STARS

University of Central Florida
STARS

Faculty Scholarship and Creative Works

11-19-2014

Data documentation & metadata

Sai Deng University of Central Florida, sai.deng@ucf.edu

Part of the Cataloging and Metadata Commons, and the Scholarly Communication Commons Find similar works at: https://stars.library.ucf.edu/ucfscholar University of Central Florida Libraries http://library.ucf.edu

This Other Presentation is brought to you for free and open access by STARS. It has been accepted for inclusion in Faculty Scholarship and Creative Works by an authorized administrator of STARS. For more information, please contact STARS@ucf.edu.

Original Citation

Deng, S. (2014). Data documentation and metadata. University of Central Florida graduate students library research workshop: Publishing in the Academy.



UCF Libraries Research Workshop

Data Documentation &Metadata



Sai Deng, Metadata Librarian University of Central Florida Libraries

What will be covered?

Part I: The Survey and Some Data Basics

The UCF Research Data Management
 Survey: Data Recording and Analysis
 Section Results (Q, D)
 Understanding Data, Research Data and Datasets
 Why data documentation (Q)

Part II: Data Documentation ABC

Data Documentation: Study-level (E)
Data Documentation: Data-level (Structured tabular data, Qualitative data) (E)

Part III: Dataset Metadata

Dataset record examples, their associated standards, and data repositories (E, D)
 Data: DOIs and Data Citation
 Controlled Vocabularies and Thesauri (Q)
 Curation Tools for Datasets

Part IV: Thoughts and Services

A Researcher's View vs. A
 Curator or Librarian's Perspective
 on Data Documentation (D)
 Dataset and Metadata Services
 at UCF

Q: w/ question. E: w/ examples. D: w/ discussion.

- O Data
- O Research data
- Dataset
- Data documentation
- Data types
- Data formats
- Project level
- File level
- Variable level
- O Label
- O Code
- Derived data
- O Data list
- O SPSS
- O SAS
- **o R**
- O Access
- Spreadsheet
- Curation tool
- O Metadata
- Metadata standards
- Metadata schemas
- Controlled vocabularies
- O Thesauri
- Funding agencies
- Research data management
- O DataCite
- O DOI
- **O** Data citation
- Data repository
- O Dataset Metadata Service
- * Word cloud generated using Tagxedo.



Part I: The Survey and Some Data Basics

Data Practices at UCF

OThe UCF Research Data Management (RDM) Survey

 The UCF Research Data Management Survey, November 2013
 Results delivered on Research Computing Day at Institute for Simulation and Training by Dr. Penny Beile on February 11, 2014
 http://www.ist.ucf.edu/hpc/rcd/Beile_datahandout.pdf

OData Recording and Analysis Section Questions and Results

O 17. Provide any technical details about the tools that you use or would like to be able to easily use for your work or research. These can be name or vendor of the software product, technical requirements of the software, special accelerators like graphical processor units (GPU), etc.

Before exposing the RDM survey results, let's see...

• Provide any technical details about the tools that you use or would like to be able to easily use for your work or research.

Olf applicable, how are you recording lab data? Please check all that apply.

- O Lab notebooks in paper
- O Excel (or other) files on computers in the lab
- O Electronic lab notebook (ELN) tool. Please specify which one.

O Do you document or record any metadata for your data or dataset?

- O Yes
- O No
- Olf you record metadata for your dataset, do you use any local, agencyspecific, or national standards or guidelines?
 - O Yes
 - O No
 - O Not sure

UCF RDM Survey Results: Data Recording and Analysis Section

Thirty-nine (39) respondents listed a variety of technical tools used or needed to perform their research.

More popular tools: SAS/SAS Enterprise version (6) MatLab (5) SPSS (5) R-project programs (4) NVivo (3) SigmaPlot (3)

...

software and databases	Processing, backup, and storage network server and cloud space
AMOS	Automated backup internal to UCF
	system (2)
Ansys/Fluent (2)	Black Armor RAID backup system
ArcGIS/GIS ((2)	Cloud storage/backup (Dropbox and
	HIPAA-compliant cloudspace
	specifically mentioned) (4)
AspenTech	DSpace
CST Microwave Studio	Personal drives
Database with graphical viewing	Replication
capabilities, basic statistics, filtering,	
custom output of datasets	
DTreg	STOKES
EndNote	
FACTSAGE	
GPower	Hardware
Gephi	EPSON Workforce Pro GT-550 scanner
Git/GitHub (2)	Tablets
Interactive Data Language	
LimeSurvey	
Lumerical FDTD	
Mathcad (Vensim) (2)	
MatLad (5)	
MS Office (2)	
NVIVO (3)	
REMARK'S OMP software	
Remain 5 OMR Software	
$\Delta S / S \Delta S$ Enterprise version (6)	
SciFinder Scholar	
SigmaPlot (3)	
SPSS (5)	Source:
SOL	<u>nttp://www.ist.ucf.edu/hpc/rcd</u>
Stata (2)	<u>/bene_dutundidout.pd</u>
Video performance analysis software	

UCF RDM Survey Results: Data Recording and Analysis Section

O18. If applicable, how are you recording lab data? Please check all that apply.

The 49 respondents selected multiple answers, with Excel (or other) files on computers in the lab the most popular choice with 48 responses (98%). This was followed by Lab notebooks in paper (n=29, 59%) and Electronic lab notebook tool (n=3, 6%).

Lab notebooks in paper	29	59 %
Excel (or other) files on	48	98 %
computers in the lab		
Electronic lab notebook	3	6%
(ELN) tool. Please specify		
which one.		

 If respondents indicated that they used an Electronic lab notebook they were asked to specify which one. The two ELNs identified were Google Docs and Word with embedded images storing NMR and other equipment data in a digital format.

UCF RDM Survey Results: Data Recording and Analysis Section

- O19. Do you document or record any metadata for your data or dataset?
 - Of the 62 people who responded, 41 (66%) indicated that they do not add metadata to their datasets while 21 (34%) noted that they do. If respondents replied to the affirmative, they were asked about specific standards or guidelines. Those responses are reported in question 20.

Yes	21	34%
No	41	66%
Total	62	100%

UCF RDM Survey Results: Data Recording and Analysis Section

○20. If you record metadata for your dataset, do you use any local, agency-specific, or national standards or guidelines?

• Twenty-one (21) respondents indicated that they assigned metadata to their data or dataset in question 19. Each of the respondents also answered the follow up question as to the type of standard or guideline applied. Of the responses, 15 (71%) do not use any specific standards or guidelines, five (24%) use identified standards, and one (5%) was not sure.

Yes (please specify)	5	24%
No	15	71%
I'm not sure	1	5%
Total	21	

 The five who use standards or guidelines provided the following types: HIPAA/FERPA, FITS standard, program specific, librarians are helping us with this, and all of the above.

Reflections from the UCF RDM Survey Results on Data Recording and Analysis

OAfter all, is data recording and documentation needed or important in your research lifecycle?

OWhat are the various ways to do data recording, documentation or analysis?

•Will you consider any standard for data documentation in your research process (e.g. local, agency-specific, or national standards or guidelines)? Is it necessary? What are these standards and where to find them?

•What are the typical tools out there that can help with data recording and analysis?

Refreshing the Concepts: What are data, research data and datasets?

OData are numerical quantities or other factual attributes derived from observation, experiment or calculation.

- National Research Council, 1992a. "Setting priorities for space research: Opportunities and imperatives."

• Data are facts, numbers, letters, and symbols that describe an object, idea, condition, situation, or other factors. Data in a database may be characterized as predominantly word oriented (e.g., as in a text, bibliography, directory, dictionary), numeric (e.g., properties, statistics, experimental values), image (e.g., fixed or moving video, such as a film of microbes under magnification or time-lapse photography of a flower opening), or sound (e.g., a sound recording of a tornado or a fire)... Data can also be referred to as raw, processed, or verified.

- Committee for a Study on Promoting Access to Scientific and Technical Data for the Public Interest, National Research Council. *A Question of Balance: Private Rights and the Public Interest in Scientific and Technical Databases* (1999). Available at:

http://www.nap.edu/openbook.php?record_id=9692&page=15

Research Data

OIn the context of these Principles and Guidelines [Principles and Guidelines for Access to Research Data from Public Funding], "research data" are defined as factual records (numerical scores, textual records, images and sounds) used as primary sources for scientific research, and that are commonly accepted in the scientific community as necessary to validate research findings.

- Organisation for Economic Co-operation and Development (OECD, 2007). OECD Principles and Guidelines for Access to Research Data from Public Funding. P.13. Available at: <u>http://www.oecd.org/science/sci-tech/38500813.pdf</u>

Research Data (Cont.)

OResearch data is often defined as the information (e.g. data sets, microarray, numerical data, clinical trial information, textual records, images, sound, etc.) generated or used as quantitative evidence in primary biomedical research. This research data is distinguished by the fact that it is accepted by the research community as a means to validate research findings, observations and hypotheses.

- HLWIKI Canada (2011). http://hlwiki.slais.ubc.ca/index.php/Data_curation

 Research data, unlike other types of information, is collected, observed, or created, for purposes of analysis to produce original research results.

- Edinburgh University Data Library Research Data Management Handbook. http://www.docs.is.ed.ac.uk/docs/data-library/EUDL_RDM_Handbook.pdf

Research Data Classification

- Research data can be generated for different purposes and through different processes. In general, it can include the following types of data:
 - **Observational:** data captured in real-time, usually irreplaceable. For example, sensor data, survey data, sample data, neuroimages.
 - **Experimental:** data from lab equipment, often reproducible, but can be expensive. For example, gene sequences, chromatograms, toroid magnetic field data.
 - Simulation: data generated from test models where model and metadata are more important than output data. For example, climate models, economic models.
 - **Derived or compiled:** data is reproducible but expensive. For example, text and data mining, compiled database, 3D models.
 - Reference or canonical: a (static or organic) conglomeration or collection of smaller (peer-reviewed) datasets, most probably published and curated. For example, gene sequence databanks, chemical structures, or spatial data portals.

Dataset

OA logically meaningful collection or grouping of similar or related data, usually assembled as a matter of record or for research, for example, the <u>American FactFinder Data</u> <u>Sets</u> provided online by the U.S. Census Bureau or the <u>National</u> <u>Elevation Dataset</u> available from the U.S. Geological Survey.

- Online dictionary for library and information science (ODLIS). http://www.abc-clio.com/ODLIS/odlis_A.aspx

A research data set constitutes a systematic, partial representation of the subject being investigated.

- Organisation for Economic Co-operation and Development (OECD, 2007). http://www.oecd.org/science/sci-tech/38500813.pdf

Data Documentation

O "Data documentation explains how data were created or digitised, what data mean, what their content and structure are, and any manipulations that may have taken place." - UK Data Archive

 The term 'documentation' encompasses all the information necessary to interpret, understand and use a given dataset or set of documents.
 Cambridge University Library

O "...a minimum requirement for closing the gap between the data producer and the secondary analyst is a high standard of data documentation." (note: the secondary analyst refers to the data user)

Nielsen, Per: How to teach data producers "the noble art" of data documentation. In: Clubb, Jerome M. (Ed.); Scheuch, Erwin K.(Ed.): Historical social research : the use of historical and process-produced data. Stuttgart : Klett-Cotta, 1980 (Historisch-Sozialwissenschaftliche Forschungen : quantitative sozialwissenschaftliche Analysen von historischen und prozeß-produzierten Daten 6). - ISBN 3-12-911060-7, pp. 477-487. URN: <u>http://nbn-resolving.de/urn:nbn:de:0168-ssoar-326298</u>

Data Documentation and Metadata

OWhat is Metadata

- Meta: Greek prefix. Means "after, behind or beyond." Data: Latin word.
 Factual information used for calculating, reasoning or measuring.
- Metadata means something behind or beyond data itself, and it includes data about its content, containers and contextual information.
- A formal definition: Metadata is data about data, data associated with an object, a document, or a dataset for purposes of description, administration, technical functionality and preservation.
- Can be embedded in the data files/documents themselves

OHow is metadata relevant in the research data cycle? For example,

"Over the life course of a survey that results in a data set - from initial conceptualization to data publication and beyond - a huge amount of **metadata** is typically **produced**. These metadata can be **recorded** in DDI format and **re-used** as the data collection, processing, tabulation, and reporting/dissemination take place."

