuraCloud] does not address fine-grained policy and access control considerations. It can be used to house entire collections of confidential data, and/or support a system which provides granular controls, but it does not do so itself. DuraCloud does support basic authentication; and you can make spaces within DuraCloud dark or light. ³⁸

DuraCloud is a service of DuraSpace³⁹ the support and development organization for DSpace and FedoraCommons repository software in common used around the globe.

Reviewer’s Analysis

DuraCloud is an attractive option for cloud-based preservation. There are various storage options, including Amazon (and Glacier), Rackspace, and EMC. Archives can be dark or light (accessible). Media can be shared or streamed from DuraCloud if the archive is light. The Colorado Alliance of Research Libraries (Alliance) conducted a pilot in 2010 as did the BioDiversity Heritage Library, and the WGBH Media Library and Archives. ⁴⁰

Considerations and Recommendations

Consortial pricing is available. DuraCloud coupled with a discovery layer could serve some research data purposes and require limited local infrastructure. More information could be obtained from current customers or those who have undertaken pilots.

³⁸ <http://www.duracloud.org/faq>

³⁹ <http://www.duraspace.org/>

⁴⁰ <http://www.digitalpreservation.gov/meetings/documents/ndiipp10/NDIPP-DuraCloud-Panel-Final.pptx>

Key Contacts: info@duracloud.org; <http://www.duracloud.org/contact>

Sponsors: DuraSpace

Funding: not-for-profit

Inception: October 2009, first pilot projects

Geographic Location: N/A

The Digital Preservation Network (DPN)

<http://www.dpn.org>

Deborah Ludwig

PRESERVATION

Brief Description of Project

The Digital Preservation Network (DPN) is building a proof-of-concept system to demonstrate digital preservation of the scholarly assets of higher education. There is no existing system yet. Members are contributing funding toward development. There is a data partnerships sub-group working on an environmental scan of research data preservation. DPN envisions university repositories as contributing nodes to federated replicating nodes. The replicating nodes are dark archives, used to restore access to content to contributing nodes in the event of data loss.

Reviewer's Analysis

Several GWLA and GPN institutions are among the 50+ contributing members to DPN. Based on a recent presentation at the Coalition for Networked Information (CNI), there is an assumption that institutions will affiliate with some content node and that different nodes may offer different services. The California Digital Library and the Texas Digital Library are members as well who represent large-scale digital collections.

Considerations and Recommendations

This merits our attention as a future preservation strategy to undergird efforts to make data available in access repositories.

Key Contacts:

Sponsors:

Funding:

Inception:

Geographic Location:

Ex Libris Rosetta

PRESERVATION

<http://www.exlibrisgroup.com/category/RosettaOverview>

Deborah Ludwig

Brief Description of Project

Rosetta is a commercial software solution for digital preservation. It was designed in response to the needs of the National Library of New Zealand. Other customers include the ETH-Bibliothek in Sweden, which is using it as a platform for preserving research data. SUNY Binghamton University and Brigham Young University are also U.S. higher education customers. There is good information available on the website.

Reviewer's Analysis

Ex Libris Rosetta is the only major commercial software solution for digital preservation. In addition to the National Library of New Zealand, the Church of the Latter Day Saints is another large customer with sizeable large digitization efforts.

Considerations and Recommendations

Ex Libris has extensive experience building software solutions for libraries. Products include Voyager, Aleph, and Alma integrated library systems. While Rosetta could preserve research data as part of a solution for digital library collections that could include research data, its purpose is much broader than research data.

Key Contacts: <http://www.exlibrisgroup.com/category/ContactUs>

Sponsors: Ex Libris

Funding: For-Profit-Company

Inception: --

Geographic Location: --

Brief Description of Project

LOCKSS stands for “Lots of Copies Keep Stuff Safe.” Private LOCKSS networks are collaborative communities that work together to preserve their institutional assets. There are a number of LOCKSS communities, including the Data Preservation Alliance for the Social Sciences (Data-PASS), which is focused on research data. Data-PASS⁴¹ is a consortium with the initial goal “to create a sustainable partnership model for preserving ‘at risk’ social science data.” LOCKSS was funded by the Library of Congress National Digital Information Infrastructure and Preservation Program (NDIIPP)⁴² and the effort is lead by Harvard’s Institute for Quantitative Social Science, which is noted in the literature review on federated and collaborative approaches to data management.

LOCKSS lists eleven collaborative communities on its website, including CLOCKSS which is focused on preservation of scholarly journal literature as an effort of publishers working with librarians.

Key Contacts: <http://www.lockss.org/contact-us/>

Sponsors: LOCKSS is a not-for-profit organization.

Funding: LOCKSS has received funding from Andrew W. Mellon Foundation, the National Science Foundation, and the Library of Congress.

Inception: 1999 at Stanford University with participation from Indiana, Emory, and the New York Public Library. Released into production, 2004.

Geographic Location: Global adoption

⁴¹ <http://www.data-pass.org>

⁴² <http://www.digitalpreservation.gov/index.php>

OAIS

PRESERVATION

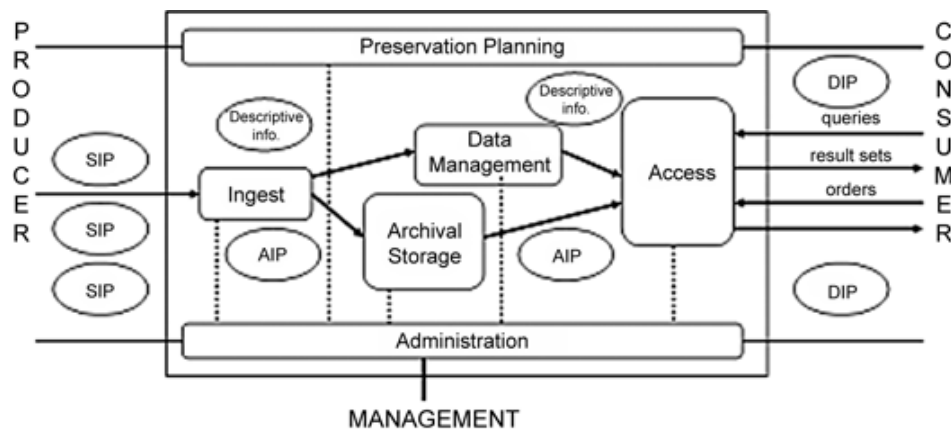
(Open Archival Information Systems Reference Model – ISO 14721:2003)

http://www.iso.org/iso/home/store/catalogue_tc/catalogue_detail.htm?csnumber=57284

Deborah Ludwig

Brief Description of the Project

OAIS is not a computer system. OAIS is a reference model (ISO 14721:2012) that defines an open archival information system (OAIS), an archive, with people and systems responsible for to preserve information and making it available for a designated community. OAIS is not a particular system or repository, but a standard, which provides a model for understanding what must be in place for long-term preservation of digital information. The classic illustration of an OAIS below identifies the components for preserving information bracketed between the roles of information producer and consumer.



Note: 3 OAIS Functional Entities Standard Diagram

Information packages within respect to the archive include the SIP (submission information package) the AIP (archival information packet) and the DIP (dissemination information packet). Functions of the archive include ingest, storage, administration, planning, and provision of access. For access, descriptive information is required.

Reviewer's Analysis

The Digital Preservation Coalition defines digital preservation as the “series of managed activities necessary to ensure continued access to digital materials for as long as necessary.”⁴³ OAIS creates a framework within which these managed activities are carried out. OAIS

⁴³ <http://www.dpconline.org/advice/preservationhandbook/introduction/definitions-and-concepts>

specifies a set of mandatory responsibilities that the people and systems comprising the archive must undertake to account for preservation of the data and its long-term access for the intended consumer. A repository utilizing an OAIS framework must:

- Negotiate for and accept appropriate information from information producers.
- Obtain sufficient control of the information provided to the level needed to ensure long-term preservation.
- Determine, either by itself or in conjunction with other parties, which communities should become the designated community and, therefore, should be able to understand the information provided.
- Ensure that the information to be preserved is independently understandable to the designated community. In other words, the community should be able to understand the information without needing the assistance of the experts who produced the information.
- Follow documented policies and procedures, which ensure that the information is preserved against all reasonable contingencies, and which enable the information to be disseminated as authenticated.

The most important implication of OAIS for a shared data management project is likely the understanding that long-term access to research data will involve more than establishing technologies. A system will require policies, agreements with data producers, designated stewards for data, specialists in developing the specifications and workflows for the “packages” of information that must accompany a research data set as it is ingested into a repository, stored, and disseminated to the intended consumer

Considerations and Recommendations

OAIS provides a high-level model for what needs to happen for long-term preservation of research data. It provides a reference to help us understand what happens in a preservation system. The advice of Chris Rusbridge (with the Digital Curation Centre) may provide a reasoned approach to thinking about how to implement a reasonable OAIS framework:

Investment in digital preservation is important for cultural, scientific, government and commercial bodies. Investments are justified by balancing cost against risk; they are about taking bets on the future. The priorities in those bets should be: first, to make sure that important digital objects are retained with integrity, second to ensure that there is adequate metadata to know what these objects are, and how they must be accessed, and only third to undertake digital preservation interventions.

Key Contacts: ISO, the International Standards Organization

Sponsors: Consultative Committee for Space Data Systems (CCSDS)

Funding: NA

Inception:

Geographic Location: NA

Additional Resources

- Allinson, Julie. *OAIS as a reference model for repositories*. November 21, 2006 accessed online at www.ukoln.ac.uk/repositories/.../oais.../Drs-OAIS-evaluation-0.5.pdf
- Digital Preservation Coalition_website <http://www.dpconline.org>
- Digital Preservation Coalition_handbook with definitions.
<http://www.dpconline.org/advice/preservationhandbook/introduction/definitions-and-concepts>
- Rusbridge, Chris. *Excuse Me ... Some Digital Preservation Fallacies?* *Ariadne*, 46, February 8, 2006. <http://www.ariadne.ac.uk/issue46/rusbridge>

Brief Description of the Project

Claims of trustworthiness are easy to make but are thus far difficult to justify or objectively prove. As Clifford Lynch has stated, 'Stewardship is easy and inexpensive to claim; it is expensive and difficult to honor, and perhaps it will prove to be all too easy to later abdicate' (2003) Establishing more clear criteria detailing what a trustworthy repository is and is not has become vital. ⁴⁴

TRAC is not a computer system for digital preservation. TRAC is a set of criteria applicable to a range of digital repositories and archives, from academic institutional preservation repositories to large data archives and from national libraries to third-party digital archiving services. TRAC is now under the management of the US Center for Research Libraries In general terms, TRAC:

- Provides tools for the audit, assessment, and potential certification of digital repositories
- Establishes documentation requirements required for audit
- Delineates a process for certification
- Establishes appropriate methodologies for determining the soundness and sustainability of digital repositories

TRAC provides tools for the audit, assessment, and potential certification of digital repositories; establishes the documentation requirements required for audit; delineates a process for certification; and establishes appropriate methodologies for determining the soundness and sustainability of digital repositories. Only a few digital repositories have pursued the certification process to date, including: *Portico*, *HathiTrust*, *Chronopolis*, the US National Archives and Records Administration (*NARA*), and Michigan's *Interuniversity Consortium for Political and Social Research (ICPSR)*. Many more have utilized the TRAC checklist as a planning tool for building trusted repositories of digital information.

Reviewer's Analysis

How does the project address, or potentially address, key facets of lifecycle management? The TRAC checklist is used by institutions for planning long-term preservation of cultural heritage and

⁴⁴ TRAC Certification Checklist. http://www.crl.edu/sites/default/files/attachments/pages/trac_0.pdf

other digital resources and has been used in combination with the OAIS reference model as a digital preservation planning tool. TRAC audit specifications consider 3 major areas:

- Organizational infrastructure, including governance, accountability, and staffing
- Digital object management
- Technologies, technological infrastructure, & security

An approach to data management that is federated across a group of partners represents an opportunity to share the work and develop as partners to implement standards and best practices. At the same time, it represents a level of complexity that can only be helped by using time-tested approaches and tools for planning and implementation that have worked for a variety of other organizations.

Considerations and Recommendations

Used in conjunction with the OAIS reference model, the *Trustworthy Repositories Audit & Certification: Criteria and Checklist (TRAC)* can provide a useful guide to understanding identifying the organizational and the technical components required for developing an approach to long-term access to research data.

Sponsors: Center for Research Libraries (<http://www.crl.edu/>).

Inception: The TRAC checklist was published in 2007 by the National Archives and Records Administration, Research Libraries Group, and Center for Research Libraries.

Additional Resources

- Digital Curation Centre Web Site on Trusted Repositories:
<http://www.dcc.ac.uk/resources/repository-audit-and-assessment/trustworthy-repositories>
- Center for Research Libraries Web Site on Archiving and Preservation:
<http://www.crl.edu/archiving-preservation>

3. Archiving & Repository Services

Data archiving definitions vary widely. In the simplest terms and in accordance with the National Science Foundation Data Archiving Policy ⁴⁵, we consider data archiving to be the sharing of data through an appropriate archive or library to encourage data sharing with other researchers. Data archiving is something more than simply storing data and backing it up and may include some common elements with digital preservation.

For the purposes of this report, we have focused more narrowly on systems for making data available, most often through some type of repository software. Initiatives to curate data will generally rely upon some repository system as a basis for making data available to appropriate consumers.

Examples of Data Archiving Platforms:

1. The DataFlow Project (with DataBank and DataShare)
2. DataSpace
3. DuraSpace
4. Hydra
5. OneShare (with DataUP)

⁴⁵ <http://www.nsf.gov/sbe/ses/common/archive.jsp>

DataFlow Project [DataFlow, DataStage, DataBank]

ARCHIVING

Greg Monaco, Great Plains Network (GPN)

<http://www.dataflow.ox.ac.uk/>

Brief Description of the Project

DataFlow is a two-stage data management infrastructure: locally, the user installs *DataStage* and the user's institution installs *DataBank*. The user then saves a dataset, locally, to her/his hard drive. That dataset will be copied to the cloud (similar to dropbox). The software is open source.

From the website:

Rather than storing datasets on external hard drives in the lab, *DataFlow* lets researchers save their work in institutional memory banks. The system will be lightweight (nothing for researchers to install; just save data to a mapped drive on their computer), with best-practice standards to make sure data is well looked after.

DataStage is a secure personalized 'local' file management environment for use at the research group level, appearing as a mapped drive on the end-user's computer. It can be deployed on a local server, or on an institutional or commercial cloud. Once the software has been installed on the server, there is no additional software for the end-user to install.

Users save files to *DataStage* just as they would on ordinary C: drive -- but with added extras:

- Private, shared and collaborative directories, with password-controlled access
- Web access – work with stored files over the web, anywhere in the world
- Users can add richer metadata via the web interface, using free-text "notes" fields
- All files can be automatically backed up via your usual backup service
- Users can invite colleagues to access group files, via password control
- Repository submission interface makes it easy for researchers to define data packages, enter minimal metadata, and deposit them in a repository of choice. The minimal metadata is in RDF format; additional (non-RDF) metadata can be added via free-text fields at the submission stage
- Packaging done using BagIt file packaging specification, soon to be SWORD-2 compliant
- Flexibility to dynamically invoke additional cloud storage as required

DataBank is a scalable data repository designed for institutional deployment that is designed to

- provide a definitive, sustainable, reference-able location for (potentially large) research datasets
- allow researchers to store, reference, manage and discover datasets

DataBank instances will expose both human- and machine-readable metadata describing their datasets, and will assign Digital Object Identifiers (DOIs) to hosted datasets, obtained

automatically using the DataCite API, to aid discovery and citation...By default, all objects are assigned a DOI and a cc-zero Open Data Waiver, and all RDF-format metadata is visible to the outside world, but other licensing/secretcy arrangements can be accommodated. Users can define an optional embargo period (making metadata visible but withholding the underlying data), add richer metadata to make their data easier to find (when searching within DataBank, or via web crawlers like Google), and users can revise datasets that have already been submitted (new DOI issued for each version, all versions kept in perpetuity). DataBank can also be run as a “dark” archive with metadata and data invisible to the outside world.

- Institutions can have their own *DataBank* instances hosted within an external cloud (e.g. Eduserv), or can choose to deploy *DataBank* on local hardware, at institutional, departmental or individual research group level.
- *DataBank* can be used together with *DataStage*, or separately.
- DataBank is a virtualized, cloud-deployable version of the databank created by Oxford's Bodleian Libraries. We are actively pursuing a variety of sustainability options for *DataBank*, but at minimum, the software will be maintained and developed for use by the Bodleian Libraries, with their code made available open-source under an MIT license.

Reviewer's Analysis

This project presents a different approach to management of research data, and encompasses all stages from creation to archiving and curation to sharing and reuse. This projects attempts to tackle the consortial/collaborative approach and promises to be quite flexible.

Considerations and Recommendations

Merits further review and possible download and testing.

Key Contacts:

Sponsors: **University of Oxford**

Funding: JISC (see <http://www.dataflow.ox.ac.uk/index.php/about/51-project-funding>)

Inception: Version 0.1 of the software packages were released on March 2, 2012. The project appears to have started in 2011.

Geographic Location: UK

DataSpace

ARCHIVING

<http://dataspace.princeton.edu/jspui/>

Michael Bolton, Texas A&M University, Sterling C. Evans Library

Brief Description of the Project

DataSpace is a digital repository meant for archiving and publicly disseminating digital data that are the result of research, academic or administrative work. DataSpace has a cost model for Princeton users; however, it is billed as a one-time charge for long-term storage of digital data. The repository also provides persistent URLs, which aid in dissemination of works. DataSpace accepts a wide range of datasets from research data to student projects and reports, conference and workshop proceedings, technical reports and digital collections of images and other digital assets. It is an Open Access repository and users can subscribe to news feeds that will keep them informed of new submissions to Communities in DataSpace. The web site says future plans do include making metadata available to services such as OAIster.org to help with discovery services. Princeton has developed a LibGuide for this service (<http://libguides.princeton.edu/content.php?pid=211802&sid=1763060>) and I have seen presentations on the service (Web Seminar for EDUCAUSE).

Reviewer's Analysis

DataSpace is built on DSpace, which would be convenient for DSpace-base institutions. I believe this is a separate instance of DSpace, intended solely for dataset management. It has the same concept of communities and collections, just as with an Institutional Repository. For a DSpace shop, the paradigm would be easy to follow and has appeal since it would leverage existing knowledge and expertise. The **About** page includes information on the licenses as well as a cost model for storing data.

DataSpace is a data repository and does not by default include a preservation model. That would have to be added. However, if you add preservation to your DSpace IR, you now have preservation for the datasets as well. In this case, the DPN project would work well for us.

How does the project address, or potentially address, key facets of lifecycle management? This model leverages existing expertise in DSpace and an institution's repository. It can be easily linked to research or back to items in the IR. It does not have a preservation model as yet and that would have to be added. It does support easy data sharing via the persistent URLs; however, it must be coupled with a stronger discovery service for broader access.

Does this project have implications for a collaborative or consortial approach to defining a regional (GWLA/GPN) strategy for lifecycle management of research data? If so, what does this project offer in terms of approaches that might be incorporated or services that might be acquired? Given the number of

DSpace instances in use this would seem to be a good place to start in testing technologies. DSpace is very robust when it comes to sharing data (via the OAI interfaces). With SWORD deposits possible, deploying and using DataSpace could be fairly straightforward and easy.

Considerations and Recommendations

What I like about this project is it represents a quick way to get into dataset management for existing DSpace shops. I think it would need a metadata management tool since DSpace is not particularly strong in that area. If coupled with DataUP, and DataUP was able to export to DSpace, this would be a quick win.

Key Contacts: Serge Goldstein, Associate CIO and Director of Academic Services at Princeton

Sponsors: Princeton initiative

Funding:

Inception: 2010?

Geographic Location: Princeton

DuraSpace (DSpace & Fedora)

ARCHIVING

<http://www.duraspace.org/>

Philip Konomos, Arizona State University.

Brief Description of the Project

DuraSpace is a not-for-profit umbrella organization that manages, develops, coordinates and supports three important open source digital repository projects:

- DSpace (<http://www.dspace.org/>)
- Fedora (<http://www.fedora-commons.org/>)
- DuraCloud (<http://www.duracloud.org/>)

DSpace is a free, open source digital repository application. The software allows institutions to manage, preserve, and provide long-term access to many common types of digital files. DSpace is probably one of the most commonly used digital repository / institutional repository applications at research universities around the world because it is configured to work “out-of-the-box” with minimal need for custom programming. A large number of universities use DSpace including MIT, University of Texas, and the University of Illinois (click here [for a comprehensive list](#)).

Fedora is also a free, open source, digital repository application, although it is much more complex than DSpace to configure. Fedora is an acronym for Flexible, Extensible, Digital, Object Repository Architecture, which indicates that this is really a sophisticated framework, rather than an out-of-the-box repository package. Fewer universities use Fedora, but those include Cornell, Indiana University, University of Virginia, and Rutgers. Fedora offers much more robust functionality and complete customization but at the cost of significant programming staff expertise.

DuraCloud is an outgrowth of a 2009 initiative between the Library of Congress National Digital Information Infrastructure and Preservation Program (NDIIPP) and DuraSpace. It is designed to offer a hosted solution (Software as a Service, or SAAS) for institutions that want to ensure perpetual access to their digital collections using cloud computing.

DSpace and Fedora were originally separate initiatives, but in 2009 the two organizations saw the wisdom of joining forces to leverage development ideas and technical expertise under the rubric of the DuraSpace Foundation. DuraCloud was one of the first new initiatives begun by the newly formed DuraSpace Foundation.

