

Research Data Management in Environmental Science



Ali Krzton, Research Data
Management Librarian
RBD Library

Overview

- Research data management is at the heart of the scientific method
- It has become an institutional priority
- Environmental sciences are particularly data intensive
- Changing research practices require new data management strategies
- Upfront investment in good systems has huge payoffs later
- Our valuable environmental data is worth it

But first, some history

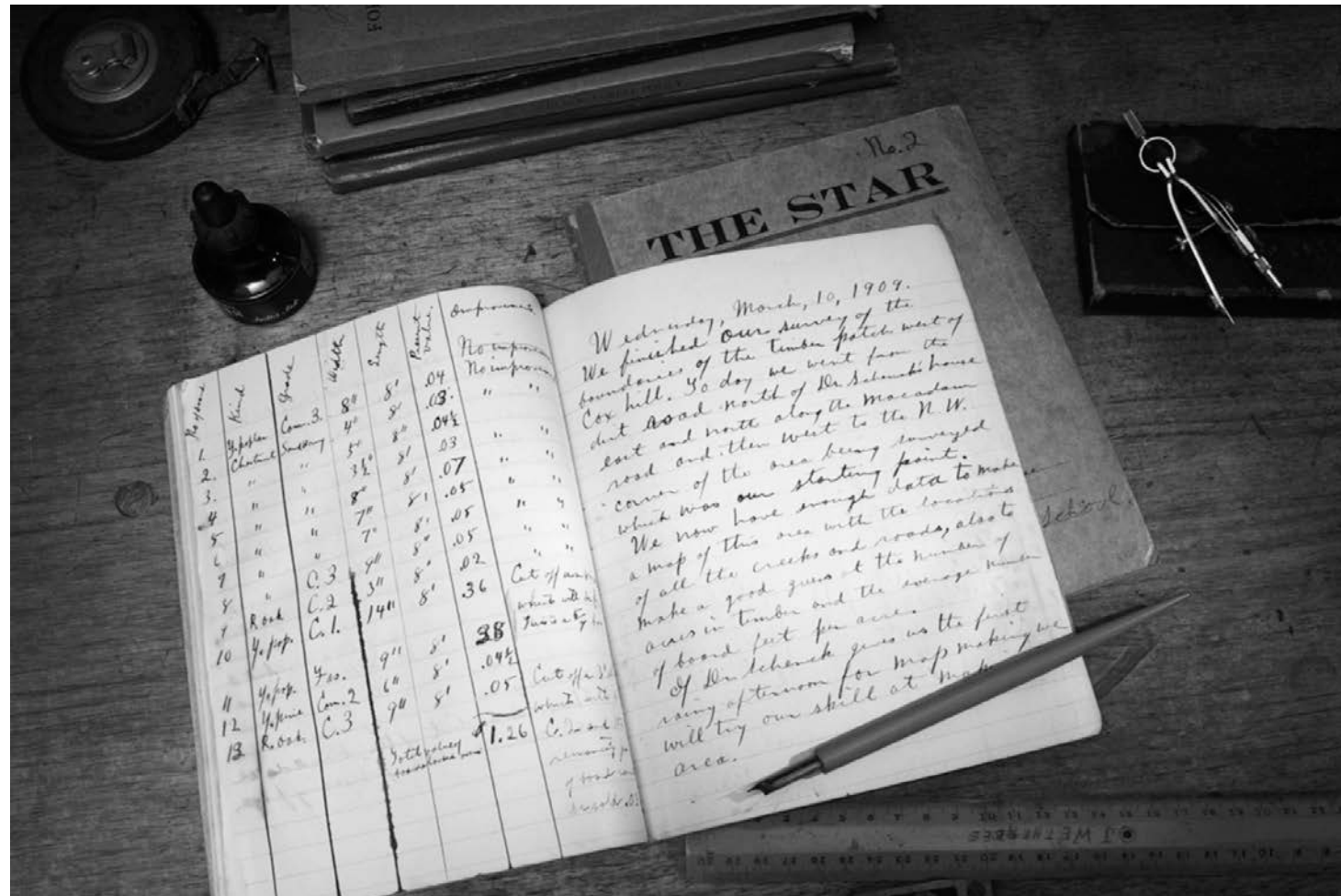


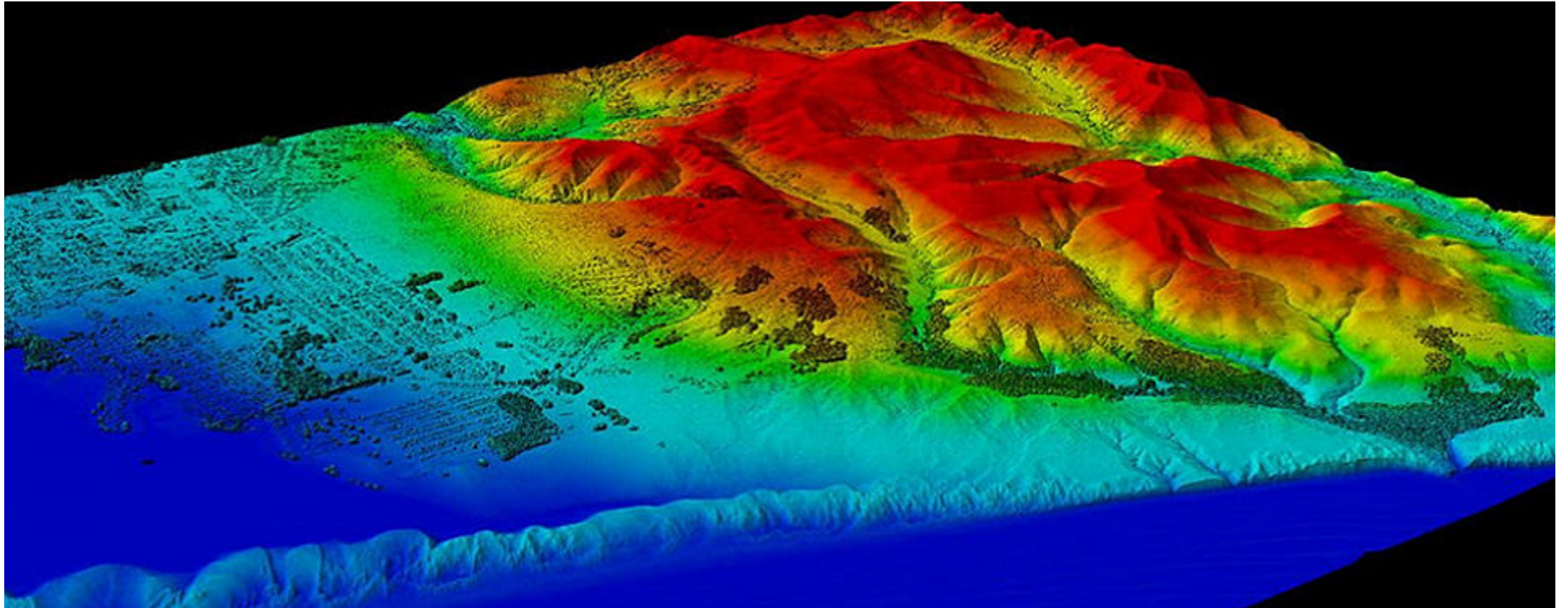
Figure reproduced from Lee (2012)

Features of environmental science

- Research conducted at a wide variety of scales – spatial, temporal, micro to macro levels of analysis
- Aggregated datasets with multiple inputs are the norm
- Data often cannot be recreated (e.g. historical data)
- Decision-makers rely upon our data