- Arofan Gregory, Open Data Foundation (2011). The Data Documentation Initiative (DDI): An Introduction for National Statistical Institutes. Available at: <u>http://odaf.org/papers/DDI_Intro_forNSIs.pdf</u>

Data Documentation and Metadata (Cont.)

 Documentation and metadata are different things. However, metadata can be taken as a type of documentation.

ODocumentation is meant to be read by humans; some metadata is designed more for machine processing than human readability.

Research data can be documented at various levels: Project level,
 File or database level and Variable or item level.

• To make your data easy to understand and analyze through your research lifecycle and in the long term, it is considered good practice to document your data. Data documentation is part of the data curation process.

Why Data Documentation

OWhy data documentation? (from Nielsen, Per: How to teach data producers "the noble art" of data documentation)

- **Reliability aspect:** in hard sciences, research results are verified by repetition of the experiment; in social sciences, measuring unique phenomena, control of results and conclusions are possible only if data and full documentation are available
- Methodological aspect: "we ask that all methodological considerations and decisions be reported at the time and place they are relevant"
- Economical aspect: it can be "cheaper to clean and document data files for general use before the primary analysis is started," "reports on new issues can be based on existing well-documented files"

• Historical aspect: archive and preserve information for future generations

• Additional aspect: to meet funder requirements

Data Definitions and Specific Requirements from the Funding Agency: National Science Foundation (NSF)

O The term "data" is used in this report to refer to any information that can be stored in digital form, including text, numbers, images, video or movies, audio, software, algorithms, equations, animations, models, simulations, etc. Such data may be generated by various means including observation, computation, or experiment.

-National Science Foundation (2005). Long-Lived digital data Collections: enabling Research and education in the 21st Century. P.9. Available at: http://www.nsf.gov/pubs/2005/nsb0540/nsb0540.pdf

• As stated in NSF's "Information about the Data Management Plan Required for all Proposals" for Biological Sciences, the Federal government defines data (OMB Circular A-110) as: "...the recorded factual material commonly accepted in the scientific community as necessary to validate research findings." This definition includes both original data (observations, measurements etc.) as well as metadata (e.g., experimental protocols, software code for statistical analysis etc.).

Data Definitions and Specific Requirements from the Funding Agency: National Science Foundation (NSF) (Cont.)

- The NSF Grant Proposal Guide recommends the inclusion of a "data management plan" that explains how your proposal will comply with NSF's data sharing policies. The data management plan may include:
 - The types of data, samples, physical collections, software, curriculum materials, and other materials to be produced in the course of the project;
 - The standards to be used for data and metadata format and content (where existing standards are absent or deemed inadequate, this should be documented along with any proposed solutions or remedies);
 - Policies for access and sharing including provisions for appropriate protection of privacy, confidentiality, security, intellectual property, or other rights or requirements;
 - O Policies and provisions for re-use, re-distribution, and the production of derivatives;
 - O Plans for archiving data, samples, and other research products, and for preservation of access to them.
- See <u>NSF's Grant Proposal Guide</u> for more information.
- Search Data Management Plan requirements of different funders at DMPTool (<u>https://dmptool.org/guidance</u>).

Part II: Data Documentation ABC

Data Documentation ABC

• Ensure that all data collected and generated through your research lifecycle is documented.

•At the beginning of your research, check what kind of documentation is available or necessary, and identify needed documentations which will enable data preservation and reuse in the future.

- The various kinds of documentation may include:
 - Embedded documentation (included within the data, e.g., code, field and label descriptions, descriptive headers or summaries, transcripts, in document properties)
 - Supporting documentation (in separate file, e.g., working papers, lab books, questionnaires or interview guides, project reports, publications)
 - Catalog Metadata (for data archiving, identification and locating)

Data Documentation ABC (Cont.)

• The different types of documentations may include:

- OLaboratory notebooks & experimental protocols
- OQuestionnaires, code books with full variable and value labels & data dictionaries
- OInformation about equipment settings & instrument calibration
- OSoftware syntax & output files
- ODatabase schema
- OMethodology reports
- OAssumptions made during analysis
- Provenance information about sources of derived data, different versions of the dataset

Data Documentation ABC (Cont.)

ODuring your research, document all research data formats utilized by your project. Research data comes in many varied formats, such as (by broad categories):

OText - flat text files, Word, PDF, RTF, XML.

- ONumerical Statistical Package for the Social Sciences (SPSS), Stata, Excel.
- OMultimedia jpeg, tiff, dicom, mpeg, quicktime.
- OModels 3D, statistical.
- OSoftware Java, C programs.
- ODiscipline specific Flexible Image Transport System (FITS) in astronomy, Crystallographic Information File (CIF) in chemistry.
- OInstrument specific Olympus Confocal Microscope Data Format, Carl Zeiss Digital Microscopic Image Format (ZVI).

Research Data Types and Formats (in detail)

Type of data	Acceptable formats for sharing, reuse and preservation	Other acceptable formats for data preservation
Quantitative tabular data with extensive metadata a dataset with variable labels, code labels, and defined missing values, in addition to the matrix of data	SPSS portable format (.por) delimited text and command ('setup') file (SPSS, Stata, SAS, etc.) containing metadata information some structured text or mark-up file containing metadata information, e.g. DDI XML file	proprietary formats of statistical packages e.g. SPSS (.sav), Stata (.dta) MS Access (.mdb/.accdb)
Quantitative tabular data with minimal metadata	comma-separated values (CSV) file (.csv) tab-delimited file (.tab) including delimited text of given	delimited text of given character set - only characters not present in the data should be used as delimiters (.txt)
column headings or variable names, but no other metadata or labelling	character set with SQL data definition statements where appropriate	widely-used formats, e.g. MS Excel (.xls/.xlsx), MS Access (.mdb/.accdb), dBase (.dbf) and OpenDocument Spreadsheet (.ods)
	ESRI Shapefile (essentialshp, .shx, .dbf, optionalprj, .sbx, .sbn)	ESRI Geodatabase format (.mdb) MapInfo Interchange Format (.mif) for vector
Geospatial data	geo-referenced TIFF (.tif, .tfw)	data
vector and raster data	CAD data (.dwg)	Keyhole Mark-up Language (KML) (.kml)
	tabular GIS attribute data	Adobe Illustrator (.ai), CAD data (.dxf or .svg)
		binary formats of GIS and CAD packages
	eXtensible Mark-up Language (XML) text	Hypertext Mark-up Language (HTML) (.html)
Qualitative data	according to an appropriate Document Type Definition (DTD) or schema (.xml)	widely-used proprietary formats, e.g. MS Word (.doc/.docx)
textual	Rich Text Format (.rtf)	some proprietary/software-specific formats
	plain text data, ASCII (.txt)	e.g. NUD*IST, NVivo and ATLAS.ti

Research Data Types and Formats (in detail, cont.)

Type of data	Acceptable formats for sharing, reuse and preservation	Other acceptable formats for data preservation
Digital image data	TIFF version 6 uncompressed (.tif)	JPEG (.jpeg, .jpg) but only if created in this format TIFF (other versions) (.tif, .tiff) Adobe Portable Document Format (PDF/A, PDF) (.pdf) standard applicable RAW image format (.raw) Photoshop files (.psd)
Digital audio data	Free Lossless Audio Codec (FLAC) (.flac)	MPEG-1 Audio Layer 3 (.mp3) but only if created in this format Audio Interchange File Format (AIFF) (.aif) Waveform Audio Format (WAV) (.wav)
Digital video data	MPEG-4 (.mp4) motion JPEG 2000 (.mj2)	
Documentation and scripts	Rich Text Format (.rtf) PDF/A or PDF (.pdf) HTML (.htm) OpenDocument Text (.odt)	plain text (.txt) some widely-used proprietary formats, e.g. MS Word (.doc/.docx) or MS Excel (.xls/.xlsx) XML marked-up text (.xml) according to an appropriate DTD or schema, e.g. XHMTL 1.0

Source: http://www.data-archive.ac.uk/create-manage/format/formats-table

Data Documentation ABC (Cont.)

- Keep the wide variety of materials that are generated or collected in your research. Research data (traditional and electronic research) may include all of the following:
 - O Documents (text, Word), spreadsheets
 - Laboratory notebooks, field notebooks, diaries
 - O Questionnaires, transcripts, codebooks
 - O Audiotapes, videotapes
 - O Photographs, films
 - Test responses
 - Slides, artifacts, specimens, samples
 - Collection of digital objects acquired and generated during the process of research
 - O Data files
 - Database contents (video, audio, text, images)
 - Models, algorithms, scripts
 - Contents of an application (input, output, log files for analysis software, simulation software, schemas)
 - Methodologies and workflows
 - Standard operating procedures and protocols

Other research records:

- Correspondence
- Project files
- Grant applications
- Ethics applications
- \circ Technical reports
- Research reports
- \circ Master lists
- Signed consent forms

Source: <u>How to manage research data</u>, Research Support Services, University of Edinburgh Information Services

Data Documentation ABC (Cont.)

ODocument research data at different levels:

- OStudy-level
- **OData-level**
 - OStructured tabular data
 - OQualitative data

OUtilize software to create embedded documentation for the data (if applicable), and make separate supporting documentation (e.g. readme text files) to describe the list of files and documentations in a folder.

O In addition, provide unique identifier for the dataset (e.g. doi, purl, handle...);

• Further, make sure that your data meets citation requirement (if applicable), and discuss with relevant personnel on how data can be archived and shared in a data center, or a library digital repository for others to search, locate and reuse.

Data Documentation: Study-level

OInformation in the Data Documentation Study-level and Data-level section is from UK Data Archive (<u>http://www.data-archive.ac.uk/create-manage/document</u>)

- Study-level information: the research context and design, data collection methods, data preparation and results or findings
 - the context of data collection: project history, aims, objectives and hypotheses
 - data collection methods: data collection protocols, sampling design, instruments used, hardware and software used, data scale and resolution, temporal coverage and geographic coverage, and digitization or transcription methods
 - structure of data files, number of cases, records, variables and relationships between files
 - O data sources used and provenance of materials, e.g. for transcribed or derived data
 - data validation, checking, proofing, cleaning and other quality assurance procedures carried out, such as checking for equipment and transcription errors, calibration procedures, data capture resolution and repetitions, or editing, proofing or quality control of materials
 - modifications made to data over time since their original creation and identification of different versions of datasets
 - for time series or longitudinal surveys, changes made to methodology, variable content, question text, variable labelling, measurements or sampling
 - information on data confidentiality, access and use conditions, where applicable

Data Documentation: Data-level

 Descriptions and annotations at the variable, data item or data file level.

- Onames, labels and descriptions for variables, records and their values
- Oexplanation of codes and classification schemes used
- Ocodes of, and reasons for, missing values
- Oderived data created after collection, with code, algorithm or command file used to create them
- Oweighting and grossing variables created and how they should be used
- Odata list describing cases, individuals or items studied, for example for logging qualitative interviews

Data Documentation: Structured tabular data

OStructured, tabular data should have cases or records and variables adequately documented with:

- •Names, labels and descriptions for all variables, fields, records and their values. Variable labels should:
 - Obe brief with a maximum of 80 characters
 - Oindicate the unit of measurement, where applicable
 - Oreference the question number of a survey or questionnaire, where applicable

How to name the variable to document the survey result for "Q11: hours spent taking physical exercise in a typical week"? For example: q11hexw

Data Documentation: Structured tabular data

• Code labels

How to name the variable for female respondents?