Reviewer's Analysis

These are important repository options for universities to consider since selection of the appropriate digital repository platform is a critical decision. Fedora and DSpace have a long history of use by a wide variety of institutions and repository managers are well aware of the strengths and weaknesses of each platform. DuraCloud, as a more recent service, has less of a track record for interested users to use for evaluation.

Considerations and Recommendations

Since both DSpace and Fedora have been widely adopted by American universities, it might be worthwhile to invite GWLA members who use one or the other option to present a demonstration and evaluation.

Key Contacts: The project email is info@duraspace.org. Michele Kimpton, Chief Executive Officer; DuraSpace, 28 Church Street, Unit #2, Winchester, MA 01890

Sponsors: Fedora was originally developed (1997) by Professor Sandy Payette and graduate student Carl Lagoze at Cornell University. By 2007, as it's community of users grew, faculty at several of the universities that used Fedora formed the Fedora Commons not-for-profit organization. DSpace was originally developed (2002) by MIT Libraries and HP Labs, and later (2007), as the number of users grew, supported by the non-profit DSpace Foundation. In both cases, a large community of universities and other institutions committed significant time and programming expertise to support the growth of each software project. The non-profit organizations were a way to organize and manage development and support user institutions. Currently, DuraSpace is run by a [staff of nine](#) and governed by a [Board of Directors](#) of seven.

Funding: Fedora was funded by grants from the Andrew W. Mellon Foundation and the Gordon and Betty Moore Foundation. Original funding for **DSpace** came from the Andrew W. Mellon Foundation and HP (Hewlett Packard). **DuraCloud** was developed with support from the Gordon and Betty Moore Foundation, the Andrew W. Mellon Foundation and the Library of Congress NDIIPP program.

Inception: Fedora was developed by computer science faculty at Cornell University in 1997. DSpace was started in 2002, and remained a separate project until 2009. DuraCloud was begun by DuraSpace as a hosting option in 2009-1010.

Geographic Location: N/A

Globus Online

ARCHIVING

www.globusonline.org

Greg Monaco, Great Plains Network (GPN)

Brief Description of the Project

The *Globus Project*'s original focus was on middleware tools to enable grid computing. *Globus Online* is the Globus Project's response to the issue of researchers needing to manage huge datasets. Globus Online provides the researcher with a secure method to organize, access, move and share with collaborators data that is located at multiple, distributed sites. Access management interfaces with Shibboleth and InCommon (see separate reviews). Rather than an institutional solution, Globus Online provides the individual researcher with a way to organize all the datasets to which s/he needs access from the Globus Online portal.

Globus Connect and *Globus Connect Multi-User* are software versions that allow one to make local storage (desktop, laptop) and shared storage (servers, data repositories) resources available to a potential community of users via Globus Online.

Note: The Globus Project team has historically been interested in collaborating with others making novel use of their tools.

From the website:

Globus Online is a fast, reliable file transfer service that makes it easy for any user to move any data anywhere. Recommended by HPC centers and user communities of all kinds, Globus Online automates the time-consuming and error-prone activity of managing file transfers, so users can stay focused on what's most important: their research.

Reviewer's Analysis

Globus Online provides individual researchers with a way to manage access to their personal data sets, from data creation, processing and analysis to data sharing and reuse. The ability to create groups and share data means that users may create teams who initially add to and reuse data.

With Globus Connect Multi-User it is possible to make test bed resources available to potential users via Globus Online. This offers the advantages of providing

- a common interface for access to resources without having to reinvent it,
- an interface that can be branded for this project,
- integration with other key technologies (Shibboleth, InCommon)

Considerations and Recommendations

I recommend contacting the Globus Project Team to gauge their interest in participating in the further development of this project.

Key Contacts: Steve Tuecke, Ian Foster

Sponsoring entities: University of Chicago, Argonne

Funding Source, if Known: DOE (Energy), NSF, NIH

Inception Date, if Known: Globus Project (1995), Globus Alliance (2003)

Geographic Location, if Applicable: Headquartered in Chicago, IL

Hydra

ARCHIVING

<http://projecthydra.org/>

Jason Stirnaman, University of Kansas Medical Center

Brief Description of the Project

Hydra adds interfaces for discovery, workflow, and access control to Fedora-based data through a Ruby on Rails framework, Solr, and Blacklight. Enables powerful use of Fedora's capabilities through a familiar lightweight, Ruby-based toolkit. "Hydra's flexibility means we use the same infrastructure, but can generate individual solutions."

Hydra is not just a repository software solution. Rather, it is three complementary components:

- A vibrant, highly active [community](#) supporting the work of the project which shares an [underlying philosophy](#) behind all that it does
- [Design \(and other\) principles](#) involved in constructing a successful Hydra "head" for use with compatible digital objects, and,
- The [software components](#), the Ruby gems, that the Hydra community has constructed which are combined together to provide a local installation

Reviewer's Analysis

How does the project address, or potentially address, key facets of lifecycle management? [Data Creation, Data Processing, Data Analysis, Data Preservation, Data Access for Others (data sharing), and Data Reuse] Hydra, although extensible, is primarily concerned with the facets of preservation, access, and reuse of digital content objects. Hydra expects that objects in the repository:

- follow a model pattern asserted by the objects themselves.
- have an associated rights schema that Hydra may enforce.
- have accompanying metadata.
- may have digital content for delivery.

Ruby on Rails provides models, controllers, and interfaces to create, read, update, and delete objects and their associations. Fedora stores objects in one of two standard models. Apache Solr indexes objects and their metadata, and Blacklight provides a discovery interface to the index of objects and metadata.

Examples of research data implementations:

- History DMP Project, University of Hull
 - Interaction between Hydra and DataCite has been explored to enable the additional benefits of this widely used citable standard identifier

- “The work thus far has enabled us to exploit our institutional repository’s flexibility to use it for datasets. We shall also be using it as a data catalogue. We know we will need to exploit this flexibility further as different research data needs emerge.” [1]
- Datasets in Penn State’s Scholarsphere repository.
https://scholarsphere.psu.edu/catalog?f%5Bgeneric_file_resource_type_facet%5D%5B%5D=Dataset

Does this project have implications for a collaborative or consortial approach to defining a regional (GWLA/GPN) strategy for lifecycle management of research data? If so, what does this project offer in terms of approaches that might be incorporated or services that might be acquired? Yes, as a repository solution, Hydra has implications for a collaborative or consortial strategy for data management. The Hydra/Fedora approach to [Rights Enforcement](#) and Access Control may be of particular interest. A major advantage of Hydra is the commitment to sharing solutions that can be easily adopted by other Hydra sites. Hydra is a fully open framework built on familiar, lightweight tools and supported by a growing community. Hydra repositories are scalable, flexible, and modular.

Recommendations

Does this project merit further review as part of a site visit, phone interview, or further analysis? (and, if so, what do you recommend?) Yes, this project merits further review if a repository solution or support is considered. I recommend a phone interview with core partners or further analysis that addresses specific GWLA/GPN use cases.

Additional Resource

- Using Hydra’s flexibility to manage datasets. 2012.
<https://hydra.hull.ac.uk/catalog/hull:6010>

Key Contacts: [Matt Zumwalt](#), Mediashelf; [Bess Sadler](#), Stanford University; [various communication channels listed on Duraspace wiki](#).

Sponsors: Stanford University, University of Virginia, University of Hull, DuraSpace, MediaShelf LLC, University of Notre Dame, Northwestern University, Columbia University, Penn State University, Indiana University, London School of Economics and Political Science, The Royal Library of Denmark. [Governance model](#)

Funding: “Hydra is *not* (and has never been) grant funded. It *is* distributed, robust and open. Any single developer could walk away. Any single institution could walk away. People ask what’s your sustainability plan? We say we’ve already passed the first hurdle—more than four years of self-funded productivity, and a growing code, contributor and user base, not dependent on a transition plan.”

Inception: 2008

Geographic Location: Distributed.

Islandora Open-Source Digital Asset Management Framework

ARCHIVING

<http://islandora.ca/>

Susan Matveyeva, Wichita State University

Brief Description of the Project

Islandora is an open source framework developed by the University of Prince Edward Island's Robertson Library. Islandora combines the Drupal and Fedora open software applications to create a robust digital asset management system that can be fitted to meet the short and long term collaborative requirements of digital data stewardship. Additional open source applications are added to this core stack to create what we call Solution Packs. Islandora operates under a GNU license. Islandora may be used to create large, searchable collections of digital assets of any type and is domain agnostic in terms of the type of content it can steward. (from the Islandora website)

Reviewer's Analysis

Islandora supports data lifecycle. Islandora implemented by the Colorado Research Alliance Repository⁴⁶. The infrastructure is described on the website.⁴⁷

Features of the Colorado Research Alliance Digital Repository:

Fully Hosted Service: ADR Services centrally manages all the hardware, software, updates, and backups for the repository.

Customization with Drupal: Customize the look and feel of the repository front end with the Drupal web content management system.

Access and Authentication: Manage restricted, embargoed, or dark archive content with user accounts and security metadata.

Content Loading: Add content to the repository with easy web forms for adding metadata and attaching files.

Search: Solr indexes metadata and the full text of PDFs and other documents, which can be searched with simple and advanced search functions.

Streaming and Viewing: Embedded viewers and players display common formats of document, image, audio, and video content without users having to download anything from the repository.

Cloud Storage: Back up objects in the cloud and perform fixity checks with DuraCloud.

Content Sharing: Share OAI metadata for harvesting and aggregation into other sites. Every object is assigned a Handle link for persistent identification.

⁴⁶ <http://adrresources.coalliance.org/>

⁴⁷ http://adrresources.coalliance.org/?page_id=13

Reusability: FedoraCommons open-source repository software supports management, reuse, migration, and transformation activities on digital objects for access and preservation.

Hardware: ADR hardware is stored at a collocation facility to meet power, cooling, and security needs. Data is backed up to disk and tapes. For more information about the ADR hardware, please contact adr@coalliance.org.

Software: The ADR Basic repository platform runs on Islandora, an open-source Drupal-based repository system, and uses Fedora Commons as its core repository software.

Cloud Storage: In 2010, ADR Services participated in a pilot project for DuraCloud, cloud repository management software from DuraSpace, the parent organization of Fedora Commons. In fall 2011, the ADR joined the DuraCloud service and will begin using its cloud services for remote backup and bit integrity in 2012.

Descriptive Metadata: All the objects in the repository have a basic set of descriptive information in MODS (Metadata Object Description Schema), which is a set of metadata designed to describe resources commonly found in libraries. The ADR and its members follow best practices for MODS as set out by the Digital Library Federation's Implementation Guidelines for Shareable MODS and MODS Levels of Adoption.

Many ADR members do not have pre-existing MODS metadata for their records. Members can create a MODS record for an object by filling out a web form in the repository software, which then builds MODS. Alternatively, members can send metadata to ADR Services in MARC, Dublin Core, or even a spreadsheet and we will transform the metadata into MODS that can be used by the repository. The MODS metadata is currently being used in the search index and object display, and is being transformed to Dublin Core for OAI harvesting. ADR Basic is highly flexible with what metadata it can accept and display. Avenues for future development include objects described in VRA Core 4.0 or Darwin Core.

Security Metadata: The ADR uses XACML security metadata to control access to collections, objects, and data streams in the repository.

Considerations and Recommendations

I believe, the project and its implementation by a consortium worth a closer look

Key Contacts: Mark Leggott, University Librarian and developer of Islandora islandora@upei.ca and the founder of a company DiscoveryGarden Inc. that provides services around Islandora software <http://www.discoverygarden.com/> Mark Leggott website: <http://loomware.typepad.com/about.html>

Sponsoring entities: University of Prince Edward Island

Funding Source, if Known: multiple: 2,3 mini grants (2010); over 750,000 research projects, donations, library operational budget (see: Leggott presentation:

http://loomware.typepad.com/docs/S2I2_Islandora.pdf)

Inception: 2008

Geographic Location: Charlottetown, Prince Edward Island

ONEShare Repository [& DataUp]

ARCHIVING

<http://www.cdlib.org/cdlinfo/2012/10/02/california-digital-library-and-partners-launch-dataup/>

Michael Bolton, Texas A&M University, Sterling C. Evans Library

Brief Description of the Project

ONEShare is a repository build specifically for DataUP. ONEShare is built on the Merritt Repository software developed by Stephen Abrams at CDL.⁴⁸ Merritt is used extensively by CDL and a special instance was developed to support the DataUP project. DataONE in turn harvests metadata from ONEShare via API.

Reviewer's Analysis

DataUP and ONEShare should be viewed as a single entity for our purposes since developing a generic export facility to support other repositories could take a while. In this case, CDL relied heavily on technologies it knew well and used as part of their infrastructure, which is a smart and safe move. It does mean we would have to adapt the tools for the general use case we are developing. Merritt is an open source project with the code available on Bitbucket.

Considerations and Recommendations

If we were to test DataUP, ONEShare would need to come along for the ride.

Key Contacts: California Digital Library (CDL)

Sponsors: CDL UC3

Funding:

Inception: Launched with DataUP, October 2, 2012

Geographic Location: Open source software

⁴⁸ <http://www.cdlib.org/services/uc3/merritt/index.html>

4. Storage Systems

Storage for digital data is foundational to any program of curating, preserving or archiving data. Storage systems are becoming more sophisticated and able to do more than act as simple file directories with backup. The line between storage and archives is blurring. iRODS stands out as something more than storage. With a rules engine and metadata catalog, it shares characteristics with repositories.

Three storage platforms are reviewed in this section.

1. iRods
2. Quantum StorNext
3. SGI DMF & LiveArc

iRODS (Integrated Rule-Oriented Data System)

STORAGE+++

<http://www.irods.org/>

Michael Bolton, Texas A&M University, Sterling C. Evans Library

Brief Description of the Project

Based on the web site, iRODS is a data grid software system developed by the Data Intensive Cyber Environments (DICE) research group, which developed the Storage Resource Broker (SRB). iRODS is scalable, released under an open source license, freely available and has a very active user community. The Service Resource Broker (SRB) was the forerunner of iRODS and was described as a one-size fits all system in that policies used to manage the data were hard-coded into the system. In the development of iRODS, a more adaptive approach was used. Now policies could be applied to data based upon rules.

The rule engine is not only very powerful; it is equally flexible as well. This allows for automated data management such as replicating data across zones. iRODS is accessed via a set of uniform APIs and GUIs. This allows for iRODS storage to be accessed through a number of mechanisms such as a web frontend or it can be integrated into more complex systems as the storage subsystem. For instance, iRODS can be the backend storage for repository software including DSpace and Fedora Commons. iRODS employs a data catalog to not only track where data is stored, but to collect a rich set of metadata for the object. Schemas of various sorts can be incorporated into the catalog to provide a more controlled vocabulary and the rules engine can be used to ensure metadata is added when objects are stored in the file system.

iRODS has gained popularity in areas such as numerically intensive computing and big data, primarily because of its ease of access. Like SRB, iRODS abstracts certain aspects of the storage subsystems by providing an API that can be used to access the system. Access to the data is via the MCAT, or iRODS metadata catalog, which tracks the movement of data through the system. Users do not need, nor care to know, where or how the data is stored. The iRODS server can find the dataset, read it from whatever media and transfer it to the client using standard protocols. A growing collection of micro-services help the iRODS server carry out its automation tasks, applying policies and rules as defined by the user and system. iRODS servers can easily be linked together to form large, geographically dispersed clusters. And the rule engine can be used to control and manage replication of data across the various nodes in a cluster.

The iRODS user group meets annually and develops a list of requested features and services, which is used by the DICE team to enhance and further develop the product. The user group meetings also showcase customers who use the service, along with their use cases and applications.

Reviewer's Analysis

I have always thought of iRODS as an infrastructure service, that is, the file system used by higher-level applications. What it offers is very attractive for us in that it abstracts the real storage layer. iRODS provides a robust set of APIs and interfaces that allow administrators to use whatever physical storage is necessary and appropriate for the service without entangling the application in the mechanics of using that storage. On the customer side, a number of different clients can interact with the iRODS servers to gain access to data. The iRODS server uses an authentication system to protect the data and manage access privileges. iRODS also has a very well developed JAVA-based SDK (Jargon) to allow developers to integrate iRODS into their applications. So we could actually develop applications that have UIs customized for our service while exploiting the features of iRODS. DPN (the Digital Preservation Network) will be using iRODS as one of its preservation services and the Texas Advanced Computing Center (TACC) is a large iRODS deployment.

Considerations and Recommendations

While standalone iRODS implementations would be feasible for this project I believe our effort would be better served by integrating iRODS with easier to use frontends. This would include web interfaces, possibly iDROP or WebDAV. Integrating iRODS into a broader service offering, such as that delivered in the iPlant Data Store would be better still.

I do like the rules-based management of data provided by iRODS. As we develop our service, we could write rules to ensure data is replicated across geographically dispersed nodes. I also like the metadata model supported by the iRODS catalog. I am not in favor of supporting multiple metadata management tools so we would have to develop a process to populate the MCAT with data from our primary metadata management tool. We then augment that metadata with information collected and managed by iRODS, such as checksums, file locations, number of copy, age of replicas, etc.

Additional Notes

Information related to the various grants awarded for SRB and iRODS

- NSF ITR 0427196, Constraint-Based Knowledge Systems for Grids, Digital Libraries, and Persistent Archives (2004–2007)
- NARA supplement to NSF SCI 0438741, Cyberinfrastructure; From Vision to Reality—Developing Scalable Data Management Infrastructure in a Data Grid-Enabled Digital
- NARA supplement to NSF SCI 0438741, Cyberinfrastructure; From Vision to Reality—Research Prototype Persistent Archive Extension (2006–2007)

- ❖ NSF SDCI 0721400, SDCI Data Improvement: Data Grids for Community Driven Applications (2007–2010)
- ❖ NSF/NARA OCI-0848296, NARA Transcontinental Persistent Archive Prototype (2008–2012)

Some key sites or organizations using iRODS as part of their data management plan or system.

- iPlant Data Store - <https://pods.iplantcollaborative.org/wiki/display/start/Storing+Your+Data+with+iPlant+and+Accessing+that+Data>
- University of Michigan Office of Research – Cyberinfrastructure - Data Management - <http://orci.research.umich.edu/resources-services/data-management/>

Key Contacts: PI – Reagan Moore – Director of the Data Intensive Cyber Environments Center (DICE)

Sponsors: National Science Foundation and the National Archives and Records Administration⁴⁹

Funding: NSF and NARA

- First NSF grant awarded Fall 2004 for SRB
- Latest NSF Grant awarded September 28, 2011 for prototype national data management infrastructure⁵⁰

Inception: iRODS began in 2006 with release 0.5 released December 20, 2006. It is a continuation of the SRB project, which started in 2004. The first iRODS users group meeting was held in 2009.

Geographic Location: Widely used in a number of systems and applications.

⁴⁹ Sponsors of iRODS - <https://www.irods.org/index.php/Sponsors>

⁵⁰ Press release - <http://www.renci.org/news/releases/nsf-datamet>

Quantum StorNext

STORAGE

<http://www.quantum.com/products/software/stornext/index.aspx>

Toby Axelsson, University of Kansas, Information Technology

Brief Description of the Project

Quantum StorNext is a multi-tier (tape & disk) archival system providing access to a clustered file system via NFS, CIFS or direct mount using provided software. StorNext is well suited for storing large amounts of data (Petabytes).

Reviewer's Analysis

How does the project address, or potentially address, key facets of lifecycle management? Quantum StorNext provides means of storing and archiving large amounts (Petabytes) of data.

- Access can be provided through common file access protocols (NFS & CIFS) as well as direct network mount using provided software (Distributed LAN Clients). The Distributed LAN Client software supports both Windows and Linux.
- Data can easily be accessed after it has been archived, with delays if the data needs to be re-hydrated from tape (transparent to user).
- StorNext can be configured to handle moderate HPC workloads.
- The system is agnostic of what type of disk subsystem (SAN attached) and tape library is used.

Does this project have implications for a collaborative or consortial approach to defining a regional (GWLA/GPN) strategy for lifecycle management of research data? If so, what does this project offer in terms of approaches that might be incorporated or services that might be acquired? Not sure.

Considerations and Recommendations

Does this project merit further review as part of a site visit, phone interview, or further analysis? (and, if so, what do you recommend?) This is an interesting platform for long time storage of large amounts of data. I don't believe it does data lifecycle management to the degree that we are interested in. Coupled with a product for data lifecycle management, it may be a worthy candidate for long-term storage of research data, in particular if the total expected amount of data would exceed a Petabyte. If a solution like this fits well in the overarching vision, this is something that should be investigated further.

Key Contacts: Brian Morsch Brian.Morsch@quantum.com Sr. Account Executive

Sponsors: Quantum

Funding: Commercial

Inception: At least since the 1990's.

Geographic Location: N.A

SGI DMF & LiveArc

STORAGE

<http://www.sgi.com/products/storage/software/dmf.html>

Toby Axelsson, University of Kansas Information Technology

Brief Description of the Project

SGI DMF is a multi-tier (tape & disk) archival system providing access to a clustered file system via NFS, CIFS or direct mount using their kernel modules. DMF is well suited for storing large amounts of data (Petabytes).

Reviewer's Analysis

How does the project address, or potentially address, key facets of lifecycle management? [Data Creation, Data Processing, Data Analysis, Data Preservation, Data Access for Others (data sharing), and Data Reuse]

- SGI DMF provides means of storing and archiving large amounts (Petabytes) of data.
- Access can be provided through common file access protocols (NFS & CIFS).
- Data can easily be accessed after it has been archived, with delays if the data is stored only on tape.
- DMF can be configured to handle moderate HPC workloads.
- The system is agnostic of what type of disk subsystem (SAN attached) and tape library is used.

LiveArc provides Digital Asset Management by classifying data and creating searchable metadata based on file content (API exists for adding file types not already supported).