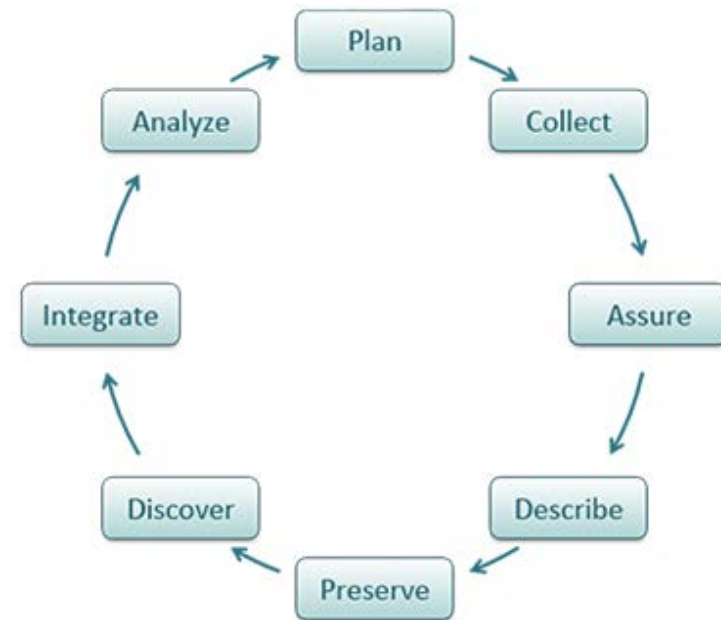


Our work is becoming more computational



Why research data management now?



- Funding mandates (NIH, NSF, OSTP)
- Journal mandates
- Reproducibility crisis
- Collaborative research models



A sample of RDM issues

- Formatting
- QA/QC
- Software interoperability
- Documentation
- Loss/corruption



ORIGINAL RESEARCH | [Open Access](#)  

Temporal degradation of data limits biodiversity research

Geiziane Tessarolo , Richard Ladle, Thiago Rangel, Joaquin Hortal

First published: 27 July 2017 | <https://doi.org/10.1002/ece3.3259> | Cited by:2

[Article Linker: Find Full Text](#)


[Read the full text >](#)

 PDF  TOOLS  SHARE

 [Figures](#)  [References](#)  [Related](#)  [Information](#)

Metrics

Citations: 2

 score 14

Details

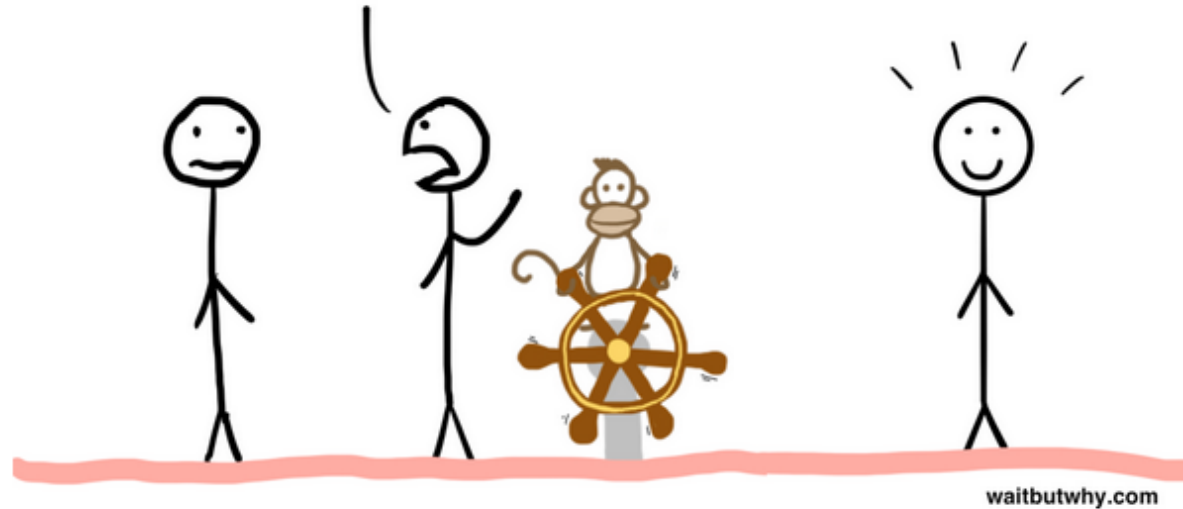
© 2017 The Authors. *Ecology and Evolution* published by John Wiley & Sons Ltd.