For example: p1sex (with codes '1=female ', '2=male', '-8=don't know', '-9=not answered')

 Coding or classification schemes used, ideally with a bibliographic reference

Where to find a list of codes to classify respondents' jobs?

Reference: Standard Occupational Classification 2000

Where to get the country codes?

Reference: ISO 3166 alpha-2 country codes

O Codes of, and reasons for, missing data

How to document missing data?

For example: '99=not recorded', '98=not provided (no answer)', '97=not applicable', '96=not known', '95=error'

http://ukdataservice.ac.uk/managedata/document/data-level.aspx

Data Documentation: Structured tabular data (cont.)

OData-level descriptions can be embedded within a data file

OStatistical e.g. SPSS

Ovariable descriptions and attributes (codes, data type, missing values) of each variable in the data file can be documented in 'Variable View' or via syntax, whereby embedded data documentation is then contained in the SPSS command file

	Name	Type	WiD	c Label	Values	Missing
1	Q1	Numeric	8 0	Have you previously shared research data	(0, No}	99, 98, 9
2	Q2a	Numeric	8 0	Share data - Send data to researchers upon request	(0, No)	None
3	Q2b	Numeric	8 0	Share data - Make data available through a website or other public resource	(0, No)	None
. 4	Q2c	Numeric	8 0	Share data - Submit data to a data bank, data centre or data archive for pres	(0, No)	None
5	Q2d	Numeric	8 0	Share data - Via publications	(0, No)	None
6	Q2e	Numeric	8 0	Share data - Exchange data with trusted researchers	(0, No}	None
7	Q2g	Numeric	8 0	Share data- made available to colleagues	(0, No)	None
8	Q2f	String	3 0	Other way to share data	None	None
9	Q3a	Numeric	8 0	Reason to share data - Sharing research data is an important part of knowle	(0, No)	None
10	Q3b	Numeric	8 0	Reason to share data - Research data should be used to the maximum by a	(0, No)	None
11	Q3c	Numeric	8 0	Reason to share data - Sharing data and knowledge advances science	(0, No)	None
12	Q3d	Numeric	8 0	Reason to share data - I personally know the researcher I shared data with	(0, No)	None
13	Q3e	Numeric	8 0	Reason to share data - I was required to do so by my institute / employer /	(0, No)	None
14	Q3f	Numeric	8 0	Reason to share data - I was required to do so by the research funder	(0, No)	None
15	Q3a	Numeric	8 0	Reason to share data - Sharing research data may lead to new research coll	(0, No)	None
Data Documentation: Structured tabular data (cont.)

OData-level descriptions can be embedded within a data file

ODatabases e.g. MS Access

Ovariable descriptions and attributes can be documented in 'Design View' and relationships between tables and files can be created

Ble Edit View Ir 12 15 2	isert Iools Window	- Help Adobe PDF	1 3 2102	
Field Name	Data Type	Description	1	
Farmercode	Text	Farmer Identification Code		
Scheme	Text	Drainage scheme name		
Transcript	Hyperlink	Transcript of interview		
Date	Date/Time	Date of interview		
Interviewer	Text	Name of interviewer		
Incomplete data	Yes/No	Observation has data missing		
Nb holdings	Number	Number of holdings farmed by farmer		
Tot Farmsize	Number	Total size business unit (ha)		
% area owned	Number	Area owned of business unit, in percentages (%)		
Farmtype	Text	Main enterprise of business unit		
Organic	Yes/No	Farm is organic		
	Fit	ald Properties		
General Lookup				
Field Size	Decimal	5	The Field	
Format	Percent		description is	
Precision	2		optional. It	
icale	0		helps you	
Decimal Places	0		describe the	
Input Mask			field and is also	
Caption			displayed in the	
Default Value	0		status bar whe	
Validation Rule	120		you select this	
/alidation Text				
Required	Yes	es		
Indexed	No		descriptions	
Smart Tank			descriptions.	

Data Documentation: Structured tabular data (cont.)

Data-level descriptions can be embedded within a data file

OSpreadsheets e.g. MS Excel

Oan additional worksheet within the data file can contain datarelated documentation

100	And the cost and the		2.1		1	
-	1 2 2 2 2 2 4		ply with Change	s End Review.		
1						
and so the	G18 -	fx:				
-	A	В	С	D	E	F
1	Site 🔽	Location	Туре	Instrument Numb	From	
2	Beckingham	Beckingham & Idle Baro	Barometer	73937	7/2/2007	18/10
3	Beckingham	Beckingham Ditch	Diver	80137	7/2/2007	16/1
4	Beckingham	Beckingham Fld Centre	Diver	80136	7/2/2007	16/1.
5	Beckingham	Beckingham Fld Edge	Diver	80129	7/2/2007	16/1
6	Bushley	Bushley Barometer	Barometer	77599	14/2/2007	4/11
7	Bushley	Bushley Ditch	Diver	63017	14/2/2007	23/1
8	Bushley	Bushley Fld Centre	Diver	53632	14/2/2007	23/1
9	Bushley	Bushley Fld Edge	Diver	53194	14/2/2007	12/4
10	Cuddyarch Sough	Cuddyarch Sough Baro	Barometer	62943	10/5/2007	30/1
11	Cuddyarch Sough	Cuddyarch Sough Fld Centre	Barometer	62963	10/5/2007	30/1
12	Cuddyarch Sough	Cuddyarch Sough Fld Edge	Barometer	62959	10/5/2007	30/1
13	Cuddyarch Sough	Wedholme Sough (River)	Diver	48432	10/5/2007	30/1.
14	Idle	Idle Ditch	Diver	80133	7/2/2007	7/1.
15	Idle	Idle Fld Centre	Diver	80131	7/2/2007	16/1
16	Idle	Idle Fld Edge	Diver	80132	7/2/2007	16/1
17	Idle	Idle Upland	Barometer	77531	8/2/2007	18/10
18	Morda	Morda Baro	Barometer	62975	31/5/2007	29/1
19	Morda	Morda Ditch	Barometer	62970	31/5/2007	29/1.

Data Documentation: Structured tabular data (cont.)

Data-level descriptions can be embedded within a data file
 OGIS e.g ArcGIS

Oshapefiles (layers) and tables can be organised in a geo-database with rich metadata created in ArcCatalog

OA dataset may also be accompanied with a Codebook detailing all variables and their values

•Variable naming

OFull variable name

Omeaningful abbreviations (e.g. oz%=percentage ozone; moocc=mother occupation)

Oquestion number system (Q1a, Q1b, Q2, Q3a,...)

Onumerical order system (V1, V2, V3,...)

Source:

http://ukdataservice.ac.uk/managedata/document/data-level.aspx

Structured Metadata: XML Schemas

- XML schema brings documentation into a single document; creates structured content about the data and allows data interoperability and sharing.
- Olt can document comprehensive variable level information, such as basic data dictionary, question text and question routing instructions.
- Data Documentation Initiative (DDI): a metadata specification for the social and behavioral sciences. It is an XML metadata standard for documenting numeric data. Detailed information is available at: <u>http://www.ddialliance.org/</u>

O Projects using the DDI (http://www.ddialliance.org/ddi-at-work/projects)

ODDI-compliant data repository:

- O ICPSR Inter-university Consortium for Political and Social Research
 - Data deposit form: <u>https://www.icpsr.umich.edu/cgi-bin/ddf2</u>
 - UCF is a member of ICPSR.

O UKDA - UK Data Archive

Social Science Dataset Display





The survey instrument used for this study was developed from an initial focus grouop involving six homeless women conducted in November 2002 and was finalized in April 2003. Field Labels: Study description Citation Funding Scope of study

- Subject terms
- Smallest
 geographic unit
- Geographic coverage
- Time period
- Date of collection
- Unit of observation
- Universe
- Data types
- Data collection notes

Methodology

- Study purpose
- Study design

To select women for Part 1 (Female Interviews) of the study, researchers entered into a cooperative agreement with a large, general-purpose shelter for the homeless in each of Jacksonville, Miami, Orlando and Tampa. All of the shelters where respondents were solicited were general-purpose homeless facilities, not battered-women's facilities, and not special-purpose facilities devoted exclusively to teens, to the addicted, or to the mentally ill.

Researchers attempted to interview the first 200 women who came "through the door" of the participating facilities during the data collection period. Recognizing the logistical difficulties of implementing any specific sampling plan in a social service context often characterized by crisis and relative chaos, researchers allowed for some deviation from th desideratum. Efforts were made to interview every woman who sought services at the respective facilities until the q of 200 interviews per site was reached.

Realizing also that interviewing each woman that came through the door would not always be possible, researchers guidelines for interviewers to randomly select from multiple women.

For Part 2 (Male Interviews), men were selected in similar fashion as females but only at the Orlando facility and only the quota of 100 was met.

(2 datasets; 57,930 KB) You can find more information via the sample characteristics utility:

Sample:

Table of Contents

Quick Download

Top of page Access Notes Dataset(s) Study Description Chation Funding Scope of Study Methodology Version(s) Related Publications Variables Utilities Metadata Exports

Download Statistics

Mode of Data Collection: Nace-to-face interview

Description of Variables

For Part 1 (Female Interviews), the data include information related to the respondent's living conditions in the past month, lifetime experience with homelessness, the respondent's partner, and living conditions as a child. There are a variables related to childhood experience with violence, experience with forced sexual situations, adult experience wi violence, and basic demographic information. As well, there is information covering such areas as experience with stalking, self-image, use of alcohol and drugs, current health, and financial status. There is also a self-report of crimin history, information related to how the respondent spent her days and evenings, and the physical environment surrounding the respondent during the day and evening. Finally, there are a small number of questions answered by interviewer regarding the respondent and the interview itself.

For Part 2 (Male Interviews), the data include much of the same information as was collected in Part 1. Information from Part 1 not included in Part 2 primarily includes questions pertaining to experience with forced sexual situations and questions related to pregnancy and children.

Response Rates: Not applicable.

Presence of Common Scales: The Conflict Tactics Scale (CTS) The Personal History Form (PHF) The Addiction Severity Index (ASI)

Extent of Processing: ICPSR data undergo a confidentiality review and are altered when necessary to limit the risk of discource_ICPBR also routinely creates ready-to-go data files along with setups in the major statistical software formats as well as standard codebooks to accompany the data. In addition to these procedures, ICPSR performed the following processing steps for this data collection:

- Created variable labels and/or value labels.
- Standardized missing values.
- Checked for undocumented or out-of-range codes.

Field Labels:

- Sample
- Mode of data collection
- Description of variables
- Response rates
- Presence of common scales
- Extent of processing



Privacy Policy



Privacy Policy

Social Science Dataset in DDI Format

🕘 www.icpsr.umich.edu/icpsrweb/ICPSR/ddi2/studies/20363

This XML file does not appear to have any style information associated with it. The document tree is shown below.

- <codebook id="ICPSR20363" version="1.2.2"></codebook>		cons
- <docdscr></docdscr>		bibli
- <citation></citation>		וטוט
+ <titlstmt></titlstmt>		info
- <prodstmt></prodstmt>		
- <producer abbr="ICPSR"></producer>		desc
<extlink <="" td="" title="ICPSR Logo" uri="http://www.icpsr.umich.edu/images/icpsr-logo.gif"><td>role="image"/></td><td>-וחם</td></extlink>	role="image"/>	-וחם
Inter-university Consortium for Political and Social Research		
<extlink title="URL of ICPSR Web Site" uri="http://www.icpsr.umich.edu/ICPSR/"></extlink>		docu
+ <copyright></copyright>		itsei
		who
+ <verstmt></verstmt>		*****
<holdings uri="http://www.icpsr.umich.edu/icpsrweb/ICPSR/ddi2/studies/20363"></holdings>		
	Included F	ields.
	metaded	ierus.
+ <stdydscr< del="">></stdydscr<>	citation	
- <filedscr id="F1"></filedscr>		C
- <filetxt id="Part1"></filetxt>	• title	eStmt
<filename>Female Interviews</filename>		
	• pro	astmt
		Stmt
+ <filedscr id="F2"></filedscr>	• ver	SUIIL
	a hal	dinge

docDscr The **Document** Description consists of bibliographic information describing the **DDI-compliant** document itself as a whole.

http://www.icpsr.umich.edu/icpsrweb/ICPSR/ddi2/studies/20363

tleStmt odStmt

- erStmt
- holdings

Social Science Dataset in DDI Format (Cont.)