Does this project have implications for a collaborative or consortial approach to defining a regional (GWLA/GPN) strategy for lifecycle management of research data? If so, what does this project offer in terms of approaches that might be incorporated or services that might be acquired? Not sure.

Considerations and Recommendations

Does this project merit further review as part of a site visit, phone interview, or further analysis? (and, if so, what do you recommend?) This is an interesting platform for long time storage large amounts of data as well as the classification of such (metadata collection). I don't believe it does data lifecycle management to the degree that we are interested in. Coupled with a product for data lifecycle management, it may be a worthy candidate for long-term storage of research data, in particular if the total expected amount of data would exceed a Petabyte. If a solution like this fits well in the overarching vision, this is something that should be investigated further.

Key Contacts: Timothy Belcher tbelcher@sgi.com (Federal Channel Sales Manager, West)

Sponsors: SGI

Funding: Commercial

Inception: Not sure, but oldest still running implementation is 21 yrs. old (going through multiple generations of hardware and software revisions).

Geographic Location: N/A

5. Enabling Technologies, Services, and Components for Data Management

The final group of technical reports on data management include some examples of various technologies or services that reviewers looked at for potential to contributing to data curation, either directly or conceptually. These have been sorted into 5 categories as follows.

1. Access Management
 - InCommon
 - Shibboleth
2. Discovery
 - Blacklight
 - Databib
 - Linked Data
 - Mercury
 - Researcher Networks: VIVO / Profile / Bibapp
 - VuFind
3. Identifiers for Digital Objects and Researchers
 - ARK
 - DOI
 - ORCID
4. Licensing for Data
 - Creative Commons
5. Planning for Data Management
 - Data Curation Profiles
 - DMP Tool

InCommon

<http://www.InCommon.org>

Greg Monaco, Great Plains Network (GPN)

ENABLERS: Access Mgt.

Brief Description of the Project

InCommon provides certificates that can be used in a trusted identity and service infrastructure in conjunction with Shibboleth for secure single sign on by users to web services in a trusted/federated identity framework. It is scalable and can work in partnership with universities, schools, libraries and government agencies. There are currently 215 higher education members.

From the website:

The mission of InCommon is to create and support a common trust framework for U.S. education and research. This includes trustworthy shared management of access to on-line resources in support of education and research in the United States. To achieve its mission, InCommon will facilitate development of a community-based common trust fabric sufficient to enable participants to make appropriate decisions about the release of identity information and the control of access to protected online resources. InCommon is intended to enable production-level end-user access to a wide variety of protected resources.

Reviewer's Analysis

The InCommon Certificate Service provides the necessary layer for secure access to shared (inter-institutional/federated) web services via Shibboleth. A successful project will need to resolve the issue of access to various web services at multiple institutions without creating an administrative nightmare. InCommon and Shibboleth provide a secure method for single sign on to web services across multiple administrative domains.

Considerations and Recommendations

InCommon and Shibboleth (see separate review) are central to a streamlined federated access management approach. They have been developed by and for the higher education and research communities and will be critical for the GWLA/GPN project. It is recommended that all participants to the project undertake implementation of Shibboleth and join InCommon as entry criteria to becoming part of the project.

Shibboleth

<http://shibboleth.net>

Greg Monaco, Great Plains Network (GPN)

ENABLERS: Access Mgt.

Brief Description of the Project

Shibboleth was originally conceived as an institutional solution (with the library as a focus) to the problem of sharing information via the Internet among institutions of higher education without creating an overwhelming burden on users to have multiple credentials for multiple web services and on organizations to manage separate identity providers for each web service.

From the website:

The Shibboleth System is a standards based, open source software package for web single sign-on across or within organizational boundaries. It allows sites to make informed authorization decisions for individual access of protected online resources in a privacy-preserving manner.

The Shibboleth software implements widely used federated identity standards, principally OASIS' Security Assertion Markup Language (SAML), to provide a federated single sign-on and attribute exchange framework. Shibboleth also provides extended privacy functionality allowing the browser user and their home site to control the attributes released to each application. Using Shibboleth-enabled access simplifies management of identity and permissions for organizations supporting users and applications. Shibboleth is developed in an open and participatory environment, is freely available, and is released under the Apache Software License.

Reviewer's Analysis

Shibboleth allows users from trusted organizations to access web services provided by those organizations for data storage, data access and so forth using sign on credentials (user name and password) from their home organization. This greatly reduces the administrative overhead associated with managing a web service and it reduces the burden on the user to manage multiple IDs, one for each web service.

A successful project will need to resolve the issue of access to various web services at multiple institutions without creating an administrative nightmare. Shibboleth is a standards-based approach that offers a proven method for single sign on to web services across multiple administrative domains.

Recommendations

Shibboleth and InCommon (see separate review) are central to a streamlined federated access management approach. They have been developed by and for the higher education and research communities and will be critical for the GWLA/GPN project. It is recommended that

all participants to the project undertake implementation of Shibboleth and join InCommon as entry criteria to becoming part of the project.

Key Contacts:

Sponsors: The Shibboleth Consortium is sponsored by Internet2, JISC and SWITCH (US, UK, Swiss Higher Ed consortiums).

Funding: Originally funded by NSF Middleware Initiative and Internet2

Inception: Began as an Internet2 Middleware project in 2000.

Geographic Location:

Blacklight

<http://projectblacklight.org>

Jason Stirnaman, University of Kansas Medical Center

ENABLERS: DISCOVERY

Brief Description of the Project

Blacklight is an open source Ruby on Rails gem that provides a discovery interface for any [Solr](#) index. Blacklight provides a default user interface that is customizable via standard Rails (templating) mechanisms. Blacklight accommodates heterogeneous data, allowing different information displays for different types of objects. Community-contributed [add-ons](#) offer additional features.

Reviewer's Analysis

Blacklight primarily addresses the facet of data access by providing faceted search and filtering to indexed content and metadata. Blacklight adds an elegant and customizable user interface on top of a Solr search index. Blacklight is also used as the default discovery interface in Hydra (see separate review of Hydra)

It should be noted that Solr includes robust tools for text analysis and can be employed as a tool for data analysis, statistics, and normalization. For example, see Erik Hatcher's presentation on [rapid prototyping Data.gov with Solr](#). In that sense, Blacklight + Solr may address data processing or data analysis needs.

Blacklight + Solr offer a pure and modern discovery layer for data and metadata.

Recommendations

I recommend further analysis that addresses specific GWLA/GPN use cases for: large-scale data analysis, data discovery, metadata discovery.

Key Contacts: Bess Sadler, Stanford University; Jonathan Rochkind, Johns Hopkins University
Sponsors: The University of Virginia, Stanford University, Johns Hopkins University, and WGBH are the principal contributors to the code base and use it heavily at their institutions. There are dozens of sites worldwide that use Blacklight.

Funding: Blacklight is supported by community members through their technical leadership and contributions.

Inception: 200x? Blacklight was originally developed at the University of Virginia Library and is made public under an Apache 2.0 license.

Geographic Location: Distributed.

Brief Description of the Project

Developed by [Purdue University Libraries](#), Data**bib** is a free, online database that contains short bibliographic records describing disciplinary repositories that hold / accept research data. Their website states that Data**bib** is “a searchable catalog of research data repositories.” Data**bib** is designed to help faculty, librarians, and others answer questions such as:

- What repositories are appropriate for a researcher to submit his or her data to?
- How do users find appropriate data repositories and discover datasets that meet their needs?
- How can librarians help patrons locate and integrate data into their research or learning?

The entries are international, though much more representative of the U.S. than other countries. Currently (December, 2012) Data**bib** includes 501 repository descriptions. Additional entries are added as new information becomes available. Users can [submit the names of new data repositories](#) for consideration. Members of the Purdue University Library Data**bib** editorial board then review each suggestion and create the new entry if appropriate. Users may also volunteer to be editors for a subject field (or fields), curating existing Data**bib** records and suggesting new entries.

Data**bib** is guided by an international advisory board with representation from Europe, Australia, Asia, Africa, and North America. An editorial board is being assembled to increase the coverage and continue the curation of records in Data**bib**, according to their web site.

Reviewer’s Analysis

This is a valuable service because it provides up-to-date information about available online digital data repositories across a broad spectrum of academic disciplines. Data**bib** is easy to use, with clear, intuitive search and browse functions. Users searching Data**bib** from the [home page](#) can enter a search term(s) in the search box, or browse an alphabetical list of all repositories. For example, a search for “archaeology” returns a list of 6 international repositories:

Databib
Find Repositories | Submit | Connect | About Login/Register

Search

Search results for: **archaeology** Total number of results: **7**

...

Digital Archaeological Record, The
The Digital Archaeological Record enables researchers to contribute knowledge about human history, a...

Archaeology Data Service
Archaeology Data Service (ADS) provides archaeological data from the early prehistoric to present in...

Arts and Humanities Data Service (AHDS)
The Arts and Humanities Data Service [AHDS] was a UK national service aiding the discovery, creation...

Data Archiving and Network Services (DANS)
Promotes sustained access to digital research data relating to scientific databases and e-publicatio...

Open Context
Open Context is an open source discovery tool for the publication of data collected in Archeological...

Edinburgh DataShare
Edinburgh DataShare collects research datasets produced at the University of Edinburgh, across a var...

Selecting the first one (tDAR) opens a new page with standard descriptive information and the URL for the repository:

Databib
Find Repositories | Submit | Connect | About Login/Register

[Go Back]

Title: Digital Archaeological Record, The
URL: <http://www.tdar.org/>
Authority: Digital Antiquity
Subjects: Anthropology Archaeology History

Description:
The Digital Archaeological Record enables researchers to contribute knowledge about human history, as well as allowing resource managers to preserve and protect archaeological resources. tDAR provides access to current and historic digital data, as well as the tools to analyze that data, through databases, spreadsheets, documents, and images.

Access: Open
Location: United States
Reuse: Open
Deposit: Instructions for deposit can be found at <http://www.tdar.org/support/contribute/>
Type: disciplinary

0 0 0

Annotations
Add an Annotation:

Considerations and Recommendations

Library staff among GWLA members will benefit from becoming familiar with Data**ib**. This is a valuable resources that is easy to use. The fact that Purdue staff and affiliates actively search for new content and manage the citations ensures that information in the catalog is up-to-date. While most universities now have an in-house institutional or digital repository, disciplinary repositories provide a number of advantages. Disciplinary repositories often use richer metadata, frequently containing elements specific to that particular field, for example.

Key Contacts:

The project email is databib@gmail.com The key contact is:

[Michael Witt](#), Interdisciplinary Research Librarian

Assistant Professor of Library Science

Purdue University Libraries 765-494-8703 mwitt@purdue.edu

Sponsors: Data**ib** was developed by Purdue University Libraries' [Distributed Data Curation Center \(D2C2\)](#), and is managed by [Michael Witt](#).

Funding: Creation and development of Data**ib** was funded through a 2011 grant from the Institute for Museum and Library Services (IMLS).

Inception: 2011

Geographic Location: Located at Purdue University but containing international information.

Linked Data

ENABLERS: Discovery

Jason Stirnaman, University of Kansas Medical Center

<http://www.w3.org/standards/semanticweb/data>, <http://linkeddata.org>,

<http://www.w3.org/wiki/SweoIG/TaskForces/CommunityProjects/LinkingOpenData>

Brief Description of the Project

Linked Data lies at the heart of what Semantic Web is all about: large scale integration of, and reasoning on, data on the Web.[4] Linked Data describes a recommended best practice for exposing, sharing, and connecting pieces of data, information, and knowledge on the Semantic Web using URIs and RDF." [1]

The Semantic Web isn't just about putting data on the web. It is about making links, so that a person or machine can explore the web of data. With linked data, when you have some of it, you can find other, related, data. "[2]

Linked Data functions through links between arbitrary things described by RDF. The URIs identify any kind of object or concept, but regardless of HTML or RDF, the same expectations apply to make the web grow:

1. Use URIs as names for things
2. Use HTTP URIs so that people can look up those names.
3. When someone looks up a URI, provide useful information, using the established standards (e.g. RDF, SPARQL)
4. Include links to other URIs, so that more things can be discovered

Reviewer's Analysis

As of 2010, the so-called Linked Open Data cloud covers more than an estimated 50 billion facts from many different domains like geography, media, biology, chemistry, economy, energy, etc. The data is of varying quality and most of it can also be re-used for commercial purposes.

Data Management programs should consider Linked Data from the perspective of Provider as well as Consumer. Careful attention to established standards, while also allowing different subject domains to dictate which standards (e.g. vocabularies) apply to their work, is crucial to making Linked Data work. Recognized standards:

- [RDF](#)
- Microdata (embeddable data) standards: [Schema.org](#), [RDFa](#)

Linked Data has significant implications for discoverability and reuse of data. Researchers will often be both providers and consumers of Linked Data. Linked Data affords researchers and data providers an opportunity to increase the impact and reuse of their data. Linked Data enables consumers to more efficiently analyze data, recognize correlations between datasets, and make new discoveries.

Providing Linked Data: See “Ingredients for high quality Linked (Open) Data” by the W3C Linked Data Cookbook [5]

Consuming Linked Data: Linked Data consumers can integrate and provide high quality information and data collections to mix their own data with. Such integration enables better decision making, disaster management, knowledge management and/or market intelligence solutions. Tools such as [ontology](#) reasoners and SPARQL enable LD consumers to analyze semantically-enriched data.

Further examples and tools on the [LinkingOpenData wiki](#). WebSchemas group provides a [vocabulary and proposal](#) for extending schema.org for dataset description. Tools for publishing, consuming, and integrating semantic data are readily available for most modern programming languages and web frameworks.

Considerations and Recommendations

I recommend phone interviews with current practitioners or further analysis that addresses specific GWLA/GPN use cases for: data discovery, data reuse, data analysis.

Sources:

1. <http://linkeddata.org>
2. Berners-Lee, T. Design Issues: Linked Data. <http://www.w3.org/DesignIssues/LinkedData.html>
3. Bauer, F and Kaltenbock, M. Linked Open Data: The Essentials. <http://www.semanticweb.at/LOD-TheEssentials.pdf>
4. <http://www.w3.org/standards/semanticweb/data>
5. http://www.w3.org/2011/gld/wiki/Linked_Data_Cookbook#Ingredients_for_High_Quality_Linked_Data

Key Contacts:

Sponsors: W3C

Funding:

Inception: 2007

Geographic Location: Distributed.

DISCOVERY

<http://mercury.ornl.gov/> &

<http://www.dataone.org/software-tools/mercury-metadata-editor>

Michael Bolton, Texas A&M University, Sterling C. Evans Library

Brief Description of the Project

Borrowing heavily from the ORNL (Oak Ridge National Laboratory) website, Mercury is a web-based system for searching metadata and retrieving the associated data. It is open source software and is based on a Service Oriented Architecture. It appears to be fairly flexible in how the service is provided supporting RSS, Geo-RSS, OpenSearch, Web Services and JSR-168 Portlets. This allows for easy integration into just about any interface or application. Mercury has the ability to extract or harvest metadata from HTML pages or XML files which makes participating in the service very easy. All the data provider needs to do is post the content on a Web server and Mercury will pick it up. The data is then incorporated into a centralized index where users can search the metadata either via a simple search mechanism or more intricate advanced search services. Mercury supports a number of metadata standards including XML, Z39.50, FGDC, Dublin-Core, Darwin-Core, EML and ISO-19115. Presentation slides are available at http://mercury.ornl.gov/slides/Mercury_presentation_05152012.pdf. ONEMercury is a specialized Mercury server for the DataONE network.

Reviewer's Analysis

In the area of finding data, Mercury seems to be a strong tool and service. The harvesting capabilities make the service very easy to use by researchers – there are no special processes to run or servers to install. Just publish on a web site that Mercury can access. How a researcher publishes this metadata may vary depending on discipline. That is, some communities provide metadata management tools, or cataloging services to be used by researchers wishing to publish their work.

A number of agencies and services are already using Mercury, in particular, the DataONE project where it is included in the software tools catalog. Continuing to use the DataONE model, contributors would use a service such as DataUP to create and edit metadata. It would then be published to a repository, such as ONEShare, which is a DataONE member node. The metadata is automatically harvested and added to the Mercury index. This same model could work for our consortia. Within our consortia we would have a number of contributing nodes that would feed our Mercury-based central index. I think this would be one of the real

strengths of Mercury for us, that is, all the consortia members contributing to a centralized search index.

One option would be to join DataONE as a contributing member. DataONE is currently focused on the Earth, environment, atmospheric and ecological sciences. Repositories that fall into this general area are welcome to join as member nodes. While that would be good for researchers in those disciplines, it would leave other disciplines in a lurch. Maybe our niche could be to cover what DataONE does not.

I have found several references to a Mercury Metadata Editor tool; however, I never could find much information on the product. I did see links to the ORNL Online Metadata Editor (OME) but access to the page is restricted by a login.

Recommendations

Testing of individual tools will be difficult without a fairly robust test bed incorporating a metadata editor, repository and search engine. We can use an existing test system, such as the one being developed and deployed by DataONE or develop one of our own. From my research, I have found a number of “centers” that focus on a particular discipline, such as Earth Sciences. While that would be acceptable for initial testing, I can foresee our needs being much broader and we will need to support a cross-discipline service. That will mean we need to customize tools, schemas and best practices. Deploying our own test service appears to be about the only way we can control all the variables and make a truly informed decision on the suitability of the toolkits.

Key Contacts: mercury-support@ornl.gov

Sponsors: NASA, USGS, Office of Science – U. S. Department of Energy

Funding: Ongoing support is through the Oak Ridge National Laboratory Distributed Active Archive Center (ORNL DAAC).

Inception: First reference seems to be 2008 but it may be older than that. There are a number of deployments and great many data providers indicating this software is fairly mature.

Geographic Location: Open Source software

Brief Description of the Project

Generally, *Research Networking* systems are web-based tools for discovering and using research and scholarly information about people and resources. Though underlying architectures vary widely, most “RN” systems share a high-level set of features:

- Compile data about research-related works (book citations, article citations, grants, presentations) and resources (labs, equipment).
- Provide tools to import, harvest, curate, and enrich the data, e.g. harvesting data from enterprise sources or PubMed, providing access to full-text through localized OpenURL resolution.
- Perform some amount of machine learning to disambiguate personal names in citation data.
- Attribute this data, unambiguously, to the people, often referred to as “experts”, who are responsible for the works.
- Provide a single discovery interface for the data,
- Show collaborations and relationships among experts, e.g. through co-authorship, shared subjects, shared lab space.
- Provide data useful in measuring research output or decision-making.
- Provide data output using common format standards, e.g. use VIVO Ontology and RDF to express data and relationships through the VIVO Ontology, export styled citations, expose data objects through additional APIs as JSON, XML, etc.

Reviewer’s Analysis

Research Networking systems address the facets of data access and data reuse. They aggregate and network data that may otherwise be siloed by database vendors or enterprise systems. They offer a discovery user interface for experts, collaborators, and their works. They may display metrics or some measure of impact for a work. They provide Linked Data/Semantic Web data sources and endpoints for scholarly works.

VIVO may currently be the most widely used RN system. VIVO is a semantic web, n-store (triples and quads)-based approach to gathering and sharing data about research activity. VIVO is developed by a consortium in the US. The project is based on the code base, VIVO, and the VIVO ontology for describing research.

VIVO is modeled to provide for a federated hub-node network of metadata about research activities. As a semantic web application, VIVO data is accessible as standard RDF, enabling the

possibility of repurposing or sharing of data over the web. Each VIVO instance, e.g. at the institution-level, provides structured Linked Data that can be harvested and aggregated into a broader network. VIVO Searchlight is a practical application of this. Designed to accommodate data about multiple kinds of resources, typical bibliographic and enterprise data, as well as metadata about physical spaces and equipment. VIVO provides a Harvester Framework to facilitate ingesting data from external systems.

VIVO is primarily a Java application consisting of a Java-based UI, an underlying Jena SDB semantic store implemented with a MySQL database, and Solr for searching.

BibApp is used by a small, but growing number of institutions. BibApp specializes in the collection of bibliographic data and displaying the output of experts and groups. One unique feature of BibApp directly related to repositories is that it comes packaged with a SWORD client for allowing a user to archive a copy of a work directly into any SWORD-compliant repository, e.g. DSpace. BibApp uses the Sherpa/RoMEO API to allow an organization to monitor, which articles are open access and could be archived. BibApp exposes VIVO-compliant RDF Linked Data using the VIVO ontology, allowing an institution's BibApp data to be discoverable alongside VIVO nodes.

BibApp is a Ruby on Rails application. It is compatible with PostgreSQL, MySQL, and likely most other RDBMS. BibApp also uses Solr for searching.

Reviewer's Analysis

Research networks provide a useful discovery interface for associating research output with people, institutions, and cross-disciplinary groups. Most repository systems are concerned with digital objects, not people, as "first-class" objects. A research network can provide a useful and attractive complementary interface to a repository system.

Considerations and Recommendations

Any or all of these projects merit further review if a repository solution or support is considered. I recommend further analysis that addresses specific GWLA/GPN use cases.

Key Contacts: VIVO: ?, Profiles: ?; BibApp: Sarah Shreeves, UIUC & Jason Stirnaman, KUMC

Sponsors: VIVO: Various partner institutions, DuraSpace; Profiles: Harvard U. ; BibApp: University of Illinois, Urbana-Champaign

Funding: VIVO: NIH grant; Profiles: ; BibApp: University of Illinois, Urbana-Champaign and various institutions contributing development resources

Inception:

Geographic Location: Distributed.

VuFind

<http://vufind.org>

Jason Stirnaman, University of Kansas Medical Center

ENABLERS: Discovery

Brief Description of the Project

VuFind is an open source PHP + Solr application that provides a discovery interface for a Solr index. The goal of *VuFind* is to enable users to search and browse through all of your library's resources by replacing the traditional OPAC to include:

- Catalog Records
- Digital Library Items
- Institutional Repository
- Institutional Bibliography
- Other Library Collections and Resources

VuFind is completely modular, so you can implement just the basic system or all of the components. A wide range of configurable options allows extensive customization without changing any code.