This is an open access article under the terms of the [Creative Commons Attribution License](#), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

Abstract

Spatial and/or temporal biases in biodiversity data can directly influence the utility, comparability, and reliability of ecological and evolutionary studies. While the effects of biased spatial coverage of biodiversity data are relatively well known, temporal variation in data quality (i.e., the congruence between recorded and actual information) has received much less attention. Here, we develop a conceptual framework for understanding the influence of time on biodiversity data quality based on three main processes: (1) the natural dynamics of ecological systems—such as species turnover or local extinction; (2) periodic taxonomic revisions, and; (3) the loss of physical and

Look at you, [REDACTED] Past Tim. You're the whole reason we're in this terrible situation now. See Future Tim over there? Why can't you be more like him? Thank god he's here to fix everything. I'd clean up your mess myself if I weren't dealing with this [REDACTED] monkey.



The biggest beneficiary of good research data management...is future you.

Solutions at the micro level

- Consistency within and between
- Talk about research data management as part of scientific procedure
- Set policies within the working group – get buy-in!
- Match training to responsibilities
- Make sure digitized data is machine-readable

Poor data entry – common errors

1	Site	Date	Plot	Species	Weight	Acult									
2	DeepWell	2/13/2010		1 DIPO	12.1	j									
3	Deep Well	Feb-10		2 Pero	13.22	j									
4	rioSalado	2/13/2010	1a	pero	16	N									
5	riuSladu	"	1*	CleGap	18.92	gut away									
6				Mean1	15.06										
7															
8															
9															
10															
11															
12	Rodent Trapping		MJK & ALN	10-Apr-10											
13	Site	Plot	Adult	Species	grams	Comments									
14	deep well		1 y	woodrat	13										
15	riosalado		2 y	PERO	24.5										
16	riosalado		3 y	Clegap	91										
17															
18															
19															
20															

Inconsistencies abound! Location and format of date information, column headings, and column order are different, as are site spellings and codes used. Information has also jumped into other columns where it doesn't belong.

Improved data organization

data.xls

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	Site	Date	Plot	Species	Weight	Adult		Rodent Trapping 3/15/2010						
2	DeepWell	2/13/2010		1 DIPO	12.1	j		Site	Plot	Adult	RodentSp	Weight		
3	Deep Well	Feb-10		2 Pero	13.22	j		DW		1 y	Pero	12		
4	rioSalado	2/13/2010	1a	pero	16	N		RS		2 j	PERO	escaped <15		
5	riuSladu	"	1*	CleGap	18.92	gut away		RS		3 n	Clegap	91		
6				Mean1	15.06									
7														
8														
9														
10														
11														
12	Rodent Trapping			MJK & ALN	10-Apr-10									
13	Site	Plot	Adult	Species	grams	Ccmmnts								
14	deep well		1 y	woodrat	13									
15	riosalado		2 y	PERO	24.5									
16	riosalado		3 y	Clegap	91									
17														
18														
19														
20														

Sheet1

SEV_SmallMammalData_v.5.25.2010.xls

	A	B	C	D	E	F	G	H
1	Date	Site	Plot	Species	Weight	Adult	Comments	
2	2/5/2010	Deep Well		1 DIPO	13.2	y		
3	2/4/2010	Deep Well		1 CLEGAP	11.6	j		
4	2/5/2010	Rio Salado		1 DIPO	14.2	y		
5	2/5/2010	Rio Salado		2 PERO	10.1	y		
6	3/15/2010	Deep Well		1 DIPO	15.2	y	plot burned	
7	3/15/2010	Deep Well		2 DIPO	21.7	y	pregnant	
8	3/15/2010	Rio Salado		1 CLEGAP	16.2	j		
9								
10								
11								
12								
13								

SmallMammalTrapping Sheet3

General rules for tabular data

- 1. Each variable must have its own column.
- 2. Each observation must have its own row.
- 3. Each value must have its own cell.