<stdyDscr> <citation> **Included Fields:** -<titlStmt> -<titl> Experience of Violence in the Lives of Homeless Persons: The Florida Four City Study, 2003-2004 </titl> <IDNo agency="ICPSR">20363</IDNo> <IDNo agency="CrossRef">10.3886/ICPSR20363.v1</IDNo> </titlStmt> -<rspStmt> <AuthEnty affiliation="University of Central Florida">Wright, James D.</AuthEnty> <AuthEnty affiliation="University of Central Florida">Jasinski, Jana L.</AuthEnty> <AuthEnty affiliation="University of Central Florida">Mustaine, Elizabeth</AuthEnty> <AuthEnty affiliation="University of North Florida">Wesely, Jennifer</AuthEnty> </rspStmt> -prodStmt> +<fundAg></fundAg> <grantNo agency="United States Department of Justice, Office of Justice Programs, National Institute of Justice">2002 </prodStmt> + <distStmt></distStmt> - <biblCit> Wright, James D., Jana L. Jasinski, Elizabeth Mustaine, and Jennifer Wesely. Experience of Violence in the Lives of Hon ICPSR20363-v1. Ann Arbor, MI: Inter-university Consortium for Political and Social Research [distributor], 2010-11-22 </biblCit> <holdings URI="http://dx.doi.org/10.3886/ICPSR20363.v1"/> </citation> -<stdvInfo> +<subject></subject> + <abstract></abstract> +<abstract>/abstract> + <abstract></abstract> + <abstract></abstract> +<sumDscr></sumDscr> </stdvInfo> + <method></method> -<dataAccs> + <setAvail media="online"></setAvail> + <useStmt></useStmt> </dataAccs>

/stdvDscr>

Citation titlStmt rspStmt prodStmt fundAg grantNo distStmt biblCit Holdings stdyInfo Subject Abstract sumDscr Method dataColl Notes anlyInfo dataAccs setAvail useStmt

stdyDscr The Study **Description** consists of information about the data collection, study, or compilation that the **DDI-compliant** documentation file describes. This section includes information about how the study should be cited, who collected or compiled the data, who distributes the data, keywords about the content of the data, summary (abstract) of the content of the data, data collection methods and processing, etc.

Social Science Dataset in DDI Format (Cont.)

Included Fields: fileDscr fileTxt fileName <u>fileDscr</u>

Data Files Description

Information about the data file(s) that comprises a collection. This section can be repeated for collections with multiple files.

Data Documentation: Qualitative Data

•Context and participant details of interviews can be:

- OA descriptive header or summary page in transcripts or field notes.
- OA structured data list
- OXML mark-up of data, for example,
 - Text Encoding Initiative (TEI) to mark up interview transcript
 - OQualitative Data Exchange Format (QuDEx) for researcher annotations and data linking

Data Documentation: Qualitative Data (Cont.)

 Anonymisation of textual data (e.g., replacing real names of people, organizations and locations with pseudonyms)

○ File naming

- Meaningful, short names; identify file types (e.g. interviews, focus groups, field notes, audio recordings); avoid space, special characters; avoid long names
- Organizing files in folders: Create uniform and structured folder names based on cases, studies, locations, data types etc., or the original, anonymized, coded or annotated versions of data

• Version control: Version numbering in file names

 Documentation: Methodology description, project plan, interview guidelines, consent form templates; data analyses and manipulation

Data Documentation: Qualitative Data Example

Morbidity and mortality - Population, vital statistics and censuses Child development and child rearing - Social stratification and groupings Elderly - Social stratification and groupings Family life and marriage - Social stratification and groupings Gender roles - Social stratification and groupings Use and provision of specific social services - Social welfare policy and systems

Depositor(s): Blaxter, M., University of East Anglia

Principal Investigator(s): Blaxter, M., University of East Anglia

Data Collector(s): Blaxter, M., University of East Anglia

Sponsor(s): Economic and Social Research Council

Abstract:

Morbidity and mortality - Population, vital statistics and censuses Child development and child rearing - Social stratification and groupings Elderly - Social stratification and groupings Family life and marriage - Social stratification and groupings Gender roles - Social stratification and groupings Use and provision of specific social services - Social welfare policy and systems

Depositor(s): Blaxter, M., University of East Anglia

Principal Investigator(s): Blaxter, M., University of East Anglia

Data Collector(s): Blaxter, M., University of East Anglia

Sponsor(s): Economic and Social Research Council

Abstract: This is an enhanced qualitative study.

The research looked at beliefs and attitudes to **health** and medical care, inter-genera and social history of members of a grandmother generation. The original study include daughters as well; this collection contains only the grandmother interviews.

Grandmothers are asked extensive questions about their own **health** and the **heal** members. Details are provided on episodes of illness and remedies used, both home as Specific topics of accidents, nutrition, dental care, and immunisation are covered.

More generally, grandmothers are asked about their views of their personal docto **health** services. They give opinions on the quality of **health** care before and after the National **Health** Service.

Grandmother-daughter relationships are explored, especially around the subject of off media medical advice concerning care for the grandchildren.

The research looked at beliefs and attitudes to **health** and medical care, inter-generational relationships, and social history of members of a grandmother generation. The original study included interviews with daughters as well; this collection contains only the grandmother interviews.

Grandmothers are asked extensive questions about their own **health** and the **health** of other family members. Details are provided on episodes of illness and remedies used, both home and **health** services. Specific topics of accidents, nutrition, dental care, and immunisation are covered.

More generally, grandmothers are asked about their views of their personal doctors and institutional **health** services. They give opinions on the quality of **health** care before and after the introduction of the National **Health** Service.

Grandmother-daughter relationships are explored, especially around the subject of offering and taking of ^f off medical advice concerning care for the grandchildren.

The collection has been enhanced by: conversion from paper to searchable RTF format by OCR and The collection has been enhanced by: conversion from paper to searchable RTF fo extensive editing and formatting of all interview transcripts. This collection will also be made available extensive editing and formatting of all interview transcripts. This collection will also through ESDS Qualidata Online.

 Example is from: A NESSTAR FOR QUALITATIVE DATA: BUILDING BLOCKS FOR DIGITAL FUTURES. By Corti Louise et al. available at: <u>http://data-archive.ac.uk/media/376907/digitalfutures_dashish_21nov2012.pdf</u>

Data Documentation: Qualitative Data Example

OData List

Study Number 6124 Being a Doctor: a Sociological Analysis, 2005-2006 Nettleton, S

(Interview			Date of	No of (Text File
		Gender	Description	Interview	Pages	Name
(x001	Male	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	07/09/2005	36	6124int001
	x002 /	Male	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	09/09/2005	41	6124int002
	x003	Male	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	09/09/2005	39	6124int003
	x004	Female	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	13/09/2005	36	6124int004
	x005	Male	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	13/09/2005	34	6124int005
	x006	Male	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	14/09/2005	50	6124int006
	x007	Male	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	21/09/2005	31	6124int007
	x008	Male	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	21/09/2005	35	
	x009	Male	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	22/09/2005	33	e lext File Name:
	x010	Female	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	22/09/2005	23	
	x011	Male	Interview with a Hospital Doctor in a Multi-Ethnic Northern City	22/09/2005	36	61 241nt 001
	:		w with a Hospital Doctor in a Multi-Ethnic Northern City	27/09/2005	41	(124:++002
Interv	iew IL):	w with a Hospital Doctor in a Multi-Ethnic Northern City	27/09/2005	21	61241Nt002
×001			w with a Hospital Doctor in a Multi-Ethnic Northern City	30/09/2005	20	(
			w with a Hospital Doctor in a Multi-Ethnic Northern City	05/10/2005	19	<u>t</u>
×002			w with a Hospital Doctor in a Multi-Ethnic Northern City	05/10/2005	27	
XUUZ			w with a Hospital Doctor in a Multi-Ethnic Northern City	07/10/2005	27	
			w with a Hospital Doctor in a Multi-Ethnic Northern City	17/10/2005	11	
•••			w with a Hospital Doctor in a Multi-Ethnic Northern City	19/10/2005	33	6124int019
			w with a Hospital Doctor in a Multi-Ethnic Northern City	07/11/2005	21	6124Int020
			W with Hospital Doctor in Northern Tourist City	19/07/2005	50	6124Int021
	-002	Mala	Interniew with Hearitel Dester in Northern Tourist City	10/08/2005	40	0124INU22
	2003	Male	Interview with Heapital Doctor in Northern Tourist City	17/08/2005	20	0124Int023
	2004	Tomolo	Interview with Hospital Doctor in Northern Tourist City	07/11/2005	21	0124INU24 6124int025
	2005	Molo	Interview with Hospital Doctor in Northern Tourist City	14/11/2005	ວ∠ วว	0124IIII023 6124int026
	2000	Male	Interview with Hospital Doctor in Northern Tourist City	16/11/2005	∠3 22	6124int020
	2007	Fomale	Interview with Hospital Doctor in Northern Tourist City	17/11/2005	23 18	6124int027
	2000	Malo	Interview with Hospital Doctor in Northorn Tourist City	18/11/2005	20	6124int020
	2003	Male	interview with hospital Doctor in Northern Tourist City	10/11/2003	20	0124111023

Part III: Dataset Metadata

Dataset Metadata

- •Create and generate metadata for your research data and datasets in your research lifecycle to preserve the data in the long run.
 - OConsider what information is needed for the data to be read and interpreted in the future.

 OUnderstand your funder requirements for data documentation and metadata. Funder requirements for NSF, GBMF, IMLS, NEH, NIH and NOAA can be found at <u>https://dmptool.org/guidance</u>.

OConsult available metadata standards in your field. You may refer to <u>Common Metadata Standards</u> and <u>Domain Specific</u> <u>Metadata Standards</u> for details.

Dataset Metadata (cont.)

 Describe data and datasets created in your research lifecycle, and use software programs and tools to assist in data documentation.
 Assign or capture administrative, descriptive, technical, structural and preservation metadata for the data. Some potential information to document:

ODescriptive metadata

Name of creator of data set
Name of author of document
Title of document
File name
Location of file
Size of file

OStructural metadata

• File relationships (e.g. child, parent)

Dataset Metadata (cont.)

OTechnical metadata

- O Format (e.g. text, SPSS, Stata, Excel, tiff, mpeg, 3D, Java, FITS, CIF)
- Compression or encoding algorithms
- O Encryption and decryption keys
- Software (including release number) used to create or update the data
- O Hardware on which the data were created
- Operating systems in which the data were created
- O Application software in which the data were created

OAdministrative metadata

- O Information about data creation (e.g. date)
- Information about subsequent updates, transformation, versioning, summarization
- Descriptions of migration and replication
- O Information about other events that have affected the files

Dataset Metadata (cont.)

OPreservation metadata

- OFile format (e.g. .txt, .pdf, .doc, .rtf, .xls, .xml, .spv, .jpg, .fits)
- **O** Significant properties
- Technical environment
- **O** Fixity information
- Adopt a **thesauri** in your field if applicable or compile a data dictionary for your dataset.
- Obtain **persistent identifiers** (e.g. doi, purl) for datasets if possible to ensure data can be found in the future.
- For your full data management plan, visit <u>UCF Libraries Data Management</u> <u>Guide.</u> Also, refer to Digital Curation Centre's Checklist for a Data Management Plan (<u>http://www.dcc.ac.uk/sites/default/files/documents/resource/DMP_Checklist_2013.pdf</u>).

Dataset Record Examples, Their Associated Standards, and Data Repositories

Common Metadata StandardsDisciplinary Metadata Standards

OActivity: Choose a dataset or a standard in your field to examine and critique OSocial Science Dataset OHumanities Dataset OBiological Sciences Dataset OBiotechnology Dataset OGeospatial Dataset **OEarth Science Dataset OPhysical Science Dataset** O0ther...

Common Metadata Standards

- Dublin Core (DC): A general metadata standard for describing a wide range of digital resources.
 - Dublin Core Metadata Element Set, Version 1.1 (<u>http://dublincore.org/documents/dces/</u>)
 - 15 Elements: Title, Creator, Subject or keyword, Description, Publisher, Type, Format, Identifier, Source, Language, Relation, Coverage, Rights
 - O DCMI Metadata Terms (<u>http://dublincore.org/documents/dcmi-terms/</u>)
 - O DC Qualifiers (<u>http://dublincore.org/documents/usageguide/qualifiers.shtml</u>)

O Encoded Archival Description (EAD)

• A standard for encoding archival finding aids with XML.

O Government Information Locator Service (GILS)

• The Global Information Locator Service defines a core element set for government information so that it can be more searchable and discoverable by the general public.

O ONIX for Books (ONline Information eXchange)

• An international standard for representing and communicating book industry product information in XML format.

Common Metadata Standards (Cont.): Image and Multi-media Metadata Standards

Categories for the Description of Works of Art (CDWA)

A conceptual framework and guidelines for the description of art objects and images.

Visual Resources Association Core Categories (VRA Core)

A data standard for the description of works of visual culture as well as the images that document them.

PBCore

The metadata standard for audiovisual media developed by the public broadcasting community.

NISO Metadata for

Digital Images

This technical metadata standard defines a set of metadata elements for raster digital images to enable users to develop, exchange, and interpret digital image files. The dictionary has been designed to facilitate interoperability between systems, services, and software as well as to support the longterm management of and continuing access to digital image collections.

Technical Metadata for Multimedia: MPEG-7

The Multimedia Content Description Interface MPEG-7 is an ISO/IEC standard and specifies a set of descriptors to describe various types of **multimedia** information and is developed by the Moving Picture Experts Group.

Disciplinary Metadata Standards: Social Sciences & Humanities

ODDI - Data Documentation Initiative

•A metadata specification for the social and behavioral sciences. Expressed in XML, the DDI metadata specification supports the entire research data life cycle.

O<u>Text Encoding Initiative (TEI)</u>: A standard for the representation of texts in digital form, chiefly in the humanities, social sciences and linguistics.

OHumanities repositories and Projects:

OProjects Using the TEI (from the official TEI website)

• See Appendix 1 for a TEI project example

Disciplinary Metadata Standards: Biological Sciences

ABCD - Access to Biological Collection Data A standard for the access to

and exchange of data about specimens and observations (a.k.a. primary biodiversity data).

Darwin Core

A metadata specification for information about the geographic occurrence of species and the existence of specimens in collections.

EML: Ecological Metadata

<u>Language</u>

A metadata specification developed by the ecology discipline and for the ecology discipline. EML is implemented as a series of XML document types that can be used in a modular and extensible manner to document ecological data.

Disciplinary Metadata Standards: Health Sciences

Health Level 7 Standards

HL7 and its members provide a framework (and related standards) for the exchange, integration, sharing, and retrieval of electronic health information. HL7 standards support clinical practice and the management, delivery, and evaluation of health services.

The Cross-Enterprise Document Sharing (XDS) Metadata

The Healthcare Enterprise (IHE) XDS profile is a protocol for sharing clinical documents in health information exchanges. IHE IT Infrastructure Technical Framework volumes can be accessed at: http://ihe.net/Resources/Technical Frameworks/

National Institute of Health (NIH) Common Data Elements (CDEs)

CDE is a data element that is common to multiple data sets across different studies. NIH encourages the use of CDEs in clinical research, patient registries, and other human subject research in order to improve data quality and opportunities for comparison and combination of data from multiple studies and with electronic health records.

<u>ClinicalTrials.gov Protocol Data</u> <u>Element Definitions</u>

It describes the registration data items (required and optional) that are entered via the Protocol Registration and Results System (PRS).

Disciplinary Metadata: Biological Sciences and Health Sciences - Repositories



Dryad (<u>https://datadryad.org/</u>) A digital repository for data underlying the international scientific publications, with an initial focus on evolutionary biology and related fields.



NIH Data Sharing Repositories page lists NIH-supported data repositories that make data accessible for reuse. Most accept submissions of appropriate data from NIHfunded investigators (and others).



GenBank is the NIH genetic sequence database, an annotated collection of all publicly available DNA sequences.



<u>GBIF - Global Biodiversity</u>

Information Facility

GBIF is a free and open access global web portal promoting and facilitating the mobilization, access, discovery and use of biodiversity data.

ClinicalTrials.gov

<u>ClinicalTrials.gov</u> is a registry and results database of publicly and privately supported clinical studies of human participants conducted around the world.

Examples:

Biological Science Dataset: See Appendix 2 Biotechnology Dataset @ GenBank:

http://www.ncbi.nlm.nih.gov/nucleotide?cmd=Retrieve&dopt=GenBank&list_uids=1293613 Biotechnology Dataset @ PubChem: http://pubchem.ncbi.nlm.nih.gov/summary/summary.cgi?cid=5760 Clinical Study Dataset @ ClinicalTrials: https://clinicaltrials.gov/show/NCT01196442

Disciplinary Metadata Standards: Earth Science

AgMES

Agricultural Metadata Element Set

AgMES is designed to include agriculture specific extensions for terms and refinements from established metadata standard such as Dublin Core and AGLS to facilitate resource discovery, interoperability and data exchange in the agriculture domain.

FGDC/CSDGM

Federal Geographic Data Committee Content Standard for Digital Geospatial Metadata Content standard for digital

geospatial metadata maintained by the Federal Geographic Data Committee (FGDC). Often referred to as the "FGDC Metadata Standard."

ISO 19115:2003

An internationally-adopted schema for describing geographic information and services. It provides information about the identification, the extent, the quality, the spatial and temporal schema, spatial reference, and distribution of digital geographic data.

DIF

Directory Interchange Format An early metadata initiative from the Earth sciences community, intended for the description of scientific data sets. It includes elements focusing on instruments that capture data, temporal and spatial characteristics of the data, and projects with which the dataset is associated.

(Climate and Forecast) Metadata Conventions

A standard for climate and forecast "use metadata" that aims both to distinguish quantities (such as physical description, units, or prior processing) and to locate the data in space-time.

Disciplinary Metadata: Earth Science- Repositories and Data Centers



AGRIS - International System for Agricultural Science and Technology

A global public domain database using the <u>AgMES</u> standard to describe structured bibliographical records on <u>agricultural</u> science and technology. <u>NCDC - National</u> Climatic Data Center

CUMATIC D.

The world's largest climate data archive, providing climatological services and data worldwide. It currently promotes the FGDC/CSDGM metadata standard for its datasets.



CEOS International Directory Network

An international effort to assist users in locating Earth science data sets, data services, and visualizations using DIF metadata. It provides free, online access to metadata on scientific data in the Earth sciences: geoscience, hydrospheric, biospheric, satellite remote sensing, and atmospheric sciences.

See a Geospatial Dataset (appendix 3) and an Earth Science Dataset (appendix 4).

Disciplinary Metadata: Physical Science Metadata Standards and Repositories

OCIF - Crystallographic Information Framework

 An extensible standard file format and set of protocols for the exchange of crystallographic and related structured data.

American Mineralogist Crystal Structure Database

<u>American</u> <u>Mineralogist Crystal</u> <u>Structure Database</u>

A CIF crystal structure database that includes every structure published in the American Mineralogist, The Canadian Mineralogist, European Journal of Mineralogy and Physics and Chemistry of Minerals, as well as selected datasets from other journals.



Physical Science Dataset Example: http://rruff.geo.arizona.edu/AMS/minerals/Abernathyite

Dublin Core Metadata	Standard DIF
Title	Entry_Title
Creator	Data_Set_Citation: Dataset_Creator
	Personnel: Role: Investigator: Last_Name
	Personnel: Role: Investigator: First_Name
	Personnel: Role: Investigator: Middle_Name
Subject and Keywords	Keyword
	Parameters: Category
	Parameters: Topic
	Parameters: Term
	Parameters: Variable
	Parameters: Detailed_Variable
	Source_Name
	Sensor_Name
	Project
	Location
Description	Summary
Publisher	Data_Set_Citation: Dataset_Publisher
	Data_Center: Data_Center_Name
	Data_Center: Data_Center_URL
	Data_Center: Data Center Contact, Last_Name
	Data_Center: Data Center Contact, First_Name
	Data_Center: Data Center Contact, Middle_Name
Contributor	Personnel: Role:
	Personnel: Last_Name
	Personnel: First_Name
	Personnel: Middle_Name
Date	Data_Set_Citation: Dataset_Release_Date
Resource Type	Data_Set_Citation: Data_Presentation_Form

Generic Dataset and Disciplinary Dataset: Metadata Mapping

- The purpose of metadata mapping is for greater data sharing, exchange, interoperability, usability and reusability.
- Example: Mapping of Dublin Core to DIF

Format	Group: Distribution
	Distribution_Media
	Distribution_Size
	Distribution_Format
	Fees
Resource Identifier	Data Center: Data_Set_ID
	Data_Set_Citation: Online_Resource
	Related_URL: URL_Content_Type
	Related_URL: URL
Source	Related_URL: URL_Content_Type
	Related_URL: URL
	Source_Name
Language	Data_Set_Language
Relation	Parent_DIF
	Data_Set_Citation: Online_Resource
	Related_URL: URL_Content_Type
	Related_URL: URL
	Reference
Coverage	Location
	Spatial_Coverage: Southernmost_Latitude
	Spatial_Coverage: Northernmost_Latitude
	Spatial_Coverage: Easternmost_Longitude
	Spatial_Coverage: Westernmost_Longitude
	Temporal_Coverage: Start_Date
	Temporal_Coverage: Stop_Date
	Paleo_Temporal_Coverage:
	Paleo_Start_Date
	Paleo_Temporal_Coverage:
	Paleo_Stop_Date
	Paleo_Temporal_Coverage:
	Chronostratigraphic_Unit
Rights Management	Use_Constraints
	Access_Constraints

Generic Dataset and Disciplinary Dataset: Metadata Mapping

 Compared to generic Dublin Core, discipline metadata standards provide more granularity.

Metadata Standards and Datasets

OCommon Metadata Standards
 (http://guides.ucf.edu/metadata/genMetaStandards)

ODisciplinary Metadata Standards
 (http://guides.ucf.edu/metadata/domMetaStandards)

OQuestions on metadata standards:

- Do they make sense to you?
- Are the standards adequate in your field? Can data be well documented?
- Have you used any standard, or, will you consider it in your future study and research?

Disciplinary Metadata and Data Repositories

For more information on disciplinary metadata standards, tools, and use cases, please refer to UK Digital Curation Centre (DCC)'s <u>Disciplinary Metadata</u> page.

For more information on data repositories and digital repositories, please refer to Databib, OpenDOAR and OAD.

Databib

DataBib: Databib is a community-driven, annotated bibliography of research data repositories. Databib is now merged with re3data.org (http://www.re3data.org/).

*Open*DOAR

OpenDOAR: An authoritative worldwide directory of academic open access repositories. http://www.opendoar.org/countrylist.php



Open Access Directory: Data Repositories A list of repositories and databases for open data. It is part of the Open Access Directory maintained by Simmons College.

http://oad.simmons.edu/oadwiki/Data_ repositories

Dataset Identifiers

Digital Object Identifier (DOI)
 oe.g. <u>http://dx.doi.org/10.3886/ICPSR20363.v1</u>

OHandles

oe.g. http://soar.wichita.edu/handle/10057/3031

OPersistent URLs (PURLs)OAll can be resolved to an internet location.
DOI

• **Digital Object Identifier (DOI):** an identifier scheme administered by the International DOI Foundation. It is built on the Handle System.

OExample:

Dataset: Experience of Violence in the Lives of Homeless Persons: The Florida Four City Study, 2003-2004 (ICPSR 20363)

http://dx.doi.org/10.3886/ICPSR20363.v1

http://dv.doi.org/	10 2006/	ICPSR20363
nup.//ax.doi.org/	10.3000/	.v1
resolver service	prefix (assigning body)	suffix (resource)

Dataset Registration and DataCite

ODataCite: A global citations framework for data with member institutions offering services and advice to researchers.

 Individuals wishing to register a DOI for their dataset normally do so via their data repository, rather than directly through DataCite.

OAny repository wishing to register DOIs needs to obtain a username and password from DataCite to gain access to the registration service.

 Alternatively, the organization can manage its DOIs through a third-party service such as EZID.

DataCite Metadata

 ICPSR (Interuniversity Consortium for Political and Social Research): an associate member of DataCite.

OICPSR's "How to prepare citation":

• Citation required basic elements:

- O Identifier
- O Creator
- O Title
- O Publisher
- O Publication Year

• For example:

- Wright, James D., Jana L. Jasinski, Elizabeth Mustaine, and Jennifer Wesely. Experience of Violence in the Lives of Homeless Persons: The Florida Four City Study, 2003-2004. ICPSR20363-v1. Ann Arbor, MI: Inter-university Consortium for Political and Social Research [distributor], 2010-11-22. doi:10.3886/ICPSR20363.v1
- Persistent URL: <u>http://dx.doi.org/10.3886/ICPSR20363.v1</u>

○ Can be exported as <u>RIS</u> (generic format for RefWorks, EndNote, etc.) or

EndNote XML (EndNote X4.0.1 or higher)

DataCite Metadata Schema

ODataCite Metadata Schema 3.1 (released 2014-10)

(http://schema.datacite.org/meta/kernel-3/doc/DataCite-MetadataKernel_v3.1.pdf)

Table 1: DataCite Mandatory Properties

ID	Property	Obligatio	n
1	Identifier (with type sub-property)	М	
2	Creator (with name identifier sub-properties)	м	Γ
3	Title (with optional type sub-properties)	м	ŀ
4	Publisher	м	ŀ
5	PublicationYear	М	-

ID	Property	Obligation
6	Subject (with scheme sub-property)	R
7	Contributor (with type and name identifier sub-properties)	R
8	Date (with type sub-property)	R
9	Language	0
10	ResourceType (with general type description sub-property)	R
11	AlternateIdentifier (with type sub-property)	0
12	RelatedIdentifier (with type and relation type sub-properties)	R
13	Size	0
14	Format	0
15	Version	0
16	Rights	0
17	Description (with type sub-property)	R
18	GeoLocation (with point and box sub-properties)	R

Table 2: DataCite Recommended and Optional Properties

DataCite XML File

- <resource xsi:schemalocation="http://datacite.org/schema/kernel-2.2 http://schema.datacite.org/schema/kernel-2.2 http://schema/kernel-2.2 http://s</th><th>rg/meta/kernel-2.2/metadata.xsd"></resource>	
<identifier identifiertype="DOI">10.3886/ICPSR20363.v1</identifier>	
+ <creators></creators>	
- <titles></titles>	FIELDS:
- <title></title>	resource
Experience of Violence in the Lives of Homeless Persons: The Florida Four City Study,	resource
	creator
	title
- <pre>publisher> Inter university Consortium for Political and Social Research </pre>	publisher
	publicationYear
<pre><pre>cypublicationYear>2010</pre></pre>	subject
+ <subjects></subjects>	data
+ <dates></dates>	uale
<resourcetype resourcetypegeneral="Dataset"> survey data </resourcetype>	resourceType
- <alternateidentifiers></alternateidentifiers>	alternativeldentifier
<alternateidentifier alternateidentifiertype="ICPSR Study Number">20363<td>version</td></alternateidentifier>	version
<version>1</version>	description
+ <descriptions></descriptions>	•••

http://www.icpsr.umich.edu/icpsrweb/ICPSR/datacite/studies/20363

Controlled Vocabulary

 Controlled vocabulary is a standardized set of terms used to organize knowledge for subsequent retrieval. It can facilitate search and browsing. It can be universally agreed on or locally created.

• What to consider in applying or designing a thesauri for your project:

- Scope of the material (core and surrounding topics, your purpose, existing thesauri and your resource);
- Your project needs and intended audience;
- Funder requirements and institutional expectation;
- What types of controlled vocabularies you may need: subject, genre, physical format, personal names, organization names, events...
- O When choosing particular terms over others, consider three warrants: literary warrant (discipline and field literature), user warrant and organizational warrant. (Gazan, CONTROLLED VOCABULARY & THESAURUS DESIGN, http://www.loc.gov/catworkshop/courses/thesaurus/pdf/cont-vocab-thes-trnee-manual.pdf)

Controlled Vocabulary (Cont.)

•For traditional library catalog:

- OMARC Code List for Countries <u>http://www.loc.gov/marc/countries/</u>
- OMARC Code List for Languages <u>http://www.loc.gov/marc/languages/</u>
- OMARC Source Codes for Vocabularies, Rules, and Schemes
 http://www.loc.gov/marc/sourcecode/form/formsource.html

•For digital and online resources:

- O Internet Media Types <u>www.iana.org/assignments/media-</u> <u>types/index.html</u>
- OMODS Note Types <u>http://www.loc.gov/standards/mods/mods-notes.html</u>
- DCMI Type Vocabulary <u>http://dublincore.org/documents/dcmi-</u> terms/index.shtml#H7

Controlled Vocabulary: Subject Thesauri

Subject Thesauri and Ontologies

- <u>AGROVOC</u> (Agricultural Organization of the United Nations Vocabulary)
- O Astronomy Thesaurus
- CAB Thesaurus (for life sciences, technology and social sciences)
- O <u>CIF dictionaries</u> (for Physics)
- o Eurovoc (European Union Thesaurus)
- O Ethnographic Thesaurus
- O Gene Ontology
- O GeoNames
- O Getty Institute Art and Architecture Thesaurus Online
- O Getty Institute Thesaurus of Geographic Names
- o ICD (International Classification of Diseases)
- o Library of Congress Authorities for subject headings
- o Library of Congress Thesaurus for Graphic Materials
- O Logical Observation Identifiers Names, and Codes (LOINC)
- o MESH (Medical Subject Headings)
- O Public Health Language
- O Rare Books and Manuscripts Section (RBMS) Controlled Vocabularies
- O <u>RxNorm</u> (for drugs)
- O SNOMED CT (Systematized Nomenclature of Medicine Clinical Terms)
- o STW Thesaurus for Economics
- o **UNBIS** Thesaurus
- o **UNESCO** Thesaurus
- O USDA National Agricultural Library Agriculture Thesaurus

Question: Have you ever used thesauri in your study and research?

Controlled Vocabulary: Name Authorities

Library of Congress Name Authority File (LC/NAF)

The LC/NAF provides authoritative data for names of persons, organizations, events, places, and titles.

<u>Getty Union List of Artist Names</u> (ULAN)

The ULAN includes proper names and associated information about artists. *Artists* may be either individuals (persons) or groups of individuals working together (corporate bodies). Artists in the ULAN generally represent creators involved in the conception or production of visual arts and architecture.

Virtual International Authority File (VIAF)

The VIAF[™] (Virtual International Authority File) combines multiple name authority files into a single OCLC-hosted name authority service. The goal of the service is to lower the cost and increase the utility of library authority files by matching and linking widely-used authority files and making that information available on the Web.

Linked Data: Metadata Structure Standards

Resource Description Framework (RDF)

RDF is a standard model for data interchange on the Web. RDF extends the linking structure of the Web to use URIs to name the relationship between things as well as the two ends of the link (this is usually referred to as a "triple"). Using this simple model, it allows structured and semi-structured data to be mixed, exposed, and shared across different applications.

MADS/RDF

The Metadata Authority Description Schema (MADS) is an XML schema for an element set that may be used to provide metadata about authorized forms of agents (people, organizations), events, and terms (topics, geographics, genres, etc.). MADS/RDF builds on MADS/XML as a knowledge organization system.

Linked data examples:

- FAST: Faceted <u>Application of</u> <u>Subject</u> Terminology;
- <u>Dewey Decimal</u> <u>Classification;</u>
- <u>Open Metadata</u> <u>Registry (RDA</u> vocabularies)

...

Library of Congress
 Linked Data
 Service

SKOS: Simple Knowledge Organization for the Web SKOS is a W3C recommendation designed for representation of thesauri, classification schemes, taxonomies, subjectheading systems, or any other type of structured controlled vocabulary.

Web Ontology Language (OWL)

The OWL 2 Web Ontology Language, is an ontology language for the Semantic Web with formally defined meaning. OWL 2 ontologies provide classes, properties, individuals, and data values and are stored as Semantic Web documents. OWL 2 ontologies can be used along with information written in RDF, and OWL 2 ontologies themselves are primarily exchanged as RDF documents.

Curation Tools: Data Preparation



Colectica for Microsoft Excel

A free tool to document your spreadsheet data using the Data Documentation Initiative (DDI) metadata format, the open standard for data documentation.

http://www.colectica.com/software/colecticaforexcel

QualAnon DSDR Qualitative Data Anonymizer Results

Transcript file: ZipTest.zip Name key file : NameKeyBatch.xls

Processing Bellos.rtf...

Akiey changed to Pepulani 8 time(s)... Beatrice changed to Luwima 3 time(s)... Dende changed to Luwima 3 time(s)... Chande changed to Malida 1 time(s)... Bafen changed to Malida 1 time(s)... Bava changed to Malida 1 time(s)... Staily changed to Shour 2 time(s)... Staily changed to Shour 2 time(s)...

QualAnon: DSDR Qualitative Data Anonymizer

This free transcript anonymization tool is designed solely to deidentify qualitative interview transcripts.

https://www.icpsr.umich.edu//icpsrweb/ DSDR/tools/anonymize.jsp

Escat / Filter	Into / Redo a		52	00	row	15		
Refresh	Reset All	Remove All	Sh	ow a	s: ro	ws records	Show: 5 10 25 50	rows
Type of Contrac		chirosi		AIE		· Contract ID	· Contractor Name	· Type of Contro
815 choices Son by	name count	Cluster			۹.	1930	ASAP SOFTWARE	Microsoft Enterprise Agreement
FFAA: Fiscal/Finar Agreement 3 FFIP 1	cial Agent	adt relate			2.	1940	BNC SOFTWARE DISTRIBUTION	Remedy Service Di Maintenance
FFP ME		0			3.	1941	GOVECNNECTION INCORPORATED	Cisco SmartHet
FFP (OPS) =					4.	1942	ITS CORPORATION	Time & Matariala
FFP (Power Supply DTFA01-92-D00004	Retrofit) Old #		2		۶.	7400	SENET INTERNATIONAL CORPORATIO	Firm Fixed Price Ci
FFP BPA		-			е.	1945		fem fixed price
					7.	1946	IT FEDERAL SALES	first fixed name

OpenRefine (ex-Google Refine) is a powerful tool for working with messy data, cleaning it, transforming it from one format into another, extending it with web services, and linking it to databases like Freebase. http://openrefine.org/



The Party of Street and Designation		TURE Last
Other and Provent	Trive Barrett	March Although
20.00	10 E	10A MART
	-2	ten Lapoberti
91	100310	un Dada an
(pade	1.00	AND BRAND
(April	1026	and the second second
Totalin .	5797	and the statement
Dates Lades	12817	113.000
and and	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	A DECEMBER OF THE OWNER OWNER OF THE OWNER OWNE
That	5279	
to balled	12005	223.00808-
These 7. To 1	200	124.404.00
and the second second		125.840
Sector and	wall?	122 544
21	20100	120 10.001
ET: Post	1,000	175 Fild
ETS To Represe	200	the state of the state from
TOA AAA MAN	1994	201-201-2
105 100-00	5400	1.000
SCN Fasts	2000	and the second se
Sector The	1000	The second second second
The Lorge	1384	1 Contraction of the second
and the second s	200	

Nesstar Publisher is a

free advanced data management program. It can be used for the preparation of data and metadata. It's DDI compliant.

http://www.nesstar.com/soft ware/publisher.html

Curation Tools: Data Analysis and Visualization



Curation Tools: XML Editing

<oXygen/>ٌ

<oXygen/> XML Editor is an XML tool that supports all the XML schema languages. The XSLT and XQuery support is enhanced with powerful debuggers and performance profilers. You can use <oXygen/> XML Editor to work with all XMLbased technologies including XML databases, XProc pipelines, and web services. http://www.oxygenxml.com/



Schematron is a rule-based validation language for making assertions about the presence or absence of patterns in XML trees. It is a structural schema language expressed in XML using a small number of elements and XPath.

http://xml.ascc.net/resource/sche matron/schematron.html

ALTOVA°

Altova XMLSpy is an advanced XML editor for modeling, editing, transforming, and debugging XML-related technologies. http://www.altova.com/xmlspy .html

Curation Tools: Lab Management



LabTrove is a free blogging platform specifically designed for use in a research environment. It aims to serve as a highly flexible electronic notebook and data management system by integrating with a lab's data-producing instruments; researchers can describe an experiment and associate it with its data output at the time of capture, rather than annotating after the fact. http://www.labtrove.org/



Kepler is a scientific workflow modeling and management system that enables users, regardless of programming experience, to set up data analysis pipelines. The software will assemble, execute, and document theof services and scripts that scientists with largescale data use to execute research. https://kepler-project.org/

Curation Tools: Data Management Plan and Persistent ID Assignment



DMPTool is an online service to enable researchers to create data management plans now required by many funding agencies, and to receive tailored institutional guidance to help them in the process.

https://dmp.cdlib.org/



DataCite

The DataCite Consortium provides a number of services to support efforts at increasing the ease and prevalence of data citation.

http://www.datacite.org

Part IV: Thoughts and Services

A Researcher's View vs. A Curator or a Librarian's Perspective on Data Documentation

OSection II addresses data documentation more from the researcher's view?

OSection III interprets data documentation more from a curator or librarian's perspective?

• What do researchers really care about?

•Will each party see the other side's points and emphases?

A life cycle approach

CDL Curation and Publishing Services

http://www.cdlib.org

Create, edit, share, and save data management plans

Open source add-in for Microsoft Excel as a data collection tool

Create and manage persistent identifiers

Curation repository: store, manage, and share research data

Open access scholarly publishing services: papers, journals, books, seminars & more

An infrastructure to publish and get credit for sharing research data



*UC3DCXL

*UC3EZID

*** UC3**Merritt



eScholarship University of California



Data Publication University of California

* This slide is by Joan Starr, California Digital Library. http://www.slideshare.net/joanstarr/dataset-metadata-tools-approaches-for-access-preservation?from_search=1



Data Set Related Services

http://library.ucf.edu/ScholarlyCommunication/UCFResearchLifecycle.pdf

Dataset Metadata Service at UCF

O"Data Set (also called 'Dataset') Metadata" provides researchers consultation on:

OProject and dataset documentation;

OMetadata standards (Common and Domain Specific);

OMetadata schemas customization;

OControlled vocabularies and thesauri;

OData curation tools and practices.

 Assists in describing basic properties of your data and enriching metadata for your datasets;

 Supports applying controlled vocabularies or optimizing keywords to enhance the search of your datasets;

OHelps to prepare your metadata and data for deposit and preservation.

Related Library Resources

OScholarly Communication: (http://library.ucf.edu/ScholarlyCommunication/)

OSC Contact Information (<u>http://library.ucf.edu/ScholarlyCommunication/Contact.php</u>)

OUCF Library Research Guides (<u>http://guides.ucf.edu</u>)
 OMetadata Guide (<u>http://guides.ucf.edu/metadata</u>)
 OData Management Guide (<u>http://guides.ucf.edu/data</u>)

OResearch and Information Services

(http://library.ucf.edu/Reference/)

OSubject Librarians (<u>http://library.ucf.edu/SubjectLibrarians/</u>)

Appendix 1: Humanities Data in TEI Format

Overall structure of an ENRICH-conformant XML document. ENRICH is "European Networking Resources and Information concerning Cultural Heritage." Examples from "The ENRICH Schema — A Reference Guide." The guide is a conformant subset of Release 1.4 of TEI P5.

<TEI>

<teiHeader>

<!-- ... metadata describing the manuscript --> </teiHeader>

<facsimile>

<!-- ... metadata describing the digital images --> </facsimile>

<text>

<!-- (optional) transcription of the manuscript --> </text>

</TEI>

The minimal required structure for teiHeader: <teiHeader>

<fileDesc>

<titleStmt>

- <title>[Title of manuscript]</title>
- </titleStmt>

<publicationStmt>

- <distributor>[name of data provider]</distributor> <idno>[project-specific identifier]</idno>
- </publicationStmt>

<sourceDesc>

<msDesc xml:id="ex5" xml:lang="en">

- <!-- [full manuscript description]-->
 - </msDesc>
 - </sourceDesc>
 - </fileDesc>

<revisionDesc>

<change when="2008-01-0

- <!-- [revision information] </change>
- </revisionDesc>
- </teiHeader>

<teiHeader> (TEI

header) supplies the descriptive and declarative information making up an electronic title page prefixed to every TEI-conformant text.

Appendix 1: Humanities Data in TEI Format (Cont.)

msDesc (manuscript

description) provides

detailed information

about a single

manuscript.

settlement>Oxford</settlement>
settlement>Oxford</settlement>
settlement>Add. A. 61</idno>
settlement;

<idno>28843</idno>

</altIdentifier>

</msldentifier>

msContents>

Examples from ENRICH (http://projects.oucs.ox.ac.uk/ENRICH/ Deliverables/referenceManual en.html)

<D> <quote xml:lang="lat">Hic incipit Bruitus Anglie,</quote> the <title xml:larg="lat">De origine et gestis Regum Angliae</title> of Geoffrey of Monmouth (Galfridus Monumetensis): beg. <quote xml:lang="lat">Cum mecum multa & amp; de multis.</quote> In Latin. </msContents> <physDesc> <material>Parchment</material>: written in more than one hand: $7\frac{1}{4} \times 5\frac{3}{8}$ in., i + 55 leaves, in double columns: with a few coloured capitals. </physDesc> <history</pre> Vritten in <origPlace>England</origPlace> in the <origDate>13th cent.</origDate> On fol. 54v very faint is <quote xn/l:lang="lat">Iste liber est fratris guillelmi de buria de ... Roberti ordinis fratrum Pred[icatorum],</guote> 14th cent. (?): <guote>hanauilla</guote> is written at the foot of the page (15th cent.). Bought from the rev. W. D. Macray on March 17, 1863, for £1 10s. </history> </msDesc>



The official TEI P5 guideline is at: http://www.tei-c.org/release/doc/tei-p5doc/en/Guidelines.pdf

More TEI projects and examples are available at the TEI website: <u>http://www.tei-</u> c.org/Activities/Projects/

Appendix 2: Biological Science Dataset

	dc.contributor.author	Crawford, Nicholas G.
	dc.contributor.author	Faircloth, Brant C.
	dc.contributor.author	McCormack, John E.
	dc.contributor.author	Brumfield, Robb T.
	dc.contributor.author	Winker, Kevin
	dc.contributor.author	Glenn, Travis C.
	dc.date.accessioned	2012-05-18T15:48:08Z
(dc.date.available	2012-05-18T15:48:08Z
	dc.date.issued	2012-05-16
	dc.identifier	doi:10.5061/dryad.75nv22qj
	dc.identifier.citation	Crawford NG, Faircloth BC McCormack JE, Brumfield RT Winker K, Glenn TC (2012) Mor than 1000 ultraconserved element provide evidence that turtles ar the sister group of archosaurs Biology Letters 8(5): 783-786.
(dc.identifier.uri	http://hdl.handle.net/10255/dryad.3 8214
	dc.description	We present the first genomic-scal analysis addressing the phylogenetic position of turtles using over 1,000 loci from representatives of all major reptile lineages including tuatara
	relation.haspart	doi:10.5061/dryad.75nv22qj/1
5	velation.haspart	doi:10.5061/dryad.75nv22qj/2
1 k	relation.haspart	
NAME AND ADDRESS OF TAXABLE PARTY.		

DRYAD



Dryad (<u>https://datadryad.org/</u>)

This is an example of full metadata view.

http://www.datadryad.org/handle/ 10255/dryad.38214?show=full

Appendix 2: Biological Science Dataset (Cont.)

_	dc.relation.isreferencedby	doi:10.1098/rsbl.2012.0331
	dc.relation.isreferencedby	PMID:22593086
	dc.subject	ultraconserved elements
(dc.subject	phylogenomic
	dc.subject	phylogenetics
	dc.subject	reptiles
	dc.subject	turtles
	dc.subject	evolution
	dc.subject	archosaurs
(dc.title	Data from: More than 1000
		ultraconserved elements
		are the sister group of
		archosaurs
	dc.type	Article
	dwc.ScientificName	Pantherophis guttata
(dwc.ScientificName	Pelomedusa subrufa
	dwc.ScientificName	Chrysemys picta
	dwc.ScientificName	Alligator mississippiensis
	dwc.ScientificName	Crocodylus porosus
	dwc.ScientificName	Sphenodon tuatara
	dwc.ScientificName	Gallus gallus
	dwc.ScientificName	Taeniopygia guttata
	dwc.ScientificName	Anolis carolinensis
	dwc.ScientificName	Homo sapiens
(dc.contributor.corresponding	Faircloth, Brant C.
	prism.publicationName	Biology Letters



Dryad (https://datadryad.org/)

- It is built upon the opensource <u>DSpace</u> repository software;
- It utilizes a combination of Dublin Core (DC) and Darwin Core (DwC) metadata standards.
- Digital Object Identifiers
 (DOIs) provided by
 DataCite through EZID.

Appendix 2: Biological Science Dataset (Cont.)

Files in this package

this data. (C) 2650 COPEN COPE

		Files in this nackage
Title	turtles-all-probes.fasta	i ites in this package
Downloaded	62 times	Title
Description	Fasta file of probe sequences commercially synthesized to target UCE loci ir reptiles. Probes are identical to those used in Faircloth et al. 2012 (Systemati Biology; doi: 10.1093/sysbio/sys004)	Downloaded
Download	turtles-all-probes.fasta (865.4Kb)	Description
Details	View File Details	Download
		Details
Title	turtles-contigs-enriched-from-taxa.fasta	Details
Downloaded	63 times	•••
Description	FASTA file of contigs assembled from raw reads generated by sequencing enriched libraries of several reptile species.	
Download	turtles-contigs-enriched-from-taxa.fasta (11.31Mb)	
Details	View File Details	
Title	turtles-lastz-matches-to-reptiles.tar.bz2	
Downloaded	58 times	
Description	LAST7 matches of probes to species-specific contins in turtle-contins-enriche	d_

Downloaded	58 times
Description	LASTZ matches of probes to species-specific contigs in turtle-contigs-enrif
Download	turtles-lastz-matches-to-reptiles.tar.bz2 (360.4Kb)
Details	View File Details

Title	turtles-probe-matches-to-reptiles.sqlite
Downloaded	49 times
Description	SQLITE database of LASTZ matches (from turtles-lastz-matches- to-reptiles.tar.bz2) between assembled contigs in turtle-contigs-enriched- from-taxa.fasta and the UCE probes in turtles-all-probes.fasta. This database is used to construct data sets containing loci that are shared across taxa and prep data for alignment.
Download	turtles-probe-matches-to-reptiles.sqlite (496.6Kb)
Details	View File Details

Appendix 2: Biological Science Dataset (Cont.)