Reviewer's Analysis

VuFind, like Blacklight (see Blacklight review), primarily addresses the facet of data access by providing faceted search and filtering to indexed content and metadata. VuFind adds an elegant and customizable user interface on top of a Solr search index.

VuFind was originally developed specifically to replace clunky vendor online catalog software for libraries and continues to be more oriented toward bibliographic and library collections. VuFind includes tools for harvesting metadata and content as well as importing metadata into Solr. Examples for harvesting and importing various "Open Data" metadata sources can be found at http://vufind.org/wiki/open_data_sources. However, no research data or dataset examples were listed at the time of viewing.

It should be noted that Solr includes robust tools for text analysis and can be employed as a tool for data analysis, statistics, and normalization. For example, see Erik Hatcher's presentation on [rapid prototyping Data.gov with Solr](#). In that sense, VuFind + Solr may address data processing or data analysis needs.

Reviewer's Analysis

VuFind + Solr offer a pure and modern discovery layer for data and metadata, however VuFind's default toolset is more oriented toward repurposing library (MARC) and Dublin Core metadata.

Considerations and Recommendations

Does this project merit further review as part of a site visit, phone interview, or further analysis? (and, if so, what do you recommend?) I recommend further analysis that addresses specific GWLA/GPN use cases for: data discovery, metadata discovery. However, a more generic, less bibliocentric solution like Blacklight may be preferable.

Key Contacts: Demian Katz, Villanova University

Sponsors: Falvey Memorial Library, Villanova University

Funding: VuFind development is funded by Villanova University. Additional technical contributions provided by various institutions.

Inception: 2007

Geographic Location: Distributed.

ARK [Archival Resource Key] Identifiers

ENABLERS: Identifiers

<https://confluence.ucop.edu/display/Curation/ARK>

Susan Matveyeva, Wichita State University

Brief Description of the Project

ARK identifier was developed by the National Library of Medicine and California Digital Library in March 2001. It supports a long-term access for information objects, both tangible and intangible. It is based on the idea that service or curation is the main condition of longevity of digital objects. An ARK is a URL created according to special rules. It answers the questions: who, what, when, and where (e.g. author – title – year – location). “The ARK (Archival Resource Key) naming scheme is designed to facilitate the high-quality and persistent identification of information objects.

A founding principle of the ARK is that persistence is purely a matter of service and is neither inherent in an object nor conferred on it by a particular naming syntax. The best that an identifier can do is to lead users to the services that support robust reference. The term ARK itself refers both to the scheme and to any single identifier that conforms to it. An ARK has five components: [http://NMAH/]ark:/NAAN/Name[Qualifier]

- (1) An optional and mutable Name Mapping Authority Hostport (usually a hostname)
- (2) the "ark:" label
- (3) the Name Assigning Authority Number (NAAN) [Institutional Identity.-S.M.]
- (4) the assigned Name [information object. -S.M.], and
- (5) an optional and possibly mutable Qualifier supported by the NMA.

The NAAN and Name together form the immutable persistent identifier for the object independent of the URL hostname.

An ARK is a special kind of URL that connects users to three things: the named object, its metadata, and the provider's promise about its persistence. When entered into the location field of a Web browser, the ARK leads the user to the named object. That same ARK, inflected by appending a single question mark ('?'), returns a brief metadata record that is both human- and machine-readable. When the ARK is inflected by appending dual question marks ('??'), the returned metadata contains a commitment statement from the current provider. Tools exist for minting, binding, and resolving ARKs.⁵¹

⁵¹ (Kunze, J. and R. Rogers, The Ark Identifier Scheme. May 22, 2008.) See also Kunze's 2012 presentation: <http://www.slideshare.net/jakkbl/the-ark-identifier-scheme-at-ten-years-old>

Reviewer's Analysis

The DOIs are more popular than the ARKs. However, the latter has around 100 registered users including the Internet Archive, Portico, MIT, DCC, the National Library of France, Google and many others. University of Kansas is a registered NAA (Name Assigning Authority) (NAAN: 25031) http://www.cdlib.org/uc3/naan_registry.txt

UC Curation Center provides EZID services for both identifiers: DOI and ARK. ARK is less expensive than DOI and can be hosted locally.

The ARK structure includes Name Assigning Authority (NAA) as prefix; each NAA has a unique number (NAAN). Practically, these numbers are the institutional identifiers. This feature may be useful in consortial environment: (1) NAANs provide an easy efficient way to maintain an institutional identity of member-institutions; (2) the institution (NAAN) and its contribution (unique ID of data) connected as parts of the ARK. "Each organization is identified by a unique NAAN, which can be used as a prefix for the object identifiers that it assigns. For example, CDL assigns ARK identifiers that begin with its NAAN, 13030, as a prefix." (see: Identity Service: Name Assigning Authority Numbers http://www.cdlib.org/services/uc3/naan_table.html). Ability to maintain the inextricable connection between a contributing organization and its contribution may help to overcome some members' anxiety that in a large consortium environment they may be "lost," or be "absorbed", or "disconnected" from their contribution.

Considerations and Recommendations

The discussion about possible advantages of the ARK identifiers in consortial environment might be useful. John Kunze, the expert in ARK, is the best person to talk. http://www.cdlib.org/contact/staff_directory/jkunze.html

Key Contacts: John A. Kunze, California Digital Library Curation Center

Sponsors: University of California Curation Center's IZID service

Funding: UC Curation Center operates EZID on a cost recovery basis.

Inception: 2001

Geographic Location: University of California

DOI (Digital Object Identifier)

ENABLERS: Identifiers

<http://www.doi.org/> (see also: <http://www.datacite.org/whatisdoi>)

Susan Matveyeva, Wichita State University

Brief Description of the Project

From the *DOI handbook*, <http://www.doi.org/hb.html>:

DOI is an acronym for "digital object identifier", meaning a "digital identifier of an object". A DOI name is an identifier (not a location) of an entity on digital networks. It provides a system for persistent and actionable identification and interoperable exchange of managed information on digital networks. A DOI name can be assigned to any entity – physical, digital or abstract – primarily for sharing with an interested user community or managing as intellectual property. The DOI system is designed for interoperability; that is to use, or work with, existing identifier and metadata schemes. DOI names may also be expressed as URLs (URIs). .. The DOI System provides a ready-to-use system of several components: a specified numbering syntax, a resolution service (based on the Handle System), a data model system (including the indices Data Dictionary), and policies and procedures for the implementation of DOI names through a federation of Registration Agencies.

Reviewer's Analysis

DOI is widely used in publishing industry for identification of research articles. In 2005, the German National Library of Science and Technology (TIB) has started to assign DOI to research data, and in 2009, the global consortium DataCite was founded with the goal to manage the DOI system for research data. There are three DOI Registration Agencies in U.S.: California Digital Library, Purdue University Library, and U.S. Dept. of Energy Office of Scientific and Technical Information (OSTI). The University of California Curation Center (UC3) at CDL and Purdue University Libraries offer "DataCite DOIs and other identifiers via the EZID service (<http://n2t.net/ezid>), developed by UC3 to support easy identifier creation and maintenance for educational, non-profit, governmental and commercial clients. .. Through the OSTI Data ID Service, DOIs are assigned to research datasets, and then registered with DataCite to establish persistence. OSTI offers this service for researchers performing U.S. Department of Energy (DOE)-funded research activities carried out at DOE labs and facilities nationwide and grantees at universities and other institutions, as well as to other U.S. federal agencies and thereby other federal government-funded researchers. .. (see: <http://datacite.org/DataCiteUS>).

Data publishers (e.g. data centers, institutional or subject repositories) at first should register for an account with a DataCite member. To obtain an account, data publisher should meet certain requirements (contact@datacite.org -- membership enquires; tech@datacite.org – technical questions). DOI structure: prefix / client ID / unique string of characters.

DOI is an ISO standard assigned to research data. Major data services use DOIs to identify research data. DataCite also offers another persistent identifier: ARK (see a separate review).

Considerations and Recommendations

It is recommended that the regional data service assign DOIs for data it will host, which requires registration with one of the DOI Registration Agency and subscription to the DOI service. For pricing see: <http://n2t.net/eid/home/pricing> .

Currently, DOIs and ARKs are offered as standard identifiers for research data; further analysis may clarify the choice of the most appropriate identifier for the regional data service as well as the possibility to use both identifiers if needed.

Key Contacts: inquiries@us.datacite.org

Sponsors: The International DOI Foundation (IDF)

Funding: IDF is an independent, not-for-profit, open membership organization funded by its members

Inception: 1998

Geographic Location: NA

ORCID [Open Researcher and Contributor ID]

Susan Matveyeva, Wichita State University

<http://orcid.org>

ENABLERS: Identifiers

Brief Description of the Project

ORCID ID is a persistent unique identifier and a method for linking to digital research object. According to the web site:

“ORCID is an open, non-profit, community-based effort to create and maintain a registry of unique researcher identifiers and a transparent method of linking research activities and outputs to these identifiers. ORCID is unique in its ability to reach across disciplines, research sectors, and national boundaries and in its cooperation with other identifier systems. ORCID works with the research community to identify opportunities for integrating ORCID identifiers in key workflows, such as research profile maintenance, manuscript submissions, grant applications, and patent applications. ORCID provides two core functions: (1) a registry to obtain a unique identifier and manage a record of activities, and (2) APIs that support system-to-system communication and authentication. ORCID makes its code available under an open source license, and will post an annual public data file under a CCO waiver for free download. The ORCID Registry is available free of charge to individuals, who may obtain an ORCID, manage their record of activities, and search for others in the Registry. Organizations may become members to link their records to ORCID identifiers, to update ORCID records, to receive updates from ORCID, and to register their employees and students for ORCID identifiers.”

Reviewer’s Analysis

The author name ambiguity is an old problem. There is the number of initiatives that attempt to solve it. For example, VIAF (Virtual International Authority File (<http://viaf.org/>), the national libraries project, currently includes over 12 million authority records; another example is the OCLC WorldCat Identities service (<http://www.oclc.org/developer/services/worldcat-identities>), which uses OpenURL technology to provide information on author’s works and the works about him/her. ResearcherID (<http://www.researcherid.com/>) is the Thomson Reuther service integrated it with the Web of Science. Scopus also offers author ID numbers. In 2012, ISO published ISO 27729 ISNI (http://www.iso.org/iso/home/store/catalogue_tc/catalogue_detail.htm?csnumber=44292) to uniquely identify contributors of media content, such as books, TV programs, and newspaper articles.

The ORCID project is the newest and the most promising method of author identification. It is built on ResearcherID source code that was donated by Thomson Reuther. It is open source and relay on open source philosophy. It has strong support of major publishers, societies, universities, and research community. It is a global in score and discipline agnostic. Researchers can control privacy of their records, as they wish from public, to protect and be completely private. The ORCID ID format (16-digit number) is compatible with the ISO 27729 ISNI. ORCID ID is easy to integrate with other systems⁵² research articles with DOI through CrossRef; the documents that do not have a DOI need to be entered manually. The intent is to use ORCID for grant applications, works of any format, and for research data. ORCID will also link all metrics and relevant sources related to scholarly contributions of a given author, including citations, usage data, etc.

Reviewer's Analysis

Researchers can delegate control of the ongoing management of their profiles to their institutions. This feature is very useful for the future consortium. Another important detail is a business model of the ORCID, which ensure its sustainability. The service is free for individual researchers, but not for organizational members. The size of membership fees will depend on the type and size of the organization, large commercial publishers will pay more than academic institutions. Pay may depend on the number of organizational members; when ORCID will have more members, fee may be decreased. Currently, ORCID offers two membership/subscription categories for organizations, basic and premium, and provides a 20% discount on the annual fee for non-profit organizations. Basic membership: \$5,000 per year. See more at: <http://about.orcid.org/about/membership>

Considerations and Recommendations

My recommendation is to implement ORCID in the GWLA/GPN consortium.

⁵² <http://about.orcid.org/about/community/launch-partners>

Brief Description of the Project

Creative Commons provides licenses, tools, metadata, and best practices that facilitate the sharing of content. CC provides tools for integrating license selection and metadata into asset or content management systems (see http://wiki.creativecommons.org/Web_Integration).

From http://wiki.creativecommons.org/CC0_use_for_data :

(1) We do recommend CC0 for scientific data — and we're thrilled to see CC0 used in other domains, for any content and data, wherever the rights holder wants to make clear such is in the public domain worldwide, to the extent that is possible (note that CC0 includes a permissive fallback license, covering jurisdictions where relinquishment is not thought possible).

(2) However, where CC0 is not desired for whatever reason (business requirements, community wishes, institutional policy...) CC licenses can and should be used for data and databases, right now (as they have been for 8 years) — with the important caveat that CC 3.0 license conditions do not extend to “protect” a database that is otherwise un-copyrighable.

Real world uses of CC for data can be found at <http://wiki.creativecommons.org/Data>

Reviewer's Analysis

Creative Commons addresses the facets of Data Access (sharing and re-use). Reasons to share data include:

- fulfilling funder requirements,
- some journals require data archiving,
- raising interest in the research conducted,
- facilitating and increasing speed of research,
- establishing priority and providing public record.

Facts alone are not copyrighted but their arrangement may be sufficient original expression to merit copyright. For databases, there may be a mix of copyright and data for a research project to consider.[1] Before making a database available under a CC license, a database provider must first make sure she has all rights necessary to do so. Often, the database provider is not the original author of the database contents, which may mean the database provider needs separate permissions from third parties before publishing the database under a CC legal tool.

Also, the database provider must consider what elements of the database she wants to be covered by the CC legal tool and identify those elements in a manner that re-users will see and understand.

A big part of the potential value of data, in particular its society-wide value, is realized by use across organizational boundaries. What are the legal mechanisms for this? Many sites give narrow permission to use data via terms of service. Much ad hoc data sharing occurs among researchers. And increasingly, open data is facilitated by sharing under public terms to manage copyright restrictions that might otherwise limit dissemination or reuse of data, e.g. CC licenses or the CC0 public domain dedication.

The current recommended application of Creative Commons to data and databases is discussed at <http://creativecommons.org/weblog/entry/26283>.

Does this project have implications for a collaborative or consortial approach to defining a regional (GWLA/GPN) strategy for lifecycle management of research data? If so, what does this project offer in terms of approaches that might be incorporated or services that might be acquired? Yes, research data and any inclusive management system should be accompanied by policies, expectations, or licenses for data re-use. Those policies, expectations, or licenses should be clearly expressed and readable by both humans and machines. Most modern content and asset management systems offer some of this functionality by integrating the services at http://wiki.creativecommons.org/Web_Integration

Considerations and Recommendations

This project merits further review if a repository solution or support is considered.

Sources:

1. <http://library.uoregon.edu/datamanagement/ip.html>
2. **Key Contacts:** Puneet Kishor, Project Coordinator for Science and Data, <http://creativecommons.org/staff#puneetkishor>

Sponsors: <https://creativecommons.net/supporters/>

Funding: Non-profit. Private and corporate donors.

Inception: 2001. CC version 1.0 released 2002. Science Commons launched in 2005 and was re-integrated with Creative Commons 2011.

Geographic Location: Distributed. Offices located at MIT.

Brief Description of the Project

The current description from the DDI web site is: "The **Data Documentation Initiative (DDI)** is an effort to create an international standard for describing data from the social, behavioral, and economic sciences. Expressed in [XML](#), the DDI metadata [specification](#) now supports the entire research data life cycle. DDI metadata accompanies and enables data conceptualization, collection, processing, distribution, discovery, analysis, repurposing, and archiving." This is the current description of the next generation of DDI:

The Data Documentation Initiative (DDI) is an international standard for describing data related to the observation and measurement of human activity. With origins in the quantitative social sciences, researchers in other disciplines are increasingly using DDI. The DDI specification is also being used to document other data types, such as social media, biomarkers, administrative data, and transaction data." (from: Developing a Model-Driven DDI Specification)

Reviewer's Analysis

DDI has two separate ongoing branches, DDI Codebook, and DDI Lifecycle. The latter aims to cover the data lifecycle as shown in the figure below. DDI Lifecycle is currently an XML based standard, defined by a set of XML schemas. A future version of DDI Lifecycle will be model-based with bindings into XML, RDF, and relational schemas.



Considerations and Recommendations

DDI is one of the major standards for Social Science research. With growing use by data archives (especially in Europe) and national statistical agencies, I think it is an important standard.

Also note that we will be having the first North American DDI Conference here at KU April 1-3, 2013 <http://www.ipsr.ku.edu/naddi/>

Key Contacts:

Mary Vardigan, ICPSR, DDI Alliance Director

DDI users [listserv](mailto:listserv@ddialliance.org) <http://www.ddialliance.org/community/listserv>

Sponsors: The DDI Alliance <http://www.ddialliance.org/alliance>

Funding: Alliance membership

Inception: The DDI (Data Documentation Initiative) began in the mid-1990s as a project to create a structured metadata standard for the social sciences.

Geographic Location: Worldwide

Brief Description of the Project

The METS schema is a standard for encoding descriptive, administrative, and structural metadata regarding objects within a digital library, expressed using the [XML schema language](#) of the [World Wide Web Consortium](#). The standard is maintained in the [Network Development and MARC Standards Office](#) of the Library of Congress, and is being developed as an initiative of the Digital Library Federation.

Mets includes elements to wrap representations of digital objects, organize content into hierarchical structures, and associate executable behaviors with content. Mets can also carry metadata regarding file groupings, provenance (preservation-related actions), rights, and related parties and their roles. Mets can point to other metadata in external formats.

Reviewer's Analysis

Mets plays a management role for digital objects in an archive.

Considerations and Recommendations

Merits further attention.

Key Contacts:

Network Development and MARC Standards Office
Library of Congress
LS/OPS/NDMSO (4402)
Washington, DC 20540-4402
ndmso@loc.gov
<http://www.loc.gov/help/help-desk.html>

Sponsors: [Network Development and MARC Standards Office](#) of the Library of Congress
ndmso@loc.gov

Funding:

Inception:

Geographic Location:

Metadata Object Description Schema (MODS)

ENABLERS: Metadata

(see also MADS)

<http://www.loc.gov/standards/mods/>

see also: <http://www.loc.gov/standards/mods/registry.php> (implementation registry)

and <http://www.loc.gov/standards/mads/>

Larry Hoyle, University of Kansas, Institute for Policy and Social Research

Brief Description of the Project

Mods is a schema for metadata at the collection and object (item) level. The MODS main page describes it as:

"Metadata Object Description Schema (MODS) is a schema for a bibliographic element set that may be used for a variety of purposes, and particularly for library applications. The standard is maintained by the [Network Development and MARC Standards Office](#) of the Library of Congress with input from users. -- [More about MODS](#)"

A companion schema (MADS – Metadata Authority Description Schema) is described as:

"The Metadata Authority Description Schema (MADS) is an XML schema for an authority element set that may be used to provide metadata about agents (people, organizations), events, and terms (topics, geographics, genres, etc.). MADS serves as a companion to the Metadata Object Description Schema (MODS) to provide metadata about the authoritative entities used in MODS descriptions. The standard is maintained by the MODS/MADS Editorial Committee with the [Network Development and MARC Standards Office](#) of the Library of Congress and input from users. "

Quoting from the Mods Schema Outline (<http://www.loc.gov/standards/mods/mods-outline.html>):

Top Level Elements:

[titleInfo](#)

[note](#)

[name](#)

[subject](#)

[typeOfResource](#)

[classification](#)

[genre](#)

[relatedItem](#)

[originInfo](#)

[identifier](#)

language	location
physicalDescription	accessCondition
abstract	part
tableOfContents	extension
targetAudience	recordInfo

Reviewer's Analysis

The MODS registry lists 34 projects currently using MODS. Some University repositories. More than half of the projects are archiving digitized objects (which might be classed as qualitative, or unstructured data). Examples include sheet music, photographs, digitized physical objects, buildings, and archived web sites. In several cases MODS is used with Fedora or DSpace. Some use MODS within a METS framework. Some use MODS as an intermediate format between other metadata schemas.

Considerations and Recommendations

There may be usable tools built on MODS.

Key Contacts: ndmso@loc.gov

Current MODS/MADS Editorial Committee Membership

- Rebecca Guenther, Library of Congress, Chair
- Jan Ashton, British Library
- Ann Caldwell, Brown University
- Reinhold Heuvelmann, German National Library
- Bill Leonard, Library and Archives Canada
- Sally McCallum, Library of Congress
- Betsy McKelvey, Boston College
- Jon Stroop, Princeton University
- Robin Wendler, Harvard University

Sponsors: Library of Congress

Funding:

Inception:

Geographic Location:

PREservation Metadata: Implementation Strategies
(PREMIS) - <http://www.loc.gov/standards/premis/>
Larry Hoyle

ENABLERS: Metadata

Brief Description of the Project

"The PREMIS Data Dictionary for Preservation Metadata is the international standard for metadata to support the preservation of digital objects and ensure their long-term usability. Developed by an international team of experts, PREMIS is implemented in digital preservation projects around the world, and support for PREMIS is incorporated into a number of commercial and open-source digital preservation tools and systems. The PREMIS Editorial Committee coordinates revisions and implementation of the standard, which consists of the Data Dictionary, an XML schema, and supporting documentation."

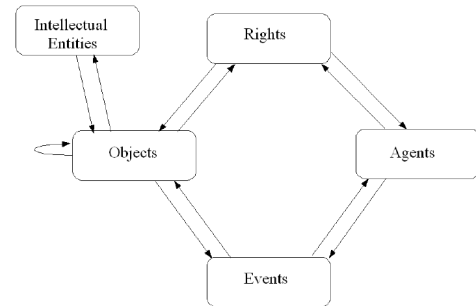


Figure 1: The PREMIS Data Model

Figure 1 at the right is taken from <http://www.loc.gov/standards/premis/v2/premis-2-2.pdf>

Objects in PREMIS are "described as a static set of bits. It is not possible to change a file (or bitstream or representation); one can only create a new file (or bitstream or representation) that is related to the source Object. "

Key Contacts:

Sponsors: Library of Congress

Funding:

Inception: June 2003 OCLC and RLG sponsored the formation of the initial working group in May 2005 with the release of *Data Dictionary for Preservation Metadata: Final Report of the PREMIS Working Group*.

Geographic Location:

Brief Description of the Project

Data Curation Profile is a free online resource for Library and Information Science professionals, Archivists, IT professionals, Data Managers, and others who want information about the specific data generated and used in research areas and sub-disciplines that may be published, shared and preserved for re-use. Data Curation Profiles capture requirements for specific data generated by a single scientist or lab, based on their reported needs and preferences for the data.

Each Data Curation Profile is essentially an outline of the “story” of a data set or collection, describing its origin and lifecycle within a research project. The website includes the *Directory* of completed Data Curation Profiles on a variety of subjects; the downloadable toolkit, and a list of resources. In 2011-2012, the profile developers conducted 12 workshops, funded by the Institute of Museum and Library Services. The Guidelines for Authors⁵³ provides information on definition, structure, and sequence of core and optional modules of the Data Curation Profile. Profile creators are encouraged to submit the completed profiles for publication in the Directory. Publication process and requirements are clearly described. The published profiles are assigned DOI; indexed by Google Scholar, major library discovery tools, and preserved with CLOCKSS and Portico.

A Data Curation Profile is a valuable tool for any data curation project. The profile addresses the data lifecycle and helps data curators to interview researchers, to become familiar with data in different disciplines and subject areas, to identify possible data services, and to plan data curation projects. A Data Curation Profile can be included into a documentation package of regional consortia as a standard tool for data curators working with researchers and planning data curation projects.

Considerations and Recommendations

This project can be included to a list of useful resources for data curators, project managers and librarians.

⁵³ <http://docs.lib.purdue.edu/dcp/guidelines.html>

Key Contacts: Jake Carlson, Associate Professor of Library Science / Data Services Specialist, Purdue University (jcarlso@purdue.edu) and D. Scott Brandt, Associate Dean of Research, Professor of Library Science, Purdue University (techman@purdue.edu)

Sponsoring entities: Purdue University Library; Distributed Data Curation Center; IMLS

Funding Source: The Institute of Museum and Library Services

Inception: 2007

Geographic Location: Purdue University, West Lafayette, Indiana

Brief Description of the Project

The DMPTool is a freely available web service. The primary goals of the tool are to allow researchers to quickly and easily produce a quality data management plan, and to inform researchers of relevant resources and support services across the community and within their institution. Features include:

- Create ready-to-use data management plans for specific funding agencies
- Meet funder requirements for data management plans
- Get step-by-step instructions and guidance for your data management plans as you build it
- In many cases, get data management advice and resources for your specific institution

The tool identifies those elements that specific funders want grant applicants to address, and it allows users to edit, save, share (if desired), print and download their data management plans. [source: CNI program, <http://www.cni.org/pbs/dmptool/>] There is also a video demo (http://dmp.cdlib.org/help/video_demo) that shows how local resources and services appear to researchers using the DMPTool.

The DMPTool site is gaining momentum. The web site reports “October [2012] was our biggest month ever. 375 new users logged into the DMPTool in October, and 336 plans were created. There [are] now a total of 3,466 users, and they’ve created almost 3,000 plans total. 2 more universities customized the DMPTool for their researchers, bringing the total to 28. 65 have configured their campus single-signon for the DMPTool.” There is a map of our participating organizations: <http://bit.ly/L85sKj>

Reviewer’s Analysis

This project addresses the overall need for planning a data management strategy at project inception, not only in response to funder requirements, but also in alignment with institutional resources if the local institution chooses to become a member institution and to customize the tool for its researchers so that the institutions choices and options are reflected in the researchers data management plan. As such, the tool potentially addresses aspects of data management across the full lifecycle of data. The ability for researchers to save, re-use, and to share (if desired) their data management plans through the plan’s repository makes this a particularly advantageous approach for the researcher and possibly for the home institutional.

Specifically, the sections of a DMP-generated plan addresses: 1) Data generated by the project, 2) Period of data retention, 3) Data format and dissemination 4) Data storage and preservation of access, and 5) additional possible data management requirements.

The implication of this project for a collaborative formed by GWLA and GPN might be the need to provide workshops to introduce this tool to member institutions, to train user or customer services staff at institutions work with researchers who are creating in the context of their local institutions, and to perhaps gather information from members about needed enhancements to the tool. Because the tool is hosted, there is no need for additional infrastructure to support the tool outside of member institutions' development of shibboleth.

Only a few GWLA and/or institutions are currently listed as institutions that have customized the tool for their researchers.

Considerations and Recommendations

Perhaps a DMP Tool workshop to future GWLA / GPN member meetings? Compare with the Data Curation Profiles developed by Purdue.

The screenshot shows the DMPTool website interface. At the top, the logo 'DMPTool' is displayed with the tagline 'Guidance and Resources for your Data Management Plan'. The user is logged in as Deborah Ludwig. A navigation menu includes Home, About DMP Tool, DMP News, My Plans, Funder Requirements, and Help. The main content area is titled 'My Data Management Plans' and features a 'Create a new plan:' dropdown menu. The dropdown menu is open, showing a list of funder categories such as 'NSF-EHR: Education and Human Resources', 'Gordon and Betty Moore Foundation', 'Gulf of Mexico Research Initiative', 'Institute of Museum and Library Services', 'National Institutes of Health', 'National Oceanic and Atmospheric Administration', 'NEH-ODH: Office of Digital Humanities', 'NSF-AGS: Atmospheric and Geospace Sciences', 'NSF-AST: Astronomical Sciences', 'NSF-BIO: Biological Sciences', 'NSF-CHE: Chemistry Division', 'NSF-CISE: Computer and Information Science and Engineering', 'NSF-DMR: Materials Research', 'NSF-EAR: Earth Sciences', 'NSF-EFRI: Emerging Frontiers in Research and Innovation', 'NSF-EHR: Education and Human Resources' (highlighted), 'NSF-ENG: Engineering', 'NSF-GEN: Generic', 'NSF-PHY: Physics', and 'NSF-SBE: Social, Behavioral, Economic Sciences'. To the right of the dropdown is a 'Go' button. Below the dropdown, there is a 'Tips' section with advice on exporting and sharing plans, and a 'Recent DMP News' section with links to new articles and demos.

Key contacts: uc3@ucop.edu; Andrew Sallans, Head of Strategic Data Initiatives, Library; Co-Lead on DMPTool Project, University of Virginia; Carly Strasser, Data Curation Specialist, California Digital Library. See <http://www.cni.org/pbs/dmptool/>

Sponsors: DMPTool is a collaborative effort of eight institutions: the California Digital Library, DataONE, the Digital Curation Centre, the Smithsonian Institution, the University of California at Los Angeles Library, the University of California at San Diego Libraries, the University of Illinois at Urbana-Champaign Library and Office of Cyberinfrastructure, and the University of Virginia Library. Other institutions can join this effort by becoming a contributing organization. Member institutions can:

- Enable Shibboleth login (“single sign-on”) as InCommon members:
https://dmp.cdlib.org/help/dmp_shibboleth
- Add links to local resources, help text, suggested answers, contact information

Funding: Part of the funding is from member institutions that contribute.

Inception: 2011

Geographic Location: Located physically on servers at UC3, California Digital Library

CHAPTER II: OUTCOMES OF DATA MEETING (NOT YET AVAILABLE)

This section to be written after Data workshop (May 28, 2013)

CHAPTER III: PLAN FOR DATA MANAGEMENT SUPPORT (NOT YET AVAILABLE)

This section to be written after Data workshop (May 28, 2013)

APPENDIX A: COMMITTEES AND STAFF

Steering Committee

Joni Blake, Executive Director, Greater Western Library Alliance

Paul Farran, Chief of Staff, Information Technology Division, University of Kansas

Judy Ganson, Director of Collection Management Services & Systems, University of Arkansas
Libraries

Deborah Ludwig, Principal Investigator, University of Kansas Libraries

Scott McEathron, Head of Data Initiatives, University of Kansas Libraries

Rick McMullen, Director, Arkansas High Performance Computing Center

Greg Monaco, Director of Research & Cyberinfrastructure, Great Plains Network

Nikki Potter, Project Coordinator, University of Kansas Libraries, Librarian & Archivist, Kansas
Geological Survey

Ann Riley, Associate Director, University of Missouri Libraries

Advisory Committee

Adrian Alexander, Dean of Libraries, University of Tulsa

Carolyn Henderson Allen, Dean of Libraries, University of Arkansas

Gary K. Allen, Vice President for Information Technology and Chief Information Officer,
University of Missouri

Martha Bedard, Dean of University Libraries, University of New Mexico

Dennis Brewer, Associate Vice Chancellor for Information Technology, University of Arkansas

David H. Carlson, Dean, of Libraries, Texas A&M

Deborah Carver, Dean of Libraries, University of Oregon

James Cogswell, Director of the Libraries, University of Missouri

Karen Cole, Director of Dykes Library, University of Kansas Medical Center

Jim Davis, Chief Information Officer, Iowa State University

Loretta Early, Chief of Information, University of Oklahoma

Don Gilstrap, Dean of University Libraries, Wichita State University

Lori Goetsch, Dean of Libraries, Kansas State University

Lorraine Haricombe, Dean of Libraries, University of Kansas

Mary Lou Hines Fritts, Vice Provost & Chief Information Officer, University of Missouri-Kansas
City

Shelia Johnson, Dean of Libraries, Oklahoma State University

Phil Konomos, Associate University Librarian & Chief Technology Officer, Arizona State
University

Bob Lim, Chief Technology Officer, University of Kansas

Bonnie Postlethwaite, Dean of Libraries, University of Missouri-Kansas City

Ken Stafford, Vice Provost & Chief of Information, Kansas State University

Kelli Trosvig, Vice President for University of Washington Information Technology and Chief Information Officer, University of Washington
Johann Van Reenen, Associate Vice President for Research, University of New Mexico
Betsy Wilson, Dean of University Libraries, University of Washington
Melissa Woo, Vice Provost for Information Services & Chief Information Officer, University of Oregon

Grant Staff

Joni Blake, Executive Director, Greater Western Library Alliance
Paul Farran, Chief of Staff, Information Technology Division, University of Kansas
Deborah Ludwig, Principal Investigator, University of Kansas Libraries
Bob Lim, Chief Information Officer, University of Kansas Information Technology Division
Greg Monaco, Director of Research & Cyberinfrastructure, Great Plains Network
Scott McEathron, Head of Data Initiatives, University of Kansas Libraries
Nicole A. Potter, Librarian, Kansas Geological Survey

Research Team A, Literature Review

Judy Ganson, Director of Collection Management Services & Systems, University of Arkansas
Andrew Johnson, Assistant Professor and Metadata Librarian, University of Colorado Boulder
Kathryn Lage, Map Librarian, University of Colorado Boulder
Deborah Ludwig, Assistant Dean, University of Kansas Libraries
Scott McEathron, Head of Data Initiatives, University of Kansas Libraries
Rick McMullen, Director, Arkansas High Performance Computing Center
Amalia Monroe-Gulick, Social Sciences Librarian, University of Kansas Libraries
Sarah Potvin, Metadata Librarian in Digital Services & Scholarly Communication, Texas A&M University Libraries
Ann Riley, Associate Director, University of Missouri Libraries
Brian Westra, Lorry I. Lokey Science Data Services Librarian, University of Oregon
Stephanie Wright, Data Services Coordinator, University of Washington Libraries

Research Team B, Key Projects & Technologies Review

Dan Andresen, Associate Professor, Kansas State University Dept. of Computer Science
Toby Axelsson, Program Director, Information Technology, University of Kansas
Michael Bolton, Director, Digital Initiatives, Texas A&M Libraries
Larry Hoyle, Senior Research Scientist, University of Kansas Institute for Policy & Social Research
Deborah Ludwig, Assistant Dean, University of Kansas Libraries
Susan Matveyeva, Assoc. Professor, Institutional Repository Librarian, Wichita State University
Scott McEathron, Head of Data Initiatives, University of Kansas Libraries
Greg Monaco, Director of Research & Cyberinfrastructure, Great Plains Network
Amalia Monroe-Gulick, Strategy and Assessment Librarian, University of Kansas Libraries
Diane Oerly, Grant Writer, University of Missouri Advanced Computing Environment
Rick McMullen, Director, Arkansas High Performance Computing Center

Jeff Perry, Deputy Technology Officer, University of Kansas
Nicole Potter, Librarian & Archivist, University of Kansas Geological Survey
Jason Stirnaman, Biomedical & Digital Projects Librarian, University of Kansas Medical Center
Jon Wheeler, Lecturer III: University Libraries, University of New Mexico

APPENDIX B: GRANT ABSTRACT

ABSTRACT: Planning for the Lifecycle Management and Long-Term Preservation of Research Data: A Federated Approach

The “data deluge” is a recent but increasingly well-understood phenomenon of scientific and social inquiry.⁵⁴ Large-scale research instruments extend our observational power by many orders of magnitude but at the same time generate massive amounts of data. Researchers work feverishly to document and preserve changing or disappearing habitats, cultures, languages, and artifacts resulting in volumes of media in various formats. New software tools mine a growing universe of historical and modern texts and connect the dots in our semantic environment. Libraries, archives, and museums undertake digitization programs creating broad access to unique cultural heritage resources for research. Global-scale research collaborations with hundreds or thousands of participants, drive the creation of massive amounts of data, most of which cannot be recreated if lost. A recent watershed report, *Harnessing the Power of Digital Data for Science and Society* produced for the National Science and Technology Council summed up the promise and the problems we now face:

*In principle, a digital data deluge can result in rapid progress in science through wider access and the ability to use sophisticated computational and analytical methods and technologies. In practice, the current landscape lacks a comprehensive framework for reliable digital preservation, access, and interoperability, so data are at risk.*⁵⁵

The University of Kansas (KU) Libraries in collaboration with two partners, the Greater Western Library Alliance (GWLA) and the Great Plains Network (GPN), seek a one year (Oct 2012-Sep 2013) IMLS National Leadership Grant designed to leverage collective strengths and create a proposal for a scalable and federated approach to the lifecycle management of research data based on the needs of GPN and GWLA member institutions. Our proposal meets the IMLS strategic goal to “Practice exemplary stewardship of collections and use the power of technology to facilitate discovery of knowledge and cultural heritage.” KU is a public university

⁵⁴Lord, P., A. Macdonald, L. Lyon and D. Giarretta (2004): “From Data Deluge to Data Curation.” *In Proceedings of the UK e-science All Hands meeting 2004*, pp. 371–375

⁵⁵“*Harnessing the Power of Digital Data for Science and Society.*” Report of the Interagency Working Group on Digital Data to the Committee on Science of the National Science and Technology Council. January 2009. www.nitrd.gov/about/harnessing_power_web.pdf

engaged in very high research activity and an active member of both GPN and GWLA. GPN is a consortium of universities in the Midwest that partners to facilitate the use of advanced cyberinfrastructure (network, storage, computation) for research computing. GWLA is a consortium of 31 research libraries. In building on the strength of the partners, the planning process we will focus on these three goals in service to the current and future generations of researchers and scholars:

- **Goal #1:** Undertake an in-depth environmental scan focused on current national and international data management initiatives and on the needs of our member universities for research data management services and infrastructure.
- **Goal #2:** Bring together a GPN and GWLA member forum and two-day workshop for the university research, library, and technology communities focused on understanding challenges and solutions in managing, sharing, and preserving research data.
- **Goal #3:** Create and disseminate a plan for a scalable multi-institutional approach to research data management to support the university members of GPN and GWLA and advance this plan for funding.

APPENDIX C: MAP OF GPN AND GWLA MEMBERSHIP

Geographic Membership, Great Plains Network/Greater Western Library Alliance



The **Greater Western Library Alliance (GWLA)** is a 33-member dynamic and project-oriented consortium of leading research libraries in the Central and Western United States.

The Great Plains Network is a consortium of over 20 universities, primarily research intensive and extensive. Researchers at these institutions participate in projects that require advanced networking and are data- and computing- extensive.⁵⁶

⁵⁶ <http://www.greatplains.net/display/Home/Data+Intensive+Projects+Across+the+GPN+Region>

APPENDIX D: BIBLIOGRAPHY (ALPHABETICAL)

- Akmon, Dharma, Ann Zimmerman, Morgan Daniels, and Margaret Hedstrom. 2011. "The application of archival concepts to a data-intensive environment: working with scientists to understand data management and preservation needs." *Archival Science* no. 11 (3-4):329-348. doi: <http://dx.doi.org/10.1007/s10502-011-9151-4>.
- Allard, Suzie. 2012. "DataONE: Facilitating eScience through Collaboration." *Journal of eScience Librarianship* no. 1 (1):3. doi: <http://dx.doi.org/10.7191/jeslib.2012.1004>.
- Bailey Jr, Charles W. 2012. "Research Data Curation Bibliography." <http://digital-scholarship.org/rdcb/rdcb.htm>.
- Baker, Karen S, and Lynn Yarmey. 2009. "Data stewardship: Environmental data curation and a web-of-repositories." *International Journal of Digital Curation* no. 4 (2):12-27. doi: <http://dx.doi.org/doi:10.2218/ijdc.v4i2.90>.
- Beagrie, Neil, Meg Bellinger, Robin Dale, Marianne Doerr, Margaret Hedstrom, Maggie Jones, Anne Kenney, Catherine Lupovici, Kelly Russell, and Colin Webb. 2002. "Trusted Digital Repositories: Attributes and Responsibilities." *Research Libraries Group & Online Computer Library Center, Report*. <https://www.oclc.org/content/dam/research/activities/trustedrep/repositories.pdf>.
- Beitz, Anthony; Kheeran Dharmawardena; Sam Searle 2012. Research Data Management Strategy and Strategic Plan 2012 – 2015. Monash University. <https://confluence-vre.its.monash.edu.au/display/rdmstrategy/Research+Data+Management+Strategy+and+Strategic+Plan+2012-2015>.
- Berman, Francine. 2008. "Got data? a guide to data preservation in the information age." *Communications of the ACM* no. 51 (12):50. doi: <http://dx.doi.org/10.1145/1409360.1409376>.
- Borgman, Christine L. 2009. "The digital future is now: A call to action for the humanities." *Digital humanities quarterly* no. 3 (4). <http://digitalhumanities.org/dhq/vol/3/4/000077/000077.html>.
- Borgman, Christine L. 2012. "The conundrum of sharing research data." *Journal of the American Society for Information Science and Technology* no. 63 (6):1059-1078. doi: <http://dx.doi.org/10.1002/asi.22634>.
- Brownlee, Rowan. 2009. "Research Data and Repository Metadata: Policy and Technical Issues at the University of Sydney Library." *Cataloging & Classification Quarterly* no. 47 (3-4):370-379. doi: <http://dx.doi.org/10.1080/01639370802714182>.

- Carlson, Jacob and D. Scott Brandt, eds. 2013. *Data Curation Profiles Directory*. Purdue University Libraries. <http://docs.lib.purdue.edu/dcp/>.
- Carlson, Jake. 2012. "Demystifying the data interview. Developing a foundation for reference librarians to talk with researchers about their data." *Reference Services Review* no. 40 (1):7-23. doi: <http://dx.doi.org/10.1108/00907321211203603>.
- Carlson, Jake, and Ruth Kneale. 2011. "Embedded librarianship in the research context Navigating new waters." *College & Research Libraries News* no. 72 (3):167-170. <http://crln.acrl.org/content/72/3/167.short>.
- Chavan, Vishwas, and Peter Ingwersen. 2009. "Towards a data publishing framework for primary biodiversity data: challenges and potentials for the biodiversity informatics community." *BMC Bioinformatics* no. 10 (Suppl 14):S2. <http://www.biomedcentral.com/1471-2105/10/S14/S2>.
- Choudhury, Sayeed, Tim DiLauro, Alex Szalay, Ethan Vishniac, Robert Hanisch, Julie Steffen, Robert Milkey, Teresa Ehling, and Ray Plante. 2008. "Digital data preservation for scholarly publications in astronomy." *International Journal of Digital Curation* no. 2 (2):20-30. doi: <http://dx.doi.org/10.2218/ijdc.v2i2.26>.
- Choudury, Sayeed. 2010. "Data curation: An ecological perspective." *College & Research Libraries News* no. 71 (4):194-196. <http://crln.acrl.org/content/71/4/194.short>.
- Cole, Gareth, Jill Evans, Jessica Gardner, Hannah Lloyd-Jones, and Stephen Trowell. 2013. "The University of Exeter's Roadmap for EPSRC's Research Data Management Expectations." <http://hdl.handle.net/10036/4377>.
- Conway, Esther, Brian Matthews, David Giaretta, Simon Lambert, Michael Wilson, and Nick Draper. 2012. "Managing risks in the preservation of research data with preservation networks." *International Journal of Digital Curation* no. 7 (1):3-15. doi: <http://dx.doi.org/10.2218/ijdc.v7i1.210>.
- Cragin, Melissa H, Carole L Palmer, Jacob R Carlson, and Michael Witt. 2010. "Data sharing, small science and institutional repositories." *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* no. 368 (1926):4023-4038. <http://rsta.royalsocietypublishing.org/content/368/1926/4023.short>.
- Crosas, Merce. 2011. "The dataverse network®: an open-source application for sharing, discovering and preserving data." *D-Lib Magazine* no. 17 (1):2. doi: <http://dx.doi.org/10.1045/january2011-crosas>.
- Dietrich, Dianne, Trisha Adamus, Alison Miner, and Gail Steinhart. 2012. "De-Mystifying the

- Data Management Requirements of Research Funders." *Issues in Science and Technology Librarianship* no. 70 (1). <http://www.istl.org/12-summer/refereed1.html>.
- Dürr, Eugène, Kees van der Meer, Wim Luxemburg, and Ronald Dekker. 2008. "Dataset Preservation for the Long Term: Results of the DareLux Project." *International Journal of Digital Curation* no. 3 (1):29-43. doi: <http://dx.doi.org/10.2218/ijdc.v3i1.40>.
- Faniel, Ixchel M, and Ann Zimmerman. 2011. "Beyond the data deluge: A research agenda for large-scale data sharing and reuse." *International Journal of Digital Curation* no. 6 (1):58-69. doi: <http://dx.doi.org/10.2218/ijdc.v6i1.172>.
- Faundeen, John L. 2003. "The challenge of archiving and preserving remotely sensed data." *Data Science Journal* no. 2:159-163. doi: <http://dx.doi.org/10.2481/dsj.2.159>.
- Ferguson, Jen. 2012. "Description and annotation of biomedical data sets." *Journal of eScience Librarianship* no. 1 (1):9. doi: <http://dx.doi.org/10.7191/jeslib.2012.1000>.
- Flach, Peter and Simon Price. 2012. University of Bristol Research Data Management Principles. Bristol, UK: University of Bristol. <http://data.bris.ac.uk/principles/>.
- Giarlo, Michael J. 2013. "Academic Libraries as Data Quality Hubs." *Journal of Librarianship and Scholarly Communication* no. 1 (3):5. doi: <http://dx.doi.org/10.7710/2162-3309.1059>.
- Gray, Jim, Alexander S Szalay, Ani R Thakar, Christopher Stoughton, and Jan Vandenberg. 2002. Online scientific data curation, publication, and archiving. Paper read at Proceedings of SPIE. doi: <http://dx.doi.org/10.1117/12.461524>.
- Green, Ann G, and Myron P Gutmann. 2007. "Building partnerships among social science researchers, institution-based repositories and domain specific data archives." *OCLC Systems & Services* no. 23 (1):35-53. doi: <http://dx.doi.org/10.1108/10650750710720757>.
- Greenberg, Jane, Hollie C. White, Sarah Carrier, and Ryan Scherle. 2009. "A metadata best practice for a scientific data repository." *Journal of Library Metadata* no. 9 (3-4):194-212. doi: <http://dx.doi.org/10.1080/19386380903405090>.
- Haendel, Melissa A., Nicole A. Vasilevsky, and Jacqueline A. Wirz. 2012. "Dealing with Data: A Case Study on Information and Data Management Literacy." *PLoS Biol* no. 10 (5):e1001339. doi: <http://dx.doi.org/10.1371/journal.pbio.1001339>.
- Hanisch, Robert, and Sayeed Choudhury. 2009. The Data Conservancy: Building a Sustainable System for Interdisciplinary Scientific Data Curation and Preservation. Paper read at Proceedings of the PV 2009 conference, European Space Agency. <https://jscholarship.library.jhu.edu/handle/1774.2/34018>.

- Hey, Anthony J.G., Stewart Tansley, and Kristin Michele Tolle. 2009. *The fourth paradigm: data-intensive scientific discovery*: Microsoft Research Redmond, WA. http://iw.fh-potsdam.de/fileadmin/FB5/Dokumente/forschung/tagungen/i-science/TonyHey_-_eScience_Potsdam_Mar2010_complete.pdf.
- Hey, Anthony J.G., and Anne E. Trefethen. 2003. "The data deluge: An e-science perspective." http://eprints.soton.ac.uk/257648/1/The_Data_Deluge.pdf.
- Higgins, Sarah. 2008. "The DCC curation lifecycle model." *International Journal of Digital Curation* no. 3 (1):134-140. doi: <http://dx.doi.org/10.2218/ijdc.v3i1.48>.
- Holden, John P. 2013. Increasing Access to the Results of Federally Funded Scientific Research. Office of Science and Technology Policy Executive Office of the President. Washington, D.C. http://www.whitehouse.gov/sites/default/files/microsites/ostp/ostp_public_access_memo_2013.pdf.
- Integrated Earth Data Applications (IEDA). *Data Management Plan Tool V.2* Lamont-Doherty Earth Observatory April 17, 2012 2011 [cited May 10, 2013. Available from <http://www.iedadata.org/compliance/plan>.
- Johnston, Lisa, Meghan Lafferty, and Beth Petsan. 2012. "Training Researchers on Data Management: A Scalable, Cross-Disciplinary Approach." *Journal of eScience Librarianship* no. 1 (2):79-87. doi: <http://dx.doi.org/10.7191/jeslib.2012.1012>.
- Jones, Sarah, Graham Pryor, and Angus Whyte. 2013. How to Develop Research Data Management Services - a Guide for HEIs. In *DCC How-to Guides*. Edinburgh, UK: Digital Curation Centre. <http://www.dcc.ac.uk/resources/how-guides>.
- Konkiel, Stacy. (February 22, 2013) "Re: IR metadata schema for ingest of data." Email to Digital Curation Interest Group listserv.
- Kuipers, Tom, and Jeffrey van der Hoeven. 2009. PARSE. Insight: Insight into issues of Permanent Access to the Records of Science in Europe. Survey Report. European Commission. http://www.parse-insight.eu/downloads/PARSE-Insight_D3-4_SurveyReport_final_hq.pdf.
- Lawrence, Bryan, Catherine Jones, Brian Matthews, Sam Pepler, and Sarah Callaghan. 2011. "Citation and peer review of data: Moving towards formal data publication." *International Journal of Digital Curation* no. 6 (2):4-37. doi: <http://dx.doi.org/10.2218/ijdc.v6i2.205>.

- Lawrence, Carolyn J, Qunfeng Dong, Mary L Polacco, Trent E Seigfried, and Volker Brendel. 2004. "MaizeGDB, the community database for maize genetics and genomics." *Nucleic acids research* no. 32 (suppl 1):D393-D397. doi: <http://dx.doi.org/10.1093/nar/gkh011>.
- Lewis, M. J. 2010. "Libraries and the management of research data." In *Envisioning Future Academic Library Services*, 145-168. London Facet Publishing.
<http://eprints.whiterose.ac.uk/11171/>.
- Linkert, Melissa, Curtis T. Rueden, Chris Allan, Jean-Marie Burel, Will Moore, Andrew Patterson, Brian Loranger, Josh Moore, Carlos Neves, Donald MacDonald, Aleksandra Tarkowska, Caitlin Sticco, Emma Hill, Mike Rossner, Kevin W. Eliceiri, and Jason R. Swedlow. 2010. "Metadata matters: access to image data in the real world." *The Journal of Cell Biology* no. 189 (5):777-782. doi: <http://dx.doi.org/10.1083/jcb.201004104>.
- Lynch, Clifford A. 2003. "Institutional repositories: Essential infrastructure for scholarship in the digital age " *ARL Bimonthly Report* no. 226:1-7.
<http://www.arl.org/storage/documents/publications/arl-br-226.pdf>.
- Lynch, Clifford A, and Joan K Lippincott. 2005. "Institutional repository deployment in the United States as of early 2005." *D-lib Magazine* no. 11 (9):1082-9873.
<http://webdoc.sub.gwdg.de/edoc/aw/d-lib/dlib/september05/lynch/09lynch.html>.
- Lyon, Liz. 2012. "The informatics transform: Re-engineering libraries for the data decade." *International Journal of Digital Curation* no. 7 (1):126-138. doi:
<http://dx.doi.org/10.2218/ijdc.v7i1.220>.
- Macdonald, Stuart, and Luis Martinez-Urbe. 2010. "Collaboration to Data Curation: Harnessing Institutional Expertise." *New Review of Academic Librarianship* no. 16:4-16. doi:
<http://dx.doi.org/10.1080/13614533.2010.505823>.
- Marchionini, Gary. 2012. *Research Data Stewardship at UNC: Recommendations for Scholarly Practice and Leadership*. University of North Carolina, Chapel Hill.
http://sils.unc.edu/sites/default/files/general/research/UNC_Research_Data_Stewardship_Report.pdf.
- Matthews, Brian, Shoaib Sufi, Damian Flannery, Laurent Lerusse, Tom Griffin, Michael Gleaves, and Kerstin Kleese. 2010. "Using a Core Scientific Metadata Model in Large-Scale Facilities." *International Journal of Digital Curation* no. 5 (1):106-118. doi:
<http://dx.doi.org/10.2218/ijdc.v5i1.146>.
- Mattmann, Chris, Daniel Crichton, Andrew Hart, Sean Kelly, and Steven Hughes. 2010. "Experiments with Storage and Preservation of NASA's Planetary Data via the Cloud." *IT Professional Magazine* no. 12 (5):28-35. doi: <http://dx.doi.org/10.1109/MITP.2010.97>.

- Mayernik, Matthew. 2011. "Metadata realities for cyberinfrastructure: Data authors as metadata creators." *Available at SSRN 2042653*. doi: <http://dx.doi.org/10.2139/ssrn.2042653>.
- Mayernik, Matthew S. 2010. "Metadata tensions: A case study of library principles vs. everyday scientific data practices." *Proceedings of the American Society for Information Science and Technology* no. 47 (1):1-2. doi: <http://dx.doi.org/10.1002/meet.14504701337>.
- McGarva, Guy, Steve Morris, and Greg Janée. 2008. "Preserving Geospatial Data." *Technology Watch Report, Digital Preservation Coalition (DPC), DPC Technology Watch Series Report: 09-01*. http://www.dpconline.org/component/docman/doc_download/363-preserving-geospatial-data-by-guy-mcgarva-steve-morris-and-greg-janee.
- Michener, William K., Suzie Allard, Amber Budden, Robert B. Cook, Kimberly Douglass, Mike Frame, Steve Kelling, Rebecca Koskela, Carol Tenopir, and David A. Vieglais. 2012. "Participatory design of DataONE – Enabling cyberinfrastructure for the biological and environmental sciences." *Ecological Informatics* no. 11 (0):5-15. doi: <http://dx.doi.org/10.1016/j.ecoinf.2011.08.007>.
- Mooney, Hailey, and Mark P Newton. 2012. "The Anatomy of a Data Citation: Discovery, Reuse, and Credit." *Journal of Librarianship and Scholarly Communication* no. 1 (1):6. <http://jlsccpub.org/jlsc/vol1/iss1/6/>.
- Moore, Richard. 2013. UCSD's Research CyberInfrastructure (RCI) Program: Enabling Research Thru Shared Services. http://escholarship.umassmed.edu/cgi/viewcontent.cgi?article=1058&context=escience_symposium.
- Murray-Rust, Peter, C Neylon, R Pollock, and JT Wilbanks. 2010. "Panton Principles-Principles for Open Data in Science." *Panton Principles*. <http://pantonprinciples.org/>.
- Neuroth, Heike, Felix Lohmeier, and Kathleen Marie Smith. 2011. "TextGrid–Virtual Research Environment for the Humanities." *International Journal of Digital Curation* no. 6 (2):222-231. doi: <http://dx.doi.org/10.2218/ijdc.v6i2.198>.
- Ogburn, Joyce L. 2010. "The Imperative for Data Curation." *Portal : Libraries and the Academy* no. 10 (2):241-246. <http://search.proquest.com/docview/216178815?accountid=14556>.
- Parham, Susan Wells, and Chris Doty. 2012. "NSF DMP Content Analysis: What Are Researchers Saying?" *Bulletin of the American Society for Information Science and Technology* no. 39 (1):37-38.
- Piwowar, Heather A, and Wendy W Chapman. 2008. A review of journal policies for sharing

research data. Paper read at ELPUB2008. <http://es.slideshare.net/hpiwowar/elpub-2008-a-review-of-journal-policies-for-sharing-research-data>.

Poschen, Meik, June Finch, Rob Procter, Mhorag Goff, Mary McDerby, Simon Collins, Jon Besson, Lorraine Beard, and Tom Grahame. 2012. "Development of a Pilot Data Management Infrastructure for Biomedical Researchers at University of Manchester—Approach, Findings, Challenges and Outlook of the MaDAM Project." *International Journal of Digital Curation* no. 7 (2):110-122. doi: <http://dx.doi.org/10.2218/ijdc.v7i2.234>.

Pryor, Graham. 2012. *Managing research data*. London: Facet Publishing.

Reznik-Zellen, Rebecca C., Jessica Adamick, and Stephen McGinty. 2012. "Tiers of Research Data Support Services." *Journal of eScience Librarianship* no. 1 (1):5. doi: <http://dx.doi.org/10.7191/jeslib.2012.1002>.

Rice, Robin. 2009. "DISC-UK DataShare Project: Final Report." <http://repository.jisc.ac.uk/id/eprint/336>.

Sallans, Andrew, and Martin Donnelly. 2012. "DMP Online and DMPTool: Different Strategies Towards a Shared Goal." *International Journal of Digital Curation* no. 7 (2):123-129. doi: <http://dx.doi.org/10.2218/ijdc.v7i2.235>.

Schaeffer, Mary L, Lisa C Harper, Jack M Gardiner, Carson M Andorf, Darwin A Campbell, Ethalinda KS Cannon, Taner Z Sen, and Carolyn J Lawrence. 2011. "MaizeGDB: curation and outreach go hand-in-hand." *Database: the journal of biological databases and curation* no. 2011. doi: <http://dx.doi.org/10.1093/database/bar022>.

Shaon, Arif, and Andrew Woolf. 2011. "Long-term Preservation for Spatial Data Infrastructures: a Metadata Framework and Geo-portal Implementation." *D-Lib Magazine* no. 17 (9):1. doi: <http://dx.doi.org/10.1045/september2011-shaon>.

Sharpe, Robert, and Martin Waller. 2006. *Mind the gap: assessing digital preservation needs in the UK*. York: Digital Preservation Coalition. <http://www.dpconline.org/docs/reports/uknamindthegap.pdf>.

Simeone, Michael, Jennifer Guiliano, Rob Kooper, and Peter Bajcsy. 2011. "Digging into data using new collaborative infrastructures supporting humanities-based computer science research." *First Monday* no. 16 (5-2). doi: <http://dx.doi.org/10.5210%2Ffm.v16i5.3372>

Simons, Natasha. 2012. "Implementing DOIs for Research Data." *D-Lib Magazine* no. 18 (5):1. doi: <http://dx.doi.org/10.1045/may2012-simons>.

Spielmann, Katherine A, and Keith W Kintigh. 2011. "The Digital Archaeological Record." *SAA*

- Archaeological Record*.
http://alexandriaarchive.org/bonecommons/archive/files/spielmann_kintigh_icz_saa_jan2011_8025c1e7b5.pdf.
- Starr, Joan, and Angela Gastl. 2011. "iscitedby: A metadata scheme for datacite." *D-Lib Magazine* no. 17 (1):9. doi: <http://dx.doi.org/doi:10.1045/january2011-starr>.
- Steinhart, Gail. 2011. "DataStaR: A Data Sharing and Publication Infrastructure to Support Research [Article and Abstract]." *Agricultural Information Worldwide* no. 4 (1):16-20. <http://journals.sfu.ca/iaald/index.php/aginfo/article/view/199/156>.
- Swan, Alma, and Sheridan Brown. 2008. *To share or not to share: Publication and quality assurance of research data outputs: Main report*. London: Research Information Network. <http://www.rin.ac.uk/system/files/attachments/To-share-data-outputs-report.pdf> .
- Tenopir, Carol, Suzie Allard, Kimberly Douglass, Arsev Umur Aydinoglu, Lei Wu, Eleanor Read, Maribeth Manoff, and Mike Frame. 2011. "Data Sharing by Scientists: Practices and Perceptions." *PLoS ONE* no. 6 (6):1-21. doi: <http://dx.doi.org/10.1371/journal.pone.0021101>.
- Tenopir, Carol, Ben Birch, and Suzie Allard. 2012. Academic Libraries and Research Data Services: Current Practices and Plans for the Future. http://0-www.ala.org.catalog.wlib.org/acrl/sites/ala.org.acrl/files/content/publications/whitepapers/Tenopir_Birch_Allard.pdf.
- Treloar, Andrew. 2009. "Design and Implementation of the Australian National Data Service." *International Journal of Digital Curation* no. 4 (1):125-137. doi: <http://dx.doi.org/10.2218/ijdc.v4i1.83>.
- Treloar, Andrew, David Groenewegen, and Catherine Harboe-Ree. 2007. "The data curation continuum: Managing data objects in institutional repositories." *D-Lib Magazine* no. 13 (9):4. <http://www.dlib.org/dlib/september07/treloar/09treloar.html>.
- University of Edinburgh, Information Services RDM Policy Implementation Committee. 2012. Research Data Management (RDM) Roadmap August 2012 – January 2014. http://www.ed.ac.uk/polopoly_fs/1.101223!/fileManager/UoE-RDM-Roadmap201121102.pdf.
- University of Wisconsin-Madison. 2012. Electronic Lab Notebook Pilot at the University of Wisconsin-Madison: Study Findings. University of Wisconsin-Madison. http://academictech.doit.wisc.edu/files/ELN_pilot_report_UWMadison.pdf.
- Vision, Todd J. 2010. "Open data and the social contract of scientific publishing." *BioScience* no.

- 60 (5):330-330. doi: <http://dx.doi.org/10.1525/bio.2010.60.5.2>.
- Wallis, Jullian C, Christine L Borgman, Matthew S Mayernik, and Alberto Pepe. 2008. "Moving archival practices upstream: An exploration of the life cycle of ecological sensing data in collaborative field research." *International Journal of Digital Curation* no. 3 (1):114-126. doi: <http://dx.doi.org/10.2218/ijdc.v3i1.46>.
- Warner, Simeon, Jeroen Bekaert, Carl Lagoze, Xiaoming Liu, Sandy Payette, and Herbert Warner. 2007. "Pathways: augmenting interoperability across scholarly repositories." *International Journal on Digital Libraries* no. 7 (1/2):35-52. doi: <http://dx.doi.org/10.1007/s00799-007-0016-7>.
- Westra, Brian, Marisa Ramirez, Susan Wells Parham, and Jeanine Marie Scaramozzino. 2010. "Selected Internet Resources on Digital Research Data Curation." *Issues in Science and Technology Librarianship* (63):9. <http://www.istl.org/10-fall/internet2.html>.
- Whyte, Angus. 2013. The sweet smell of sustainability - JISC MRD projects make the business case. In *Digital Curation Centre Blog*. Edinburgh: Digital Curation Centre. <http://www.dcc.ac.uk/blog/sweet-smell-sustainability-jisc-mrd-projects-make-business-case>.
- Williams, Robin, and Graham Pryor. 2009. Patterns of information use and exchange: case studies of researchers in the life sciences. London: Research Information Network and the British Library. http://www.publishingresearch.net/documents/RINPatterns_information_use-REPORT_Nov2009.pdf.
- Williford, Christa, Charles J Henry, and Amy Friedlander. 2012. *One Culture: Computationally Intensive Research in the Humanities and Social Sciences: a Report on the Experiences of First Respondents to the Digging Into Data Challenge*. Washington: Council on Library and Information Resources. <http://www.clir.org/pubs/reports/pub151>.
- Willis, Craig, Jane Greenberg, and Hollie White. 2012. "Analysis and synthesis of metadata goals for scientific data." *Journal of the American Society for Information Science and Technology* no. 63 (8):1505-1520. doi: <http://dx.doi.org/10.1002/asi.22683>.
- Wilson, Andrew. 2010. "How Much Is Enough: Metadata for Preserving Digital Data." *Journal of Library Metadata* no. 10 (2-3):205-217. doi: <http://dx.doi.org/10.1080/19386389.2010.506395>.
- Witt, Michael. 2012. "Co-designing, Co-developing, and Co-implementing an Institutional Data Repository Service." *Journal of Library Administration* no. 52 (2):172-188. doi: <http://dx.doi.org/10.1080/01930826.2012.655607>.

Witt, Michael, and Mike Giarlo. 2012. "Databib." *Libraries Faculty and Staff Presentations* (Paper 1).
http://docs.lib.purdue.edu/lib_fspress/1/.

Witt, Michael, and Mike Giarlo. 2012. "Databib: An Online Bibliography of Research Data Repositories." <http://0-www.ala.org.catalog.wlib.org/lita/sites/ala.org.lita/files/content/conferences/forum/2012/PosterSessionDescriptions.pdf>.

Yakel, Elizabeth. 2007. "Digital curation." *OCLC Systems & Services* no. 23 (4):335-340. doi:
<http://dx.doi.org/10.1108/10650750710831466>

APPENDIX E: BIBLIOGRAPHY (THEMATIC)

1. Introduction and general works

- Bailey Jr, Charles W. 2012. "Research Data Curation Bibliography." <http://digital-scholarship.org/rdc/rdc.htm>
- Borgman, Christine L. 2009. "The digital future is now: A call to action for the humanities." *Digital humanities quarterly* no. 3 (4).
<http://digitalhumanities.org/dhq/vol/3/4/000077/000077.html> /000077.html
- Gray, Jim, Alexander S Szalay, Ani R Thakar, Christopher Stoughton, and Jan Vandenberg. 2002. Online scientific data curation, publication, and archiving. Paper read at Proceedings of SPIE. doi: <http://dx.doi.org/10.1117/12.461524>
- Hey, Anthony JG, Stewart Tansley, and Kristin Michele Tolle. 2009. The fourth paradigm: data-intensive scientific discovery: Microsoft Research Redmond, WA. http://iw.fh-potsdam.de/fileadmin/FB5/Dokumente/forschung/tagungen/i-science/TonyHey_-_eScience_Potsdam_Mar2010_complete_.pdf
- Hey, Anthony JG, and Anne E Trefethen. 2003. "The data deluge: An e-science perspective." http://eprints.soton.ac.uk/257648/1/The_Data_Deluge.pdf
- Higgins, Sarah. 2008. "The DCC curation lifecycle model." *International Journal of Digital Curation* no. 3 (1):134-140. doi: <http://dx.doi.org/10.2218/ijdc.v3i1.48>
- Ogburn, Joyce L. 2010. "The Imperative for Data Curation." *Portal : Libraries and the Academy* no. 10 (2):241-246. <http://search.proquest.com/docview/216178815?accountid=14556>
- Pryor, Graham. 2012. *Managing research data*. London: Facet Publishing.
- Westra, Brian, Marisa Ramirez, Susan Wells Parham, and Jeanine Marie Scaramozzino. 2010. "Selected Internet Resources on Digital Research Data Curation." *Issues in Science and Technology Librarianship* (63):9. <http://www.istl.org/10-fall/internet2.html>
- Witt, Michael, and Mike Giarlo. 2012. "Databib." *Libraries Faculty and Staff Presentations* (Paper 1). <http://0-www.ala.org.catalog.wbilib.org/lita/sites/ala.org.lita/files/content/conferences/forum/2012/PosterSessionDescriptions.pdf>
- Witt, Michael, and Mike Giarlo. 2012. "Databib: An Online Bibliography of Research Data Repositories." http://docs.lib.purdue.edu/lib_fspress/1/

Yakel, Elizabeth. 2007. "Digital curation." *OCLC Systems & Services* no. 23 (4):335-340. doi: <http://dx.doi.org/10.1108/10650750710831466>

2. Assessments of researcher behavior, attitudes and needs

Carlson, Jacob and D. Scott Brandt, eds. . *Data Curation Profiles Directory*. Purdue University Libraries 2013 [cited May 10, 2013. Available from <http://docs.lib.purdue.edu/dcp/>

Feijen, Martin. 2011. "What researchers want." *Surf Foundation*.
http://www.surf.nl/nl/publicaties/Documents/What_researchers_want.pdf

Kuipers, Tom, and Jeffrey van der Hoeven. 2009. PARSE. Insight: Insight into issues of Permanent Access to the Records of Science in Europe. Survey Report. European Commission. http://www.parse-insight.eu/downloads/PARSE-Insight_D3-4_SurveyReport_final_hq.pdf

Marchionini, Gary. 2012. Research Data Stewardship at UNC: Recommendations for Scholarly Practice and Leadership. University of North Carolina, Chapel Hill.
http://sil.unc.edu/sites/default/files/general/research/UNC_Research_Data_Stewardship_Report.pdf

Scaramozzino, Jeanine Marie, Marisa L. Ramírez, and Karen J. McGaughey. 2012. "A Study of Faculty Data Curation Behaviors and Attitudes at a Teaching-Centered University." *College & Research Libraries* no. 73 (4):349-365. <http://crl.acrl.org/content/73/4/349.abstract>

Sharpe, Robert, and Martin Waller. 2006. *Mind the gap: assessing digital preservation needs in the UK*. York: Digital Preservation Coalition.
<http://www.dpconline.org/docs/reports/uknamindthegap.pdf>

Swan, Alma, and Sheridan Brown. 2008. *To share or not to share: Publication and quality assurance of research data outputs: Main report*. London: Research Information Network.
<http://www.rin.ac.uk/system/files/attachments/To-share-data-outputs-report.pdf>

Tenopir, Carol, Suzie Allard, Kimberly Douglass, Arsev Umur Aydinoglu, Lei Wu, Eleanor Read, Maribeth Manoff, and Mike Frame. 2011. "Data Sharing by Scientists: Practices and Perceptions." *PLoS ONE* no. 6 (6):1-21. doi: <http://dx.doi.org/10.1371/journal.pone.0021101>

Williams, Robin, and Graham Pryor. 2009. Patterns of information use and exchange: case studies of researchers in the life sciences. London: Research Information Network and the British Library.

http://www.publishingresearch.net/documents/RINPatterns_information_use-REPORT_Nov2009.pdf

3. Services, roles and responsibilities

- Beitz, Anthony; Kheeran Dharmawardena; Sam Searle 2012. Research Data Management Strategy and Strategic Plan 2012 – 2015. Monash University. <https://confluence-its.monash.edu.au/display/rdmstrategy/Research+Data+Management+Strategy+and+Strategic+Plan+2012-2015>
- Choudury, Sayeed. 2010. "Data curation: An ecological perspective." *College & Research Libraries News* no. 71 (4):194-196. <http://crln.acrl.org/content/71/4/194.short>
- Cole, Gareth, Jill Evans, Jessica Gardner, Hannah Lloyd-Jones, and Stephen Trowell. 2013. "The University of Exeter's Roadmap for EPSRC's Research Data Management Expectations." <http://hdl.handle.net/10036/4377>
- Flach, Peter and Simon Price. 2012. University of Bristol Research Data Management Principles. Bristol, UK: University of Bristol. <http://data.bris.ac.uk/principles/>
- Giarlo, Michael J. 2013. "Academic Libraries as Data Quality Hubs." *Journal of Librarianship and Scholarly Communication* no. 1 (3):5. doi: <http://dx.doi.org/10.7710/2162-3309.1059>.
- Haendel, Melissa A., Nicole A. Vasilevsky, and Jacqueline A. Wirz. 2012. "Dealing with Data: A Case Study on Information and Data Management Literacy." *PLoS Biol* no. 10 (5):e1001339. doi: <http://dx.doi.org/10.1371/journal.pbio.1001339>.
- Jones, Sarah, Graham Pryor, and Angus Whyte. 2013. How to Develop Research Data Management Services - a Guide for HEIs. In *DCC How-to Guides*. Edinburgh, UK: Digital Curation Centre. <http://www.dcc.ac.uk/resources/how-guides>
- Lewis, M. J. 2010. "Libraries and the management of research data." In *Envisioning Future Academic Library Services*, 145-168. London Facet Publishing. <http://eprints.whiterose.ac.uk/11171/>
- Linkert, Melissa, Curtis T. Rueden, Chris Allan, Jean-Marie Burel, Will Moore, Andrew Patterson, Brian Loranger, Josh Moore, Carlos Neves, Donald MacDonald, Aleksandra Tarkowska, Caitlin Sticco, Emma Hill, Mike Rossner, Kevin W. Eliceiri, and Jason R. Swedlow. 2010. "Metadata matters: access to image data in the real world." *The Journal of Cell Biology* no. 189 (5):777-782. doi: <http://dx.doi.org/10.1083/jcb.201004104>.

- Lyon, Liz. 2012. "The informatics transform: Re-engineering libraries for the data decade." *International Journal of Digital Curation* no. 7 (1):126-138. doi: <http://dx.doi.org/10.2218/ijdc.v7i1.220>.
- Macdonald, Stuart, and Luis Martinez-Urbe. 2010. "Collaboration to Data Curation: Harnessing Institutional Expertise." *New Review of Academic Librarianship* no. 16:4-16. doi: <http://dx.doi.org/10.1080/13614533.2010.505823>.
- Marchionini, Gary. 2012. Research Data Stewardship at UNC: Recommendations for Scholarly Practice and Leadership. University of North Carolina, Chapel Hill. http://sils.unc.edu/sites/default/files/general/research/UNC_Research_Data_Stewardship_Report.pdf
- Moore, Richard. 2013. UCSD's Research CyberInfrastructure (RCI) Program: Enabling Research Thru Shared Services. http://escholarship.umassmed.edu/cgi/viewcontent.cgi?article=1058&context=escience_symposium
- Poschen, Meik, June Finch, Rob Procter, Mhorag Goff, Mary McDerby, Simon Collins, Jon Besson, Lorraine Beard, and Tom Grahame. 2012. "Development of a Pilot Data Management Infrastructure for Biomedical Researchers at University of Manchester—Approach, Findings, Challenges and Outlook of the MaDAM Project." *International Journal of Digital Curation* no. 7 (2):110-122. doi: <http://dx.doi.org/10.2218/ijdc.v7i2.234>.
- Reznik-Zellen, Rebecca C., Jessica Adamick, and Stephen McGinty. 2012. "Tiers of Research Data Support Services." *Journal of eScience Librarianship* no. 1 (1):5. doi: <http://dx.doi.org/10.7191/jeslib.2012.1002>.
- University of Edinburgh, Information Services RDM Policy Implementation Committee. 2012. Research Data Management (RDM) Roadmap August 2012 – January 2014. http://www.ed.ac.uk/polopoly_fs/1.101223!/fileManager/UoE-RDM-Roadmap201121102.pdf
- Whyte, Angus. 2013. The sweet smell of sustainability - JISC MRD projects make the business case. In *Digital Curation Centre Blog*. Edinburgh: Digital Curation Centre. <http://www.dcc.ac.uk/blog/sweet-smell-sustainability-jisc-mrd-projects-make-business-case>
- University of Wisconsin-Madison. 2012. Electronic Lab Notebook Pilot at the University of Wisconsin-Madison: Study Findings. Madison: University of Wisconsin. http://academictech.doit.wisc.edu/files/ELN_pilot_report_UWMadison.pdf

Witt, Michael. 2012. "Co-designing, Co-developing, and Co-implementing an Institutional Data Repository Service." *Journal of Library Administration* no. 52 (2):172-188. doi: <http://dx.doi.org/10.1080/01930826.2012.655607>.

4. Sharing, reuse, publication and citation

Borgman, Christine L. 2012. "The conundrum of sharing research data." *Journal of the American Society for Information Science and Technology* no. 63 (6):1059-1078. doi: <http://dx.doi.org/10.1002/asi.22634>.

Chavan, Vishwas, and Peter Ingwersen. 2009. "Towards a data publishing framework for primary biodiversity data: challenges and potentials for the biodiversity informatics community." *BMC Bioinformatics* no. 10 (Suppl 14):S2. doi: <http://dx.doi.org/10.1186/1471-2105-10-S14-S2>

Cragin, Melissa H, Carole L Palmer, Jacob R Carlson, and Michael Witt. 2010. "Data sharing, small science and institutional repositories." *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* no. 368 (1926):4023-4038. <http://rsta.royalsocietypublishing.org/content/368/1926/4023.short>

Faniel, Ixchel M, and Ann Zimmerman. 2011. "Beyond the data deluge: A research agenda for large-scale data sharing and reuse." *International Journal of Digital Curation* no. 6 (1):58-69. <http://ijdc.net/index.php/ijdc/article/viewFile/163/231>

Holden, John P. 2013. Increasing Access to the Results of Federally Funded Scientific Research. edited by Office of Science and Technology Policy Executive Office of the President. Washington, D.C. http://www.whitehouse.gov/sites/default/files/microsites/ostp/ostp_public_access_memo_2013.pdf

Lawrence, Bryan, Catherine Jones, Brian Matthews, Sam Pepler, and Sarah Callaghan. 2011. "Citation and peer review of data: Moving towards formal data publication." *International Journal of Digital Curation* no. 6 (2):4-37. doi: <http://dx.doi.org/10.2218/ijdc.v6i2.205>.

Mooney, Hailey, and Mark P Newton. 2012. "The Anatomy of a Data Citation: Discovery, Reuse, and Credit." *Journal of Librarianship and Scholarly Communication* no. 1 (1):6. <http://jlscc-pub.org/jlsc/vol1/iss1/6/>

Murray-Rust, Peter, C Neylon, R Pollock, and JT Wilbanks. 2010. "Panton Principles-Principles for Open Data in Science." *Panton Principles*. <http://pantonprinciples.org/>

Piwowar, Heather A, and Wendy W Chapman. 2008. A review of journal policies for sharing research data. Paper read at ELPUB2008. <http://es.slideshare.net/hpiowar/elpub-2008-a-review-of-journal-policies-for-sharing-research-data>

Simons, Natasha. 2012. "Implementing DOIs for Research Data." *D-Lib Magazine* no. 18 (5):1. doi: <http://dx.doi.org/10.1045/may2012-simons>.

Vision, Todd J. 2010. "Open data and the social contract of scientific publishing." *BioScience* no. 60 (5):330-330. doi: <http://dx.doi.org/10.1525/bio.2010.60.5.2>.

5. Data management planning

Carlson, Jake. 2012. "Demystifying the data interview. Developing a foundation for reference librarians to talk with researchers about their data." *Reference Services Review* no. 40 (1):7-23. doi: <http://dx.doi.org/10.1108/00907321211203603>.

Carlson, Jake, and Ruth Kneale. 2011. "Embedded librarianship in the research context Navigating new waters." *College & Research Libraries News* no. 72 (3):167-170. <http://crln.acrl.org/content/72/3/167.short>.

Dietrich, Dianne, Trisha Adamus, Alison Miner, and Gail Steinhart. 2012. "De-Mystifying the Data Management Requirements of Research Funders." *Issues in Science and Technology Librarianship* no. 70 (1). <http://www.istl.org/12-summer/refereed1.html>.

Holden, John P. 2013. Increasing Access to the Results of Federally Funded Scientific Research. Office of Science and Technology Policy Executive Office of the President. Washington, D.C. http://www.whitehouse.gov/sites/default/files/microsites/ostp/ostp_public_access_memo_2013.pdf.

Integrated Earth Data Applications (IEDA). *Data Management Plan Tool V.2* Lamont-Doherty Earth Observatory April 17, 2012 2011 [cited May 10, 2013. Available from <http://www.iedadata.org/compliance/plan>.

Johnston, Lisa, Meghan Lafferty, and Beth Petsan. 2012. "Training Researchers on Data Management: A Scalable, Cross-Disciplinary Approach." *Journal of eScience Librarianship* no. 1 (2):79-87. doi: <http://dx.doi.org/10.7191/jeslib.2012.1012>.

Parham, Susan Wells, and Chris Doty. 2012. "NSF DMP Content Analysis: What Are Researchers Saying?" *Bulletin of the American Society for Information Science and Technology*

no. 39 (1):37-38.

Sallans, Andrew, and Martin Donnelly. 2012. "DMP Online and DMPTool: Different Strategies Towards a Shared Goal." *International Journal of Digital Curation* no. 7 (2):123-129. doi: <http://dx.doi.org/10.2218/ijdc.v7i2.235>.

Tenopir, Carol, Ben Birch, and Suzie Allard. 2012. Academic Libraries and Research Data Services: Current Practices and Plans for the Future. http://0-www.ala.org.catalog.wlib.org/acrl/sites/ala.org.acrl/files/content/publications/whitepapers/Tenopir_Birch_Allard.pdf.

University of Wisconsin-Madison. 2012. Electronic Lab Notebook Pilot at the University of Wisconsin-Madison: Study Findings. University of Wisconsin-Madison. http://academictech.doit.wisc.edu/files/ELN_pilot_report_UWMadison.pdf.

6. Policies and standards [pending]

7. Institutional repositories, approaches, and issues

Konkiel, Stacy. (February 22, 2013) "Re: IR metadata schema for ingest of data." Email to Digital Curation Interest Group listserv.

Lynch, Clifford A. 2003. "Institutional repositories: Essential infrastructure for scholarship in the digital age " *ARL Bimonthly Report* no. 226:1-7. <http://www.arl.org/storage/documents/publications/arl-br-226.pdf>.

Lynch, Clifford A, and Joan K Lippincott. 2005. "Institutional repository deployment in the United States as of early 2005." *D-lib Magazine* no. 11 (9):1082-9873. <http://webdoc.sub.gwdg.de/edoc/aw/d-lib/dlib/september05/lynch/09lynch.html>.

Macdonald, Stuart, and Luis Martinez-Uribe. 2010. "Collaboration to Data Curation: Harnessing Institutional Expertise." *New Review of Academic Librarianship* no. 16:4-16. doi: <http://dx.doi.org/10.1080/13614533.2010.505823>.

Marchionini, Gary, et al. 2012. Research Data Stewardship at UNC: Recommendations for Scholarly Practice and Leadership. University of North Carolina, Chapel Hill. http://sil.unc.edu/sites/default/files/general/research/UNC_Research_Data_Stewardship_Report.pdf.

Moore, Carol; et al. 2009. The Research Library's Role in Digital Repository Services: Final Report of the ARL Digital Repository Issues Task Force. Washington: Association of

Research Libraries. <http://www.arl.org/storage/documents/publications/repository-services-report-jan09.pdf>

Rice, Robin. 2009. "DISC-UK DataShare Project: Final Report." <http://repository.jisc.ac.uk/id/eprint/336>.

Steinhart, Gail. 2011. "DataStaR: A Data Sharing and Publication Infrastructure to Support Research." *Agricultural Information Worldwide* no. 4 (1):16-20. <http://journals.sfu.ca/iaald/index.php/aginfo/article/view/199/156>.

Treloar, Andrew, David Groenewegen, and Cathrine Harboe-Ree. 2007. "The data curation continuum: Managing data objects in institutional repositories." *D-Lib Magazine* no. 13 (9):4. <http://www.dlib.org/dlib/september07/treloar/09treloar.html>

8. Disciplinary or subject repositories, approaches, and issues

9. Federated approaches

Allard, Suzie. 2012. "DataONE: Facilitating eScience through Collaboration." *Journal of eScience Librarianship* no. 1 (1):3. doi: <http://dx.doi.org/10.7191/jeslib.2012.1004>.

Baker, Karen S, and Lynn Yarmey. 2009. "Data stewardship: Environmental data curation and a web-of-repositories." *International Journal of Digital Curation* no. 4 (2):12-27. doi: <http://dx.doi.org/doi:10.2218/ijdc.v4i2.90>.

Crosas, Merce. 2011. "The dataverse network®: an open-source application for sharing, discovering and preserving data." *D-Lib Magazine* no. 17 (1):2. doi: <http://dx.doi.org/10.1045/january2011-crosas>.

Green, Ann G, and Myron P Gutmann. 2007. "Building partnerships among social science researchers, institution-based repositories and domain specific data archives." *OCLC Systems & Services* no. 23 (1):35-53. doi: <http://dx.doi.org/10.1108/10650750710720757>.

Hanisch, Robert, and Sayeed Choudhury. 2009. The Data Conservancy: Building a Sustainable System for Interdisciplinary Scientific Data Curation and Preservation. Paper read at Proceedings of the PV 2009 conference, European Space Agency. <http://jscholarship.library.jhu.edu/handle/1774.2/34018>.

Lawrence, Carolyn J, Qunfeng Dong, Mary L Polacco, Trent E Seigfried, and Volker Brendel. 2004. "MaizeGDB, the community database for maize genetics and genomics." *Nucleic acids research* no. 32 (suppl 1):D393-D397. doi: <http://dx.doi.org/10.1093/nar/gkh011>.

- Marchionini, Gary, et al. 2012. Research Data Stewardship at UNC: Recommendations for Scholarly Practice and Leadership. University of North Carolina, Chapel Hill. http://sils.unc.edu/sites/default/files/general/research/UNC_Research_Data_Stewardship_Report.pdf.
- Michener, William K., Suzie Allard, Amber Budden, Robert B. Cook, Kimberly Douglass, Mike Frame, Steve Kelling, Rebecca Koskela, Carol Tenopir, and David A. Vieglais. 2012. "Participatory design of DataONE—Enabling cyberinfrastructure for the biological and environmental sciences." *Ecological Informatics* no. 11 (0):5-15. doi: <http://dx.doi.org/10.1016/j.ecoinf.2011.08.007>.
- Neuroth, Heike, Felix Lohmeier, and Kathleen Marie Smith. 2011. "TextGrid—Virtual Research Environment for the Humanities." *International Journal of Digital Curation* no. 6 (2):222-231. doi: <http://dx.doi.org/10.2218/ijdc.v6i2.198>.
- Schaeffer, Mary L, Lisa C Harper, Jack M Gardiner, Carson M Andorf, Darwin A Campbell, Ethalinda KS Cannon, Taner Z Sen, and Carolyn J Lawrence. 2011. "MaizeGDB: curation and outreach go hand-in-hand." *Database: the journal of biological databases and curation* no. 2011. doi: <http://dx.doi.org/10.1093/database/bar022>.
- Simeone, Michael, Jennifer Guiliano, Rob Kooper, and Peter Bajcsy. 2011. "Digging into data using new collaborative infrastructures supporting humanities-based computer science research." *First Monday* no. 16 (5-2). doi: <http://dx.doi.org/10.5210%2Ffm.v16i5.3372>.
- Spielmann, Katherine A, and Keith W Kintigh. 2011. "The Digital Archaeological Record." *SAA Archaeological Record*. http://alexandriaarchive.org/bonecommons/archive/files/spielmann_kintigh_icz_saa_jan2011_8025c1e7b5.pdf.
- Treloar, Andrew. 2009. "Design and Implementation of the Australian National Data Service." *International Journal of Digital Curation* no. 4 (1):125-137. doi: <http://dx.doi.org/10.2218/ijdc.v4i1.83>.
- Warner, Simeon, Jeroen Bekaert, Carl Lagoze, Xiaoming Liu, Sandy Payette, and Herbert Warner. 2007. "Pathways: augmenting interoperability across scholarly repositories." *International Journal on Digital Libraries* no. 7 (1/2):35-52. doi: <http://dx.doi.org/10.1007/s00799-007-0016-7>.
- Williford, Christa, Charles J. Henry, and Amy Friedlander. 2012. *One Culture: Computationally Intensive Research in the Humanities and Social Sciences: a Report on the Experiences of First Respondents to the Digging Into Data Challenge*. Washington: Council on Library and Information Resources. <http://www.clir.org/pubs/reports/pub151>.

10. Economics/costs of data curation

11. Archiving and preservation

- Akmon, Dharma, Ann Zimmerman, Morgan Daniels, and Margaret Hedstrom. 2011. "The application of archival concepts to a data-intensive environment: working with scientists to understand data management and preservation needs." *Archival Science* no. 11 (3-4):329-348. doi: <http://dx.doi.org/10.1007/s10502-011-9151-4>.
- Beagrie, Neil, Meg Bellinger, Robin Dale, Marianne Doerr, Margaret Hedstrom, Maggie Jones, Anne Kenney, Catherine Lupovici, Kelly Russell, and Colin Webb. 2002. "Trusted Digital Repositories: Attributes and Responsibilities." *Research Libraries Group & Online Computer Library Center, Report*.
<https://www.oclc.org/content/dam/research/activities/trustedrep/repositories.pdf>.
- Berman, Francine. 2008. "Got data? a guide to data preservation in the information age." *Communications of the ACM* no. 51 (12):50. doi: <http://dx.doi.org/10.1145/1409360.1409376>.
- Choudhury, Sayeed, Tim DiLauro, Alex Szalay, Ethan Vishniac, Robert Hanisch, Julie Steffen, Robert Milkey, Teresa Ehling, and Ray Plante. 2008. "Digital data preservation for scholarly publications in astronomy." *International Journal of Digital Curation* no. 2 (2):20-30. doi: <http://dx.doi.org/10.2218/ijdc.v2i2.26>.
- Conway, Esther, Brian Matthews, David Giaretta, Simon Lambert, Michael Wilson, and Nick Draper. 2012. "Managing risks in the preservation of research data with preservation networks." *International Journal of Digital Curation* no. 7 (1):3-15. doi: <http://dx.doi.org/10.2218/ijdc.v7i1.210>.
- Dürr, Eugène, Kees van der Meer, Wim Luxemburg, and Ronald Dekker. 2008. "Dataset Preservation for the Long Term: Results of the DareLux Project." *International Journal of Digital Curation* no. 3 (1):29-43. doi: <http://dx.doi.org/10.2218/ijdc.v3i1.40>.
- Faundeen, John L. 2003. "The challenge of archiving and preserving remotely sensed data." *Data Science Journal* no. 2:159-163. doi: <http://dx.doi.org/10.2481/dsj.2.159>.
- Mattmann, Chris, Daniel Crichton, Andrew Hart, Sean Kelly, and Steven Hughes. 2010. "Experiments with Storage and Preservation of NASA's Planetary Data via the Cloud." *IT Professional Magazine* no. 12 (5):28-35. doi: <http://dx.doi.org/10.1109/MITP.2010.97>.
- McGarva, Guy, Steve Morris, and Greg Janée. 2008. "Preserving Geospatial Data." *Technology Watch Report, Digital Preservation Coalition (DPC), DPC Technology Watch Series Report: 09-01*. http://www.dpconline.org/component/docman/doc_download/363-preserving-

[geospatial-data-by-guy-mcgarva-steve-morris-and-gred-greg-janee.](#)

Tenopir, Carol, Suzie Allard, Kimberly Douglass, Arsev Umur Aydinoglu, Lei Wu, Eleanor Read, Maribeth Manoff, and Mike Frame. 2011. "Data Sharing by Scientists: Practices and Perceptions." *PLoS ONE* no. 6 (6):1-21. doi: <http://dx.doi.org/10.1371/journal.pone.0021101>.

Tenopir, Carol, Ben Birch, and Suzie Allard. 2012. Academic Libraries and Research Data Services: Current Practices and Plans for the Future. http://0-www.ala.org.catalog.wlib.org/acrl/sites/ala.org.acrl/files/content/publications/whitepapers/Tenopir_Birch_Allard.pdf.

Wallis, Jullian C, Christine L Borgman, Matthew S Mayernik, and Alberto Pepe. 2008. "Moving archival practices upstream: An exploration of the life cycle of ecological sensing data in collaborative field research." *International Journal of Digital Curation* no. 3 (1):114-126. doi: <http://dx.doi.org/10.2218/ijdc.v3i1.46>.

12. Metadata and description

Brownlee, Rowan. 2009. "Research Data and Repository Metadata: Policy and Technical Issues at the University of Sydney Library." *Cataloging & Classification Quarterly* no. 47 (3-4):370-379. doi: <http://dx.doi.org/10.1080/01639370802714182>.

Ferguson, Jen. 2012. "Description and annotation of biomedical data sets." *Journal of eScience Librarianship* no. 1 (1):9. doi: <http://dx.doi.org/10.7191/jeslib.2012.1000>.

Greenberg, Jane, Hollie C. White, Sarah Carrier, and Ryan Scherle. 2009. "A metadata best practice for a scientific data repository." *Journal of Library Metadata* no. 9 (3-4):194-212. doi: <http://dx.doi.org/10.1080/19386380903405090>.

Matthews, Brian, Shoaib Sufi, Damian Flannery, Laurent Lerusse, Tom Griffin, Michael Gleaves, and Kerstin Kleese. 2010. "Using a Core Scientific Metadata Model in Large-Scale Facilities." *International Journal of Digital Curation* no. 5 (1):106-118. doi: <http://dx.doi.org/10.2218/ijdc.v5i1.146>.

Mayernik, Matthew. 2011. "Metadata realities for cyberinfrastructure: Data authors as metadata creators." Available at SSRN 2042653. doi: <http://dx.doi.org/10.2139/ssrn.2042653>.

Mayernik, Matthew S. 2010. "Metadata tensions: A case study of library principles vs. everyday scientific data practices." *Proceedings of the American Society for Information Science and Technology* no. 47 (1):1-2. doi: <http://dx.doi.org/10.1002/meet.14504701337>.

Shaon, Arif, and Andrew Woolf. 2011. "Long-term Preservation for Spatial Data Infrastructures:

- a Metadata Framework and Geo-portal Implementation." *D-Lib Magazine* no. 17 (9):1. doi: <http://dx.doi.org/10.1045/september2011-shaon>.
- Starr, Joan, and Angela Gastl. 2011. "iscitedby: A metadata scheme for datacite." *D-Lib Magazine* no. 17 (1):9. doi: <http://dx.doi.org/doi:10.1045/january2011-starr>.
- Willis, Craig, Jane Greenberg, and Hollie White. 2012. "Analysis and synthesis of metadata goals for scientific data." *Journal of the American Society for Information Science and Technology* no. 63 (8):1505-1520. doi: <http://dx.doi.org/10.1002/asi.22683>.
- Wilson, Andrew. 2010. "How Much Is Enough: Metadata for Preserving Digital Data." *Journal of Library Metadata* no. 10 (2-3):205-217. doi: <http://dx.doi.org/10.1080/19386389.2010.506395>.

APPENDIX F: FORM FOR REVIEW OF KEY PROJECTS & TECHNOLOGIES

TEMPLATE FOR CYBERINFRASTRUCTURE PROJECT or SPECIFIC TECHNOLOGY REVIEW (maybe 2 pages for each report?)

Project Name:

Web Address:

Key Contacts:

Sponsoring entities:

Funding Source, if Known:

Inception Date, if Known:

Geographic Location, if Applicable:

Brief Description of the Project (1 paragraph):

Analysis:

How does the project address, or potentially address, key facets of lifecycle management? [Data Creation, Data Processing, Data Analysis, Data Preservation, Data Access for Others (data sharing), and Data Reuse]

Does this project have implications for a collaborative or consortial approach to defining a regional (GWLA/GPN) strategy for lifecycle management of research data? If so, what does this project offer in terms of approaches that might be incorporated or services that might be acquired?

Recommendations:

Does this project merit further review as part of a site visit, phone interview, or further analysis? (and, if so, what do you recommend?)

Reviewer: [Name, Organization]

APPENDIX G: SURVEY QUESTIONS



IMLS - KU GWLA GPN Survey

Lifecycle Management of Research Data

This study surveys member institutions of the Great Plains Network and the Greater Western Library Alliance about services related to the management of research data. We are sending this survey to Deans and Directors of Libraries, to Chief Research Officers, and to Chief Information Officers. We welcome three individual responses from each member institution. Alternatively, one representative may be designated to respond on behalf of the institution.

The management of research data by institutions is an issue of growing interest and complexity for universities. Large-scale research instruments extend our observational power by many orders of magnitude but at the same time generate massive amounts of data. Researchers work feverishly to document and preserve changing or disappearing habitats, cultures, languages, and artifacts resulting in volumes of media in various formats. New software tools mine a growing universe of historical and modern texts and connect the dots in our semantic environment.

Research data management includes services that extend throughout the lifecycle of data such as organizing data; providing metadata; making data accessible for discovery and sharing; storing, archiving, and/or preserving data for long term retention and use; creation of data management plans; development of institutional policies for the appropriate management of data; training in best practices for managing data; and developing security or privacy provisions.

The purpose of the study is to assess what services GWLA and GPN institutions provide to help researchers manage research data and, within each institution, to discover which units or divisions are responsible for providing specific resources and services. We will use this information to inform an advisory council of university leaders from among the GWLA and GPN universities as they consider opportunities to leverage collective strengths and institute scalable and shared approaches to the management of research data.

This study is conducted and funded through a grant from the Institute for Museum and Library Studies and conducted on behalf of members of the Greater Western Library Alliance (GWLA), a consortium of 33 research libraries and the Great Plains Network, (GPN) a consortium of 20+ universities partnering to facilitate the use of advanced cyberinfrastructure. The University of Kansas serves as the home institution for the IMLS grant. Our project web site can be found at: <http://imls.gwla.org>.

This survey will take approximately 20 – 30 minutes of your time. We thank you in advance for helping us conduct this research by providing data from your institution. If you have any questions about this survey, please contact the principle investigator Deborah Ludwig (dludwig@ku.edu) or investigator Scott McEathron, (macmap68.ku.edu).



If you have any questions or concerns regarding this survey, please email dludwig@ku.edu.



IMLS - KU GWLA GPN Survey

Lifecycle Management of Research Data

University Name

Person completing the survey (name and position)

Organizational Unit / Department

Email address of person completing the survey

For which are you filling out the survey?

- I am responding on behalf of my organizational unit
- I am responding on behalf of my entire institution

If you have any questions or concerns regarding this survey, please email dludwig@ku.edu.



IMLS - KU GWLA GPN Survey

Lifecycle Management of Research Data

Some universities provide services to support the use and management of research data. Please indicate by checking any or all boxes that apply what research data services are available at your institution and who is providing that service.

General Support and Services for Research Data

	University Library	Central IT	University Office of Research	Department or Research Center	Other	Not Offered	Don't Know
2.01) Provides assistance to researchers in developing data management plans for research proposals	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.02) Offers formal training (e.g. workshops or courses) for researchers on data management practices	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.03) Consults with researchers on data management best practices (data archiving, data sharing, metadata, etc.)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.04) Dedicates funding resources (staffing, facilities, and services) that support long-term management of research data generated by campus researchers	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.05) Ensures university compliance for research data in accordance with commercial licenses, government regulations, and funding agency mandates	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.06) Consults with researchers on options for data licensing agreements for open or restricted access	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Storage, Archiving, Preservation, and Sharing of Data

	University Library	Central IT	University Office of Research	Department or Research Center	Other	Not Offered	Don't Know
2.07) Helps researchers decide which data are important to preserve for long-term access.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.08) Advises researchers on data and or metadata standards for research data or datasets	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.09) Prepares metadata for researchers to enhance the discovery of their research data.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.10) Provides standards and methods for de-identifying sensitive data	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.11) Provides short term networked data storage for researchers (5 years or less)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.12) Provides long term networked data storage for researchers (more than 5 years)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.13) Provides a repository on site to store metadata and data together	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.14) Provides a repository for sharing data with appropriate access controls	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.15) Provides assistance with identifying national or international data repositories for archiving research data	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.16) Provides guidance to researchers on offsite repositories for data and metadata deposit	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.17) Provides ongoing support for discovery, citation, and usability of data over the long term	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.18) Helps prepare data / data sets for deposit into a repository	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>



IMLS - KU GWLA GPN Survey

Lifecycle Management of Research Data

Some universities provide services to support the use and management of research data. Please indicate by checking any or all boxes that apply what research data services are available at your institution and who is providing that service.

Accessing and Using Research Data

	University Library	Central IT	University Office of Research	Department or Research Center	Other	Not Offered	Don't Know
2.19) Helps researchers locate sources of data	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.20) Supports the assignment of persistent identifiers to research data sets. [Digital Object Identifiers, or DOIs are one example.]	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.21) Supports the assignment of persistent identifiers for researchers. [ORCID ID's are one example]	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.22) Supports linking research data to research publications based on the research data	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.23) Provides support for the analysis of data	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.24) Provides support for visualization of data (for example, simulations, geographic information systems or GIS, or statistical visualization	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

If you have any questions or concerns regarding this survey, please email dludwig@ku.edu.



IMLS - KU GWLA GPN Survey

Lifecycle Management of Research Data

Some universities provide services to support the use and management of research data. Please indicate by checking any or all boxes that apply what research data services are available at your institution and who is providing that service.

High Performance Computing

	University Library	Central IT	University Office of Research	Department or Research Center	Other	Not Offered	Don't Know
2.25) Provides a secure data facility for research data with access controls, backup and restore facilities meeting regulatory standards such as HIPAA	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.26) Provides computing facilities for computationally intensive research, e.g. massively parallel technology	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.27) Provides computing facilities for in-place analysis of extremely large research datasets	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2.28) Provides assistance in securing advanced or high performance computing resources located off-campus.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

If you indicated above that "Other" entities provided one or more of these services, please list those other entities here.

If you have any questions or concerns regarding this survey, please email dludwig@ku.edu.



IMLS - KU GWLA GPN Survey

Lifecycle Management of Research Data

A helpful model for thinking about the different kinds of information and services needed for managing data is the Johns Hopkins University stack model for Data management described by Sayeed Choudhury. In this hierarchical model, layers depend on the layers listed below them. Preservation includes ensuring that there is enough information about the creation of the data that the data can be interpreted and reused without communication with its creators. Curation activities may include developing enhanced ability to find elements of the data or link data to other collections.

Data Management Layers

Layers	Characteristics	Sample Actions	Researcher Implications	NSF Implications
Curation	Active and ongoing management of data through its lifecycle of interest and usefulness	<ul style="list-style-type: none"> • Provide ongoing bibliographic control for data. • Link research data to publications based on the data. • Provide tools for further analysis. • Harvest metadata for the data to share with external search engines. 	<ul style="list-style-type: none"> • Feature Extraction • New query capabilities • Cross-disciplinary accessibility 	<ul style="list-style-type: none"> • Offers competitive advantage • New opportunities for data use
Preservation	Ensures that archived data can be fully used and interpreted over time	<ul style="list-style-type: none"> • Add information to maintain the <i>viability, render-ability, and understandability</i> of data long term. • Monitor format obsolescence; migrate data to new digital formats as need. • Preserve tools and/or documentation for using the data. 	<ul style="list-style-type: none"> • Ability to use own data in the future (e.g. 5 years) • Data sharing with others 	<ul style="list-style-type: none"> • Satisfies NSF needs across directorates
Archiving	Data protection is applied to stored data, including fixity checking, and assignment of data identifiers	<ul style="list-style-type: none"> • Check for viruses as data is deposited • Establish checksum snapshots over time to ensure data has not change. • Assign a persistent identifier such as a DOI or handle. • Link metadata to data 	<ul style="list-style-type: none"> • Data is better protected. • Provides persistent identifiers for locating, sharing, and referencing data. 	<ul style="list-style-type: none"> • Could satisfy most NSF requirements
Storage	Bits on disk, tape, cloud etc. Backup and restore.	<ul style="list-style-type: none"> • Place data in networked storage • Invoke backups. • Set access protocols. 	Responsibilities for <ul style="list-style-type: none"> • Restore • Sharing • Staffing 	<ul style="list-style-type: none"> • Could be enough for now, but not near-term future

Adapted from the "stack model" of Data Management Layers under development by John Hopkins University and based on the definition of data curation advanced by the University of Illinois Graduate School of Library and Information Science. See: <http://www.clir.org/initiatives-partnerships/data-curation>



IMLS - KU GWLA GPN Survey

Lifecycle Management of Research Data

For which of the following types of digital and non-digital research data are support services in place to manage and preserve them for long-term access your campus? (please check all levels of the stack model that apply)

[Click here to show the stack model description in a separate tab](#)

	Curation	Preservation	Archiving	Storage	None	Don't know
4.02) Digital texts or digital copies of texts and manuscripts	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.05) Digital images or digital copies of images	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.07) Digital audio recordings	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.09) Digital video recordings	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.10) Spreadsheets	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.11) Digital Databases [or Data Sets] (e.g. surveys, census data, government statistics, etc)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.12) Computer code	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.13) Hardware or research equipment	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.15) Spatial data	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.16) Digital gene sequences or similar digital renderings of biological/organic/inorganic samples or specimens	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.17) Artistic products	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4.18) Other (Please Specify)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Please describe the other type of data selected above:

If you have any questions or concerns regarding this survey, please email dludwig@ku.edu.



IMLS - KU GWLA GPN Survey

Lifecycle Management of Research Data

Policies For Research Data Management

Does your university have a general policy on ownership of research data?

- Yes
- No
- Don't Know

If yes, please provide a reference to that policy.

Has your university established a policy for funding to cover the costs of research data management and storage

..... for externally funded research?

- Costs are primarily expected to be folded into the direct costs of those grants and contracts
- Costs are primarily paid for by the researcher's department or center
- Costs are primarily paid for by the University from research overhead (F&A funds) or other sources
- Other (Please Specify Below)

- The University has no policy
- Don't know

... for research NOT supported by grants, contracts, or other external sources of funding?

- Costs are primarily paid for by the individual researcher or research team
- Costs are primarily paid for by the researcher's department or center
- Costs are primarily paid for by the University from overhead (F&A funds) or other sources
- Other (Please Specify Below)

- The University has no policy
- Don't know

... once external funding has expired?

- Costs are primarily paid for by the individual researcher or research team
- Costs are primarily paid for by the researcher's department or center
- Costs are primarily paid for by the University from overhead (F&A funds) or other sources
- Other (Please Specify Below)
- The University has no policy
- Don't know

What are the main institutional challenges for working with research data that your organization has identified?

What services and/or future plans has your organization developed to meet the challenges listed above?

If you have any questions or concerns regarding this survey, please email djudwig@ku.edu.



IMLS - KU GWLA GPN Survey

Lifecycle Management of Research Data

We thank you for your time spent taking this survey.
Your response has been recorded.

If you have any questions or concerns regarding this survey, please email dludwig@ku.edu.

APPENDIX H: GLOSSARY OF TERMS AND CONCEPTS

One challenge of forging a path forward is the abstract nature of many of the terms and concepts associated with the use and management of digital data. Because various professions and subject domains also have different understandings and usage of these same terms, it is important to state clear definitions to support shared understanding. A shared understanding is essential to forge a cooperative approach to support research and establish a path forward.

Data: an abstract term that forms the lowest level of abstraction from which information and then knowledge are derived; may be structured or unstructured; digital or analog; factual numbers, words, images, etc., accepted as they stand that are often used as a basis for reasoning, discussion, or calculation.⁵⁷

Research data: facts, observations or experiences on which an argument, theory or test is based. Data may be numerical, descriptive or visual. Data may be raw or analyzed, experimental or observational. Data includes: laboratory notebooks; field notebooks; primary research data (including research data in hardcopy or in computer readable form); questionnaires; audiotapes; videotapes; models; photographs; films; test responses. Research collections may include slides; artifacts; specimens; samples. Provenance information about the data might also be included: the how, when, where it was collected and with what (for example, instrument). The software code used to generate, annotate or analyze the data may also be included.⁵⁸

Data archiving or digital archiving: The library and archiving communities often use it interchangeably with *digital preservation* (see below). Computing professionals tend to use digital archiving to mean the process of backup and ongoing maintenance as opposed to strategies for long-term digital preservation.

Data management plans (DMP): the documentation of research data management practices and any responsibilities such as university policies, ethics, intellectual property, attribution etc. Its purpose is to ensure the quality of your research data and outputs, integrity and repeatability, appropriate access to data, and appropriate reuse of data for subsequent research. A DMP may be created for a department, a project or collaboration. The responsibility of

⁵⁷ Wikipedia; <http://www.merriam-webster.com/dictionary>

⁵⁸ Courtesy of The University of Melbourne, <https://policy.unimelb.edu.au/MPF1242>

implementing and following the DMP lies with the involved researchers, IT managers and data managers.⁵⁹

Digital data curation: “*digital curation*” or “*data curation*” often used interchangeably: the process of establishing and developing long term repositories of digital assets for current and future research. Steps include selection, description (via metadata), maintenance, preservation, and providing access.

Digital data management: or “*data management*,” the processes of creating, organizing, and making accessible and preserving digital research data (may include conventions for naming and structuring files and folders, version control, backing up of data; and metadata documentation of provenance).

Digital preservation: refers to the series of managed activities necessary to ensure continued access to digital materials for as long as necessary. Digital preservation is defined very broadly for the purposes of this study and refers to all of the actions required to maintain access to digital materials beyond the limits of media failure or technological change.⁶⁰

E-Research: a broader term than “*e-Science*,” research that utilizes digital technology within the research process, including sciences, social sciences and humanities.

E-Science: an area of scientific research characterized by intensive use of computing infrastructure, highly networked environments and vast amounts of digital data.

Metadata: loosely define as *data about data*; data or information about one or more aspects of the data content. For example, card catalogs of libraries are a form of metadata.

⁵⁹ Courtesy of The University of Melbourne; <https://policy.unimelb.edu.au/MPF1242>

⁶⁰Digital Preservation Coalition (2008). "Introduction: Definitions and Concepts". Digital Preservation Handbook. York, UK. Retrieved 3 December 2012
<http://www.dpconline.org/advice/preservationhandbook/introduction/definitions-and-concepts>