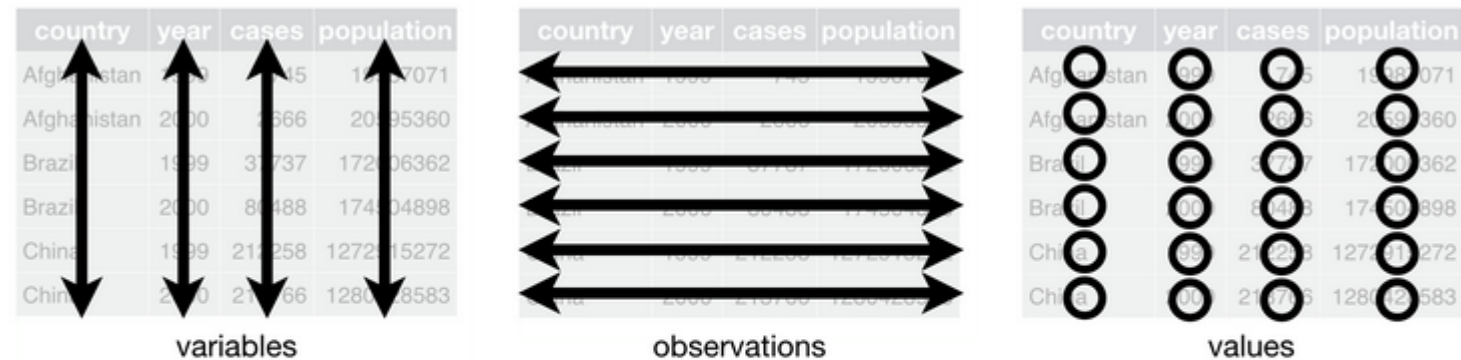


Figure 12.1: Following three rules makes a dataset tidy: variables are in columns, observations are in rows, and values are in cells.

File formats

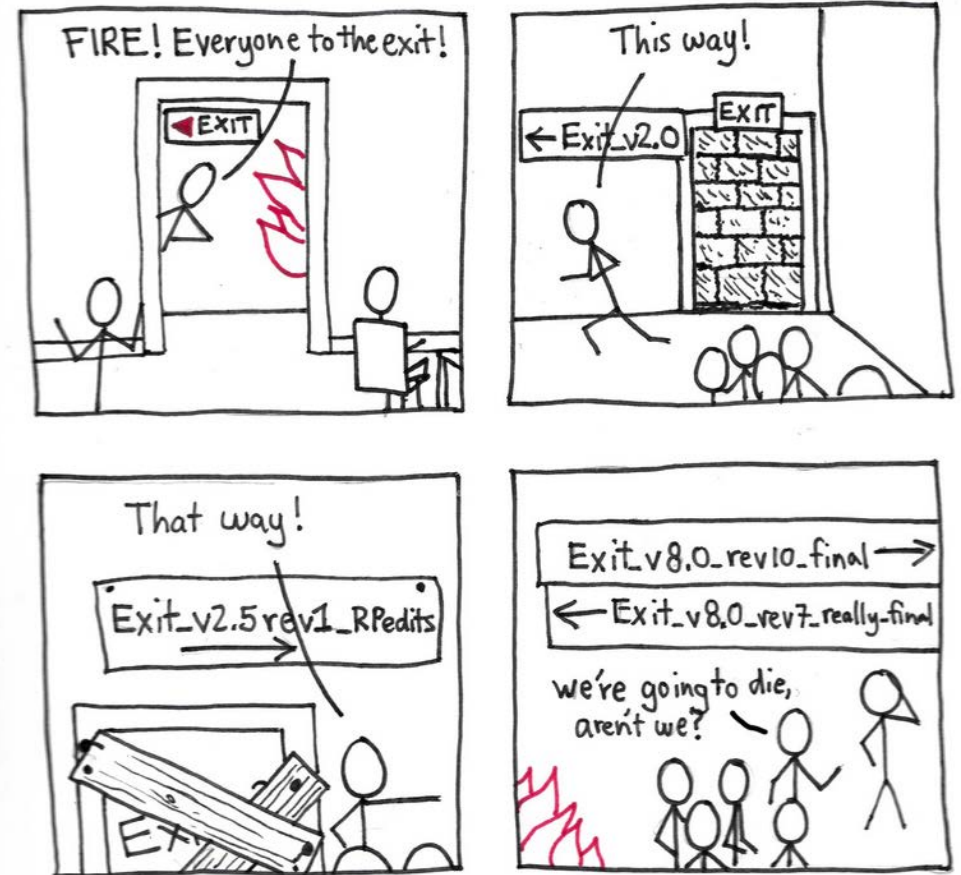
- Some file formats will be determined by equipment or software
- Proprietary file formats can cause problems
- Choose the most compatible formats when possible
- Consider saving a backup/archival copy in an open format
- READMEs should be in plain text (.txt)

File naming – best practices

- Naming needs to be descriptive and consistent
- A unique identifier system may help, but check with collaborators
- Dates and version numbers could also be included
- Shorter is usually better, and `_underscores_` or CamelCase may improve readability
- Make sure not to use special characters: offenders include hyphens, periods other than the one before the extension, spaces
- Avoid generic names and words like “new” or “final”

Backups and versioning

- Keeping multiple copies in multiple locations, if possible
- Backup and versioning processes are best automated
- Consider cloud services and/or central repositories
- Version retention: value vs. confusion
- Tools and protocols change too!



Stop, Drop, and Use a Versioning System
@redpenblackpen

Solutions at the macro level

- Working groups and institutions have developed tools designed to improve data management in environmental research
- DataONE
- Knowledge Network for Biocomplexity
- eMammal
- Community-based guidelines
- Ask your librarian – don't be shy!

DataONE Training and Education



Site DataONE Search

Search phrase

Go

About News Participate Resources Education Data

Discover

Find data from affiliated member repositories

More

Participate

Make the most of your data with supporting tools and services

More

Learn

Access materials and resources to support learning in good data management practices

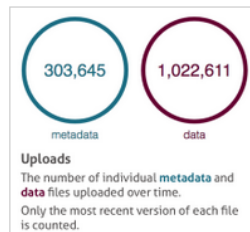
More



What is DataONE?

DataONE is a community driven project providing access to data across multiple member repositories, supporting enhanced search and discovery of Earth and environmental data. DataONE promotes best practices in data management through responsive educational resources and materials... [more about DataONE](#)

DataONE Reaches 1 Million Objects



DataONE reached a milestone the week of July 17th 2017, exceeding 1 million data objects! This number represents data uploaded to our network of 37 Member Nodes and does not include replicas. All 1 million objects are discoverable through DataONE Search.

SEARCH FOR DATA

MEMBER NODES

Tweets about @dataoneorg

christopherlortie @cjortie
Huge #CyberMonday sales on #openscience data. Cost: free! Check out @DataONEorg member nodes dataone.org/current-member ... or @figshare or the new @ESIPfed repo esip.figshare.com best kind of shopping!
Nov 27, 2017

Adam Shepherd @ashep_15
Replying to @metamattj and 3 others
Let us know if something doesn't fit or isn't extendible from your perspective. Then we can work it out together. Much appreciated!
Nov 15, 2017

Knowledge Network for Biocomplexity



Tools

Over the years, many tools have been developed to facilitate effective data management, archiving content, and retrieving data for synthetic analysis projects.

Morpho

Data management for earth, environmental and ecological scientists.

Morpho allows researchers to create metadata, (i.e. describe their data in a standardized format), and create a catalog of data & metadata upon which to query, edit and view data collections. In addition, it also provides the means to access network servers - like the KNB - in order to query, view and retrieve all relevant, public ecological data! Check the [Morpho User Guide](#) for complete details.



- Create and Edit EML metadata
- Search for existing data packages
- View and download data packages
- Verify and Edit your data
- Specify Access Control rules for your data
- Share and publish your data via the KNB

Download the Morpho data management application.

Easy-to-use installers are available for various platforms. For windows and mac versions, users need to double click the installers. For linux version, users need to run "java -jar morpho-version-linux.jar". [To run Morpho, you must have Java 1.7 or later installed on your computer.](#)

- [Windows](#) ⓘ
- [Linux](#) ⓘ
- [Mac OSX](#) ⓘ
- [Older Versions](#) »

eMammal (Smithsonian Institute)

The screenshot shows the eMammal website interface. At the top left is the eMAMMAL logo with the tagline "See Wildlife, Do Science". To the right are social media icons for Facebook and Twitter, and a "Login" button. Below this is a navigation bar with links for Home, About, View Photos, Explore Projects, Browse Data, and Resources, along with a search bar. The main content area features a "Recent Blog Posts" section with two entries: "eMammal workshop with Western Kimberley indigenous rangers" and "Camera trapping with Virginia Teachers". To the right of the blog posts is a large image of an African Buffalo with a caption "African Buffalo Syncerus caffer" and a link to "Enter Project Page". Below the main content area is a banner with the text "eMammal is a tool for collecting, archiving, and sharing camera trapping images and data". At the bottom, there are five circular icons representing different website functions: About (two people), View Photos (camera), Explore Projects (globe), Browse Data (binary code), and Get Involved (hands raised).

eMAMMAL See Wildlife, Do Science

Home About View Photos Explore Projects Browse Data Resources

Recent Blog Posts

eMammal workshop with Western Kimberley indigenous rangers
Posted by Michael Wyson on February 26, 2018

On February 22, several indigenous ranger programs from the Western Kimberley came together learn about eMammal and to test out the eMammal app on photos from the Spectacled Hare-Wallaby survey project. The workshop included an overview camera...
continue reading...

Camera trapping with Virginia Teachers
Posted by Megan Blance on February 12, 2018