O If clicking View File Details, it displays:

turtles-all-probes.fasta
When using this data, please cite the original article:
Crawford NG, Faircloth BC, McCormack JE, Brumfield RT, Winker K, Glenn TC (2012) More than 1000 ultraconserved elements provide evidence that turtles are the sister group of archosaurs. Biology Letters 8(5): 783-786. <u>doi:10.1098/rsbl.2012.0331</u>
Additionally, please cite the Dryad data package:
Crawford NG, Faircloth BC, McCormack JE, Brumfield RT, Winker K, Glenn TC (2012) Data from: More than 1000 ultraconserved elements provide evidence that turtles are the sister group of archosaurs. Dryad Digital Repository. <u>doi:10.5061/dryad.75nv22qj</u>
Cite Share
Show Full Metadata
Files in this item
Name: turtles-all-probes.fasta View/Open Size: 845.1Kb Format: application/x-fasta Description: dataset-file Checksum (MD5): d802f4f976ef085d63a7eab10d435b03
To the extent possible under law, the authors have waived all copyright and related or neighboring ri

Simple View

DOI	doi:10.5061/dryad.75nv22qj/1
Pageviews	131
Downloaded	62 times
Keywords	ultraconserved elements, phylogenomic, phylogenetics, reptiles, turtles
Date Submitted	2012-05-18T15:48:08Z
Scientific Names	Pantherophis guttata, Pelomedusa subrufa, Chrysemys picta, Alligator mississippiensis, Crocodylus porosus, Sphenodon tuatara, Gallus gallus, Taeniopygia guttata, Anolis carolinensis, Homo sapiens
Scientific Names Contained in Data Package	Pantherophis guttata, Pelomedusa subrufa, Chrysemys picta, Alligator mississippiensis, Crocodylus porosus, Sphenodon tuatara, Gallus gallus, Taeniopygia guttata, Anolis carolinensis, Homo sapiens Data from: More than 1000 ultraconserved elements provide evidence that turtles are the sister group of archosaurs.

Description

Fasta file of probe sequences commercially synthesized to target UCE loci in reptiles. Probes are identical to those used in Faircloth et al. 2012 (Systematic Biology; doi: 10.1093/sysbio/sys004)

Appendix 3: Geospatial Dataset



http://www.geoplatform.gov/node/243/bf5a5c64-085e-4c68-a489-93e8608d3ad1

Content Standard for Digital Geospatial Metadata (CSDGM)

(http://www.fgdc.gov/m etadata/geospatialmetadata-standards)

It is maintained by the Federal Geographic Data Committee (FGDC). Often referred to as the **"FGDC Metadata** Standard."

Geospatial Platform: An Internet-based capability providing shared and trusted geospatial data, services, and applications for use by the public and by government agencies and partners to meet their mission needs.

Appendix 3: Geospatial Dataset (Cont.)

Biological data of field activity 08CRD01 (B-1-08-VI) in U.S. Virgin Islands from 05/30/2008 to 06/13/2008

Metadata

File Identifier:

Metadata Language: eng; USA: utf8

Resource Type: Dataset

Responsible Party:

Individual Name: Clint Steele < http://walrus.wr.usgs.gov/staff/csteele.html>

Organisation Name: U.S. Geological Survey (USGS) <http://www.usgs.gov>, Coastal and Marine Geology (CMG) <http://walrus.wr.usgs.gov>

Position Name: InfoBank Group Leader http://walrus.wr.usgs.gov/staff/csteele.html

Role: Point Of Contact

Contact Info: ...

Metadata Date: 2013-03-03

Metadata Standard Name: ISO 19115-2 Geographic Information - Metadata - Part 2: Extensions for Imagery and Gridded Data

Metadata Standard Version: ISO 19115-2:2009(E)

FGDC/CSDGM

Metadata

Appendix 3: Geospatial Dataset (Cont.)

Data Identification

Abstract: United States Geological Survey, Saint Petersburg, Florida, Center for Coastal and Watershed Studies... FGDC/CSDGM Purpose: These data and information are intended for science researchers, students... Language: eng; USA Metadata Citation: Title: Biological data of field activity 08CRD01 (B-1-08-VI) in U.S. Virgin Islands from 05/30/2008 to 06/13/2008 Date: Date: 2013-03-03 Date Type: Publication Date Organisation Name: U.S. Geological Survey (USGS) http://www.usgs.gov, Coastal and Marine Geology (CMG) <http://walrus.wr.usgs.gov> Role: Publisher Contact Info: ... Point Of Contact: ... **Representation Type:** Vector **Topic Category**: **Keyword Collection:** Keyword: EARTH SCIENCE > OCEANS Associated Thesaurus: Global Change Master Directory (GCMD) *Keyword*: Marine Geology Associated Thesaurus: USGS CMG InfoBank **Spatial Extent:** West Bounding Longitude: -65.75000 East Bounding Longitude: -63.25000 North Bounding Latitude: 18.75000 South Bounding Latitude: 17.25000

Appendix 3: Geospatial Dataset (Cont.)

Constraints: Please recognize the U.S. Geological Survey (USGS) as the source of this information. Physical materials are under controlled on-site access. Some USGS information accessed through this means may be preliminary in nature and presented without the approval of the Director of the USGS...

Legal Constraints:

Use Constraints: Other Restrictions

Other Constraints: Use Constraints: Please recognize the U.S. Geological Survey (USGS) as the source of this information. Physical materials are under controlled on-site access...

•••

Distribution

Distribution Format:

Format Name: ASCII

Format Version:

File Decompression Technique: No compression applied

Transfer Options:

URL: http://walrus.wr.usgs.gov/infobank/b/b108vi/html/b-1-08-vi.nav.html

Distributor:

Distributor Contact: ...

Quality

Scope: Dataset

FGDC/CSDGM Metadata

Appendix 3: Geospatial Dataset (in XML View)

- <metadata></metadata>			
- <idinfo></idinfo>			
- <citation></citation>		Idin	
- <citeinfo></citeinfo>			
- <origin></origin>			
U.S. Geological Survey (USGS) http://www.usgs.gov/ >, Coastal a			
<pre><pubdate>20130303</pubdate></pre>			
- <title></title>			
Biological data of field activity 08CRD01 (B-1-08-VI) in U.S. Virgin			
+ <pre>cnubinfo></pre> /pubinfo>			
+ <onlink></onlink>			
+ <onlink></onlink>	Top level elemen	ts:	
+ <onlink></onlink>	idia facilidantificatio		
	idinto: identificatio	n	
	Information;		
- <descript></descript>	dataqual. Data Qua	lity	
+ <abstract></abstract>		urty	
+ <purpose></purpose>	Information;		
+ <supplinf></supplinf>	spdoinfo: Spatial D	ata	
	Organization		
+ <timeperd></timeperd>	Organization		
+ <status></status>	Information;		
+ <spdom></spdom>	spref: Spatial Refer	ence	
+ <keywords></keywords>	spiel. spatiat kerei	CIICC	
+ <usesenst>/usesenst ></usesenst>	Information;		
+ <pre>sptcontac></pre> / useconst/	eainfo: Entity and		
<native>Digital ASCII</native>	Attribute Informati	00.	
+ <crossref></crossref>	Attribute information	011,	
	distinfo: Distributio	n	
+ <dataqual> (dataqual></dataqual>	Information.		
+ <spdoinfo></spdoinfo>		_	
+ <spref></spref>	metalnio: Metadata	a	
+ <eainfo></eainfo>	Reference Informat	ion.	
+ <distinfo></distinfo>			
<metainfo>//metainfo></metainfo>			

CSDGM Fields (under idinfo): Idinfo

Citation

citeinfo Origin Pubdate Title Pubinfo Onlink Descript Abstract Purpose Supplinf Timeperd Status Spdom Keywords Accconst Useconst Ptcontac Native Crossref

Content Standard for Digital Geospatial Metadata (CSDGM) Record in XML View

Appendix 4: Earth Science Dataset



http://gcmd.gsfc.nasa.gov/KeywordSearch/M etadata.do?Portal=langley&KeywordPath=Par ameters%7CATMOSPHERE%7CAIR+QUALITY%7C CARBON+MONOXIDE&OrigMetadataNode=GCM D&EntryId=MOP034&MetadataView=Full&Meta dataType=0&lbnode=mdlb1

(Click for Interactive Map) Spatial coordinates N: 90.0 S: -90.0 E: 180.0 W: -180.0

Temporal Coverage Start Date: 2000-03-03

Appendix 4: Earth Science Dataset (Cont.)

Location Keywords GEOGRAPHIC REGION > GLOBAL

Science Keywords

ATMOSPHERE > AIR QUALITY > CARBON MONOXIDE ATMOSPHERE > ATMOSPHERIC CHEMISTRY > CARBON AND HYDROCARBO CARBON MONOXIDE

ISO Topic Category CLIMATOLOGY/METEOROLOGY/ATMOSPHERE

Platform <u>TERRA > Earth Observing System, TERRA (AM-1)</u>

Instrument MOPITT > Measurements Of Pollution In The Troposphere.

Project

EOSDIS > Earth Observing System Data Information System description MOPITT > Measurements Of Pollution In The Troposphere, description ESIP > Earth Science Information Partners Program description CWIC > CEOS WGISS Integrated Catalog description

Ancillary Keywords EOSDIS

Data Set Progress IN WORK

Data Center

Atmospheric Science Data Center, Science Directorate, Langley Research Cer Data Center URL: http://eosweb.larc.nasa.gov/ Dataset ID: MOPITT Gridded Daily CO Retrievals V004

Data Center Personnel

Name: <u>ASDC USER SERVICES</u> Phone: 757-864-8666 Email: support-asdc at earthdata.nasa.gov Contact Address: NASA Langley Atmospheric Science Data Center User and Data Services NASA Langley Research Center Mail Stop 157D City: Hampton Province or State: VA Postal Code: 23681-2199 Country: USA

Click to view more

Personnel

ASDC USER SERVICES Role: DIF AUTHOR Phone: 757-864-8656 Email: support-asdc at earthdata.nasa.gov Contact Address: NASA Langley Atmospheric Science Data Center User and Data Services NASA Langley Research Center Mail Stop 157D Citru Licentea Labels: Location Keywords Science Keywords ISO Topic category Platform Instrument Project Ancillary Keywords **Data Set Progress** Data Center Personnel **Extended Metadata Properties Creation and Review Dates** ...

Extended Metadata Properties V (Click to hide)

Extended Metadata Properties From GCMD UUID: 37d9e4f6-a5ba-404d-9ee3-0e3025a2ca9c

Creation and Review Dates DIF Creation Date: 2013-06-10 Last DIF Revision Date: 2013-07-19

- <u>Reformat as FGDC document</u>
 Reformat as ISO 19115 document
- View Text Only Format

Directory Interchange Format (DIF):

a descriptive and standardized format for exchanging information about scientific data sets. The DIF Writer's Guide: http://gcmd.gsfc.nasa.gov/U

ser/difguide/difman.html.

Origin: DIF was the product of an Earth Science and Applications Data Systems Workshop (ESADS) held February 24-26, 1987 on catalog interoperability (CI). (http://gcmd.gsfc.nasa. gov/add/difguide/whatisadif. html) Contact: Sai Deng, Metadata Librarian and Associate Librarian <u>sai.deng@ucf.edu</u> 407-823-4312 (Office)

Thank you!