

Data Curation Workshop 16 September 2011

Suzie Allard Ph.D. University of Tennessee, USA Edward A. Fox, Ph.D. Virginia Tech University, USA Lucia Lötter, Ph.D. Human Sciences Research Council, South Africa Lynn Woolfrey University of Cape Town, South Africa









UNIVERSITEIT STELLENBOSCH UNIVERSITY



Drs. Jegede & Habib's Talks



Welcome

- Data has many definitions. But an unifying theme is that it is research output that is not codified in journal articles, monographs etc.
- Data curation is a world-wide initiative since we recognize the importance of the primary data as well as the products that come out of our data.

Introductions

The Data Lifecycle



JISC http://www.jisc.ac.uk/whatwedo/campaigns/res3/jischelp.aspx







copyright University Corporation for Atmospheric Research, Photo by Lee Klinger



copyright University Corporation for Atmospheric Research



The researcher follows the research design to collect the data.

For example, in the environmental sciences this may be observations from individual researchers or from monitoring instruments.





The researcher reviews the data to assure its quality.

This could take many forms since data may be recorded in a variety of ways ranging from notebooks to vast amounts of electronic observations.

This step is essential to providing reliable, accurate data.





DataSNE

This step is essential to helping discovery.





The researcher or designee deposits the data in a repository.

Often this is simply another computer in the researcher's office.

It is best if it is a curated repository that will help the researcher's data in the next two steps – preservation & discovery







Data needs to be preserved whether it is analog or digital.

Preservation challenges include data loss, scattered data sources, data deluge, poor data practices.







Image courtesy of DataONE

The data are "discovered" by another researcher looking for data to assist his research.

This step may also be an opportunity for data to cross disciplinary lines.







Data from different collections are discovered and integrated around one research question.





Integrated data are analyzed and new knowledge is created.





Data may not go through all the steps in the data lifecycle.

For example data may be preserved and then be analyzed by a researcher who helped collected it.

This analysis may inform the next round of data collection.





DataONEpedia

Data Files Best Practices

Data & Metadata tools

Role of Data Curation

• Lucia Lötter

Data Sharing

• Lucia Lötter

Data Management & Workflows

• Lucia Lötter

Case Study: Virginia Tech

Presented by Edward A. Fox, Ph.D

fox@vt.edu .

What is Data Citation?

- bibliographic means to reference a data set
- Allows for proper attribution in a scholarly work.

Just as we cite papers, so too should we cite data sets

Four Outcomes from Data Citation

- Enables finding the data & repeating the researcher's analysis.
- Enhances visibility of the data producer & promotes scholarly recognition for their research.
- Increases visibility of the data repository → may increase use of the data holdings.
- Leads to increased data sharing supporting scientific discovery and the scientific enterprise.

Find Data, Repeat Research

- Scientific method founded on inquiry based on collecting observable and measurable evidence.
- identifying a question, doing background research, constructing a hypothesis, testing the hypothesis with data, and analyzing and reporting the results.
- Verification: Repeatable and allow for predicting future results.
- Full disclosure.

Data citation plays an important role in this verification process, as citation provides scientists with a link to the data used in the testing and analysis steps, which allows scientists to repeat the analysis, and verify the results.

Attribution

- The lack of data citation best practices is a hindrance to the scientific community, since it limits the ability for credit attribution (Altman & King, 2007).
- The need for data citation is becoming widely recognized in the scientific community.
- Even publishers are becoming active partners in encouraging a data archiving policy that includes preserving data sets and linking them to articles.

Attribution is an essential to encouraging scientists to share data. This is an important step towards following the scientific process.

Intellectual Property

- Gives credit to data producers and data publishers.
- Provides a link from the traditional literature to the data, enhancing the intellectual property link.
- Gives intellectual legitimacy to the creation of data → important scientific product worthy of being stored and cited.

Incentive to Publish, Share, Preserve

- Increases the incentive to publish a data set.
- Facilitates sharing of the data.
- Supports curation and preservation activities.

The outcome is that data is more likely to be accessed and re-used, a condition which leads to increased scholarly productivity from an initial data product.



Ideal Data Citation

- describe any data set, database, or data file.
- provide for all levels of granularity (table, row, cell).
- be used for any snapshot (version, e.g., in time)
- be formatted for any view: XML, HTML, CSV, etc.
- have (or not have) annotations
- link to older, newer, and latest versions
- provide actionability ("Click-through")
- provide persistence (validity into the future)
- be machine readable, thus allowing for automatic parsing

Exemplar Citations

On-Line Data Set

 Turner, D.P., W.D.Ritts, and M. Gregory. 2006. BigFoot NPP Surfaces for North and South American Sites, 2002-2004. Data set. Available on-line [<u>http://daac.ornl.gov</u>] from Oak Ridge National Laboratory Distributed Active Archive Center, Oak Ridge, Tennessee, U.S.A.doi:10.3334/ORNLDAAC/750.

Subset of larger Data Set

 Oak Ridge National Laboratory Distributed Active Archive Center (ORNL DAAC). 2009. MODIS subsetted land products, Collection 5. Available on-line [<u>http://daac.ornl.gov/MODIS/modis.html</u>] from ORNL DAAC, Oak Ridge, Tennessee, U.S.A. Accessed November 20, 2009.

Online Map

Oak Ridge National Laboratory Distributed Active Archive Center (ORNL DAAC).
 2009. FLUXNET Network Map. Available online
 [http://www.fluxnet.ornl.gov/fluxnet/Maps/Political_fluxnet_networks_cropped_s
 mall_april2009.png] from ORNL DAAC, Oak Ridge, Tennessee, U.S.A.

ORNL DAAC citation examples (available at <u>http://daac.ornl.gov/citation_policy.html</u>

Establishing a Community of Practice

• Lynn Woolfrey



Some CoP in Education





SciData

Funded by:













What you can do..

- Plan for how your institution will handle data
- Help make your researchers away of the importance of preserving their data
- Teach data users how to properly cite data sets they are using
- Network with colleagues to share experiences and expertise regarding data curation
- Make use of national & international resources about data citation.

Questions?

A Resource

Best Practices	Software Tools
Search Best Practices Contains	Search Tools Contains SEARCH
Best Practice Categories	Tool Categories
All Best Practices	<u>All Tools</u>
Content and Structure (15) Data Access and Discovery (5) Data Documentation (3) Data Preservation and Archives (7)	Analysis & Modeling (34) Data Aquisition & Modeling (30) Workflow (7)
Planning Policies and Governance (1) Quality Assurance and Quality Control (5)	Featured Tool
Vocabulary Standards and Services (2)	<u>S-PLUS (S+)</u> Primary Category: Analysis & Modeling
Featured Best Practice	S-PLUS is a commercial implementation of the S statistical programming language
Define the contents of Data Files Category: Data Documentation	that was developed by Bell Labs. S+ has a cross-platform integrated development environment (IDE) , provides the ability to analyze gigabyte class data sets on the desktop, and a package system for deployment of analytics.
Define any coded values	Cost: Cost-basis
Quality flags or qualifying valuesDefine missing values	TIBCO Spotfire S+

www.dataone.org/dataonepedia Editors: Cook, Michener Contributors: Best practices workshop participants