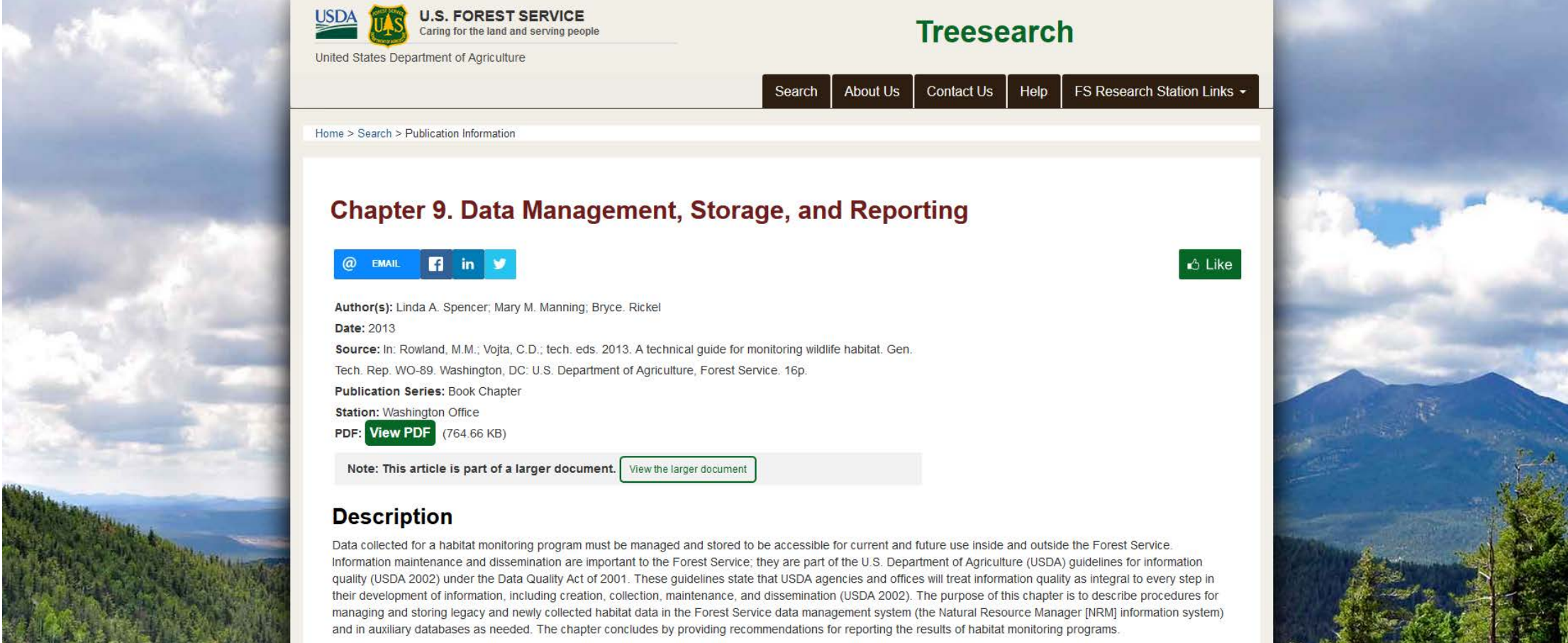
The Eyes on Wildlife project is a collaboration between the Smithsonian Conservation Biology...

African Buffalo
Syncerus caffer
Enter Project Page

eMammal is a tool for collecting, archiving, and sharing camera trapping images and data

About View Photos Explore Projects Browse Data Get Involved

Community best practices in the literature



USDA U.S. FOREST SERVICE
Caring for the land and serving people
United States Department of Agriculture

TreeSearch

Search About Us Contact Us Help FS Research Station Links ▾

Home > Search > Publication Information

Chapter 9. Data Management, Storage, and Reporting

@ EMAIL f in t Like

Author(s): Linda A. Spencer; Mary M. Manning; Bryce. Rickel
Date: 2013
Source: In: Rowland, M.M.; Vojta, C.D.; tech. eds. 2013. A technical guide for monitoring wildlife habitat. Gen. Tech. Rep. WO-89. Washington, DC: U.S. Department of Agriculture, Forest Service. 16p.
Publication Series: Book Chapter
Station: Washington Office
PDF: [View PDF](#) (764.66 KB)

Note: This article is part of a larger document. [View the larger document](#)

Description

Data collected for a habitat monitoring program must be managed and stored to be accessible for current and future use inside and outside the Forest Service. Information maintenance and dissemination are important to the Forest Service; they are part of the U.S. Department of Agriculture (USDA) guidelines for information quality (USDA 2002) under the Data Quality Act of 2001. These guidelines state that USDA agencies and offices will treat information quality as integral to every step in their development of information, including creation, collection, maintenance, and dissemination (USDA 2002). The purpose of this chapter is to describe procedures for managing and storing legacy and newly collected habitat data in the Forest Service data management system (the Natural Resource Manager [NRM] information system) and in auxiliary databases as needed. The chapter concludes by providing recommendations for reporting the results of habitat monitoring programs.

Your turn

- Questions?

Ali Krzton, Research Data Management Librarian

alk0043@auburn.edu

334-844-8268

libguides.auburn.edu/researchdata

References

- Lee, David S. (2012). The Biltmore Forest School: Poking Back into an Extraordinary Time. *The American Biology Teacher* 74:7, 464-469. doi: 10.1525/abt.2012.74.7.7
- Tessarