

DM Vitals: A Data Management Assessment Recommendations Tool

Sherry Lake, Senior Scientific Data Consultant, shlake@virginia.edu
Scientific Data Consulting Group, University of Virginia, Library,
University of Virginia, Charlottesville, Virginia, U.S.A.

Susan Borda, Graduate Student Summer Intern, shborda@syr.edu
School of Information Studies,
Syracuse University, Syracuse, New York, U.S.A



Abstract:

In an effort to develop an understanding of how science and engineering researchers at the University of Virginia (UVA) manage their research data, UVA Library's Scientific Data Consulting (UVA SciDaC) group began a series of research data interviews. The goals of the data interview process include identifying common research data problems, identifying research support needs, and providing recommendations on improving data management. In practice, however, providing objective suggestions for data management practices proved to be troublesome. It was difficult to make reliable customization recommendations; be objective in a timely fashion.

In response to these challenges, the UVA SciDaC group developed a system (DM Vitals) that would easily and objectively rate the current state of the researcher's data management practices. Using best practice statements from UVA sources (Information Technology Services' Risk Management Program and SciDaC guidelines) and the Australian National Data Service's (ANDS) long-term sustainability scoring model, the system compares the information collected during the data interview process with these data management best practice statements. The model then further correlates the researcher's data management practices with the eight data management practice components developed by the SciDaC group: File Formats and Data Types, Organizing Files, Security/Storage/Backups, Funding Guidelines, Copyright & Privacy/Confidentiality, Data Documentation & Metadata, Archiving & Sharing and Citing Data.

To provide a framework for comparing and improving departmental data management practices, we took the value resulting from the average of the data management best practice statements and compared them the Crowston and Qin Capability Maturity Model (CMM). Using this model as a basis, the data management maturity levels are defined as: Level 0: Initial (this includes current practices that can be seen as counter-productive or even "risky" from a security standpoint), Level 1: Managed (the researcher begins to uniformly apply some of the lower level/easier best practices, really starting to "manage" the situation), Level 2: Defined (the researcher is further "defining" their DM practices), Level 3: Quantitatively Manage (the researcher begins to use central and outside services to manage their data), and Level 4: Optimizing (the researcher are continually improving their data management practices).

The strength of the DM Vitals tool is in generating tasks customized to each researcher. These tasks can then be easily grouped into phases, creating a data management implementation plan for each researcher based on their personal data interview and subsequent information gathering. The DM Vitals assessment tool differs from the Digital Curation Center's (DCC) CARDIO in that its focus is not on consensus and collaboration by individuals responsible for the research data (PI, IT, Data manager, etc.).

Once this tool is fully integrated into the existing UVA SciDaC Data Interview and Data Management Plan process, it will expedite the recommendation report process by providing valuable actionable feedback that the researcher can use immediately to improve the sustainability of their data.



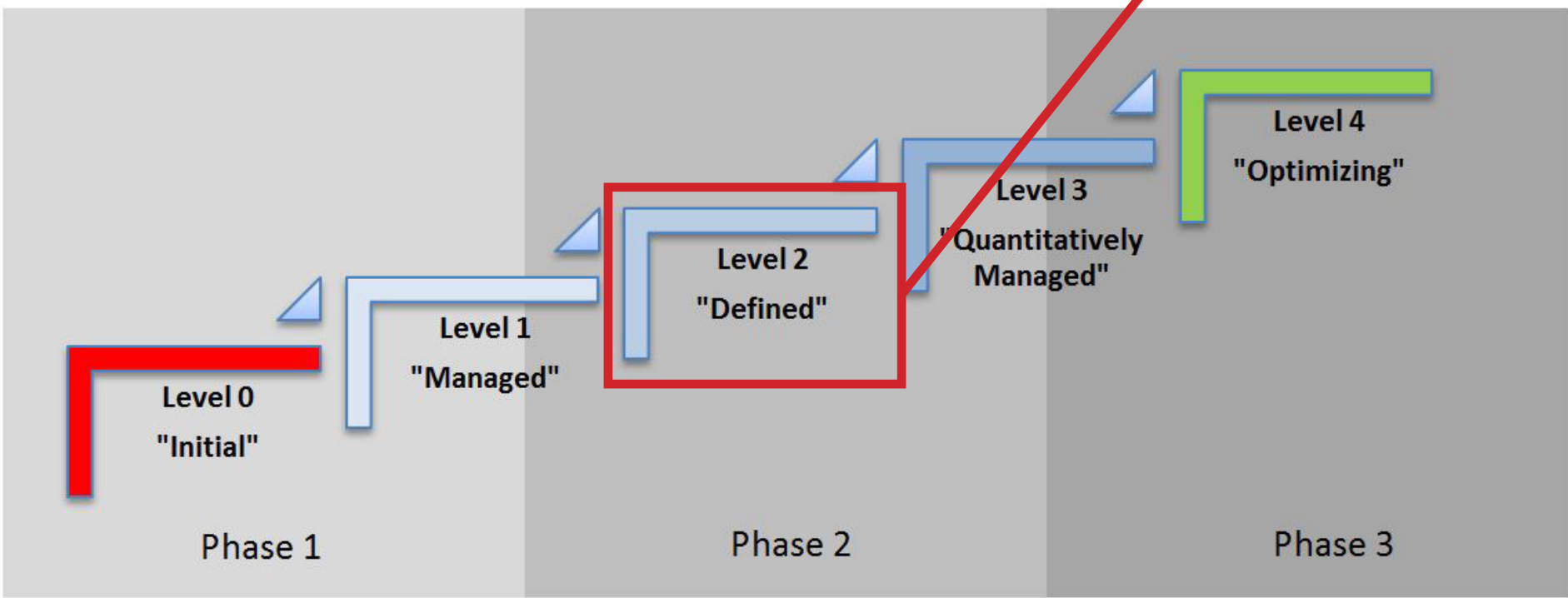
Interview Questions:

2.1 Describe your day-to-day work with regard to data. What data do you have? What kind of data, how are they created, what type and formats and what software do you use? How much data do you have?	2.2 Are these files yours or do they belong to a wider group or to the institution? Who owns the Intellectual Property rights of the data you create?	3.1 Do you have a data management plan? Who is responsible for managing the data? Are you using any filing or naming conventions for the files? How are the files organized? Is there any documentation on the files and/or data fields?	3.2 How do you share data - among lab group or other colleagues e-mail, shared drive, removable devices, CD, web pages, others? Have you had version control issues with many people working on the same data file?	4.1 What challenges have you faced in terms of storage, formats, costs, and continued access to older data?	5.1 Have you been asked to provide or share your data? Could or should your data be reused or repurposed by others, and if so, how and by whom?
2.1.1 General Category (experimental, simulation/computational, observational, ...)	no	Have you read UVA's Laboratory Notebook and Recordkeeping policy?	3.1.1 Management Plan	3.2.1 File sharing	4.1.1 Do they have older files?
2.1.2 Curation (servers, instruments, software)	no	Have you read UVA's Ownership Rights in Copyrightable Material policy?	no	3.1.2 Version control issues: Versions not maintained	4.1.2 Obsolete data formats
2.1.3 Data Type (docs, emails, databases, images, videos, etc.)	no	no	An "informal" DMP exists	4.1.3 Obsolete media	5.1.3 Restrictions (Confidentiality, Sensitivity)
2.1.4 Data Format (MS Word, Excel, spss, html, ipb, etc.)	no	no	DMP has been improved to include all 8 categories.	4.1.4 Storage space (Also see 3.1.6 & 3.1.7)	5.1.4 Documented for sharing
no Non-standard	no	no	DMP has been reviewed by SciDaC	4.1.5 Costs	
no Proprietary software file type	no	no	DMP is being followed by all research team members.		
yes Non-proprietary or non-software specific file format	yes	yes	3.1.2 Naming Conventions: Are they using file naming conventions?		
yes Standard representation (ASCII, Unicode)	yes	yes	File Naming Conventions for Specific Disciplines		
no Common, used by the research community	no	no	3.1.3 File Organization		
no Open, documented standard	no	no	2.1.5 Amount (files, files sizes, growing?)		
2.1.6 Software	no	no	Files lost or disorganized		
	yes	yes	File Structure		
	yes	yes	Use Same Structure for Backups		

- UVA's "Data Interview Question Template Sections" are split into discrete questions.
- Data management practices are added to these questions, as sub-statements with three possible answers: no, yes, or null.

Final Assessment (Given to Researcher):

	Average	
File Formats Data Types	2	Satisfactory
Organization of Files	4	More Sustainable
Security Storage Backups	3.5	Good
Copyright Privacy Confidentiality	0	Least Sustainable
Data Documentation Metadata	1.333333333	Fair
Cumulative	2.166666667	Satisfactory



- Final DMP component values and an overall sustainability grade;
- A Capability Maturity Model indicating a researcher's "grade" on the Data Management practice continuum.

References:

Australian National Data Service. (2011) ANDS and Data Storage. Available: <http://ands.org.au/guides/storage.html>. Last accessed August 30, 2011.

Crowston, K., & Qin J. (2010). A capability maturity model for scientific data management. American Society for Information Science and Technology Annual Meeting, Pittsburg, PA. Working Paper available: <http://crowston.syr.edu/content/capability-maturity-model-scientific-data-management-0>. Last accessed August 23, 2011.

Digital Curation Center. (2011). CARDIO. Available: <http://cardio.dcc.ac.uk/>. Last accessed August 30, 2011.

Information Technology Security (2010). University of Virginia Information Technology Security Risk Management (ITS-RM) Program. Available: http://its.virginia.edu/security/riskmanagement/docs/ITS-RM_3-0.pdf. Last accessed August 23, 2011

University of Virginia Library (2011). Scientific Data Consulting Data Management Home. Web Site available: <http://www2.lib.virginia.edu/brown/data/>. Last access August 30, 2011.

Data Management Practice Statement Valuation:

1 - Least Sustainable	1 - Fair	2 - Satisfactory	3 - Good	4 - More Sustainable	
Documentation	Documentation	Documentation	Documentation	Research Project Documentation	
no No DMP exists - only in researcher's mind.	no Data sources used (see Citing Data)	yes Context of data collection	no Structure, organization of data files	yes Transformations of data from the raw data through analysis	# Yes 6
no An "informal" DMP exists	yes Data collection methods	no DMP has been reviewed by SciDaC	yes DMP has been reviewed by SciDaC	yes Information on confidentiality, access & use conditions	Total Possible: 18
	no DMP has been improved to include all 8 categories.	no Data validation, quality assurance	no Dataset Documentation: Variable names, and descriptions	no DMP is being followed by all research team members.	Final Value 1.333333
	yes Dataset Documentation: File format and software (including version) used.	no Dataset Documentation: Explanation of codes and classification schemes used	no File format and software (including version) used.	no Metadata: Use a research community standard.	
	yes Understands the need for metadata	no Metadata: Informal metadata practice.			
0 Total Yes	0 Total Yes	4 Total Yes	0 Total Yes	2 Total Yes	6 Total

Each sub-statement is correlated with a "sustainability level" that rates the quality of existing DM practices from worst to best and assigns the point values to each level.

Implementation Phases and Action Statements:

Applicable?	Interview Topic	Sustainability	Phase	Action Statement
X	Data sources used	1	1	Document all data sources used.
X	An "informal" DMP exists	1	1	Have an "Informal" Data Management Plan (DMP): a DMP is the basis of all data management, and is a critical tool in protecting the continuity of your research process. Once in place, it can continually be updated, provided to new members of the lab as guidelines, and easily be applied to future grant proposals. We will work with you to develop an appropriate plan for managing your data. This is a fundamental first step in improving process.
X	Have you read UVA's Laboratory Notebook and Recordkeeping policy?	1	1	Read UVA's "Laboratory Notebook and Recordkeeping" Policy: https://policy.its.virginia.edu/policy/policydisplay?tid=RES-002
X	Have you read UVA's Ownership Rights in Copyrightable Material policy?	1	1	Read UVA's "Ownership Rights in Copyrightable Material" Policy: https://policy.its.virginia.edu/policy/policydisplay?tid=RES-001
X	DMP has been improved to include all 8 categories.	2	2	Informal DMP has been improved to include these 8 categories: File Formats and Data Types; Organizing Files; Security/Storage/Backups; Funding Guidelines; Copyright & Privacy/Confidentiality; Data Documentation & Metadata; Archiving & Sharing Data, and Citing Data.
X	Common, used by the research community	3	2	Use file types and data formats that are commonly used by your research community
X	Structure, organization of data files	3	2	Document structure and organization of data files
X	Data validation, quality assurance	3	2	Document data validation and quality assurance processes
X	Explanation of codes and classification schemes used	3	2	Document and explain any codes and classification schemes used
X	File format and software (including version) used.	3	2	Document any file formats and software (including version) used
X	Informal metadata practice	3	2	Begin applying appropriate metadata standard to data and other research materials
X	Adhering to ITC recommended backup cycles with off-site storage	3	2	Work with ITC to be sure that all servers and computers containing research data are being backed up to their specifications (with regular cycles and off-site storage)
X	DMP has been reviewed by SciDaC	3	3	Have SciDaC review DMP
X	Open, documented standard	4	3	Use an open and documented standard for file formats
X	Algorithms used to transform data	4	3	Document any algorithms used to transform data
X	In physically secure environment	4	3	Store data storage media and servers in a physically secure environment
X	Should be formally administered	4	3	Have data backups formally administered by ITC or some other system administrator
X	DMP is being followed by all research team members.	4	3	DMP is being followed by all members of the research team
X	Use a research community metadata standard.	4	3	Fully apply the appropriate research community metadata standard to research data and materials.

- Data management "action statements" are created from corresponding sub-statements with values of "no."
- Implementation phases corresponding to the complexity of implementing each action are added to each "action statement."
- Action statements are sorted according to implementation phase.

Aggregated Results: Cumulative Results for Researchers at UVA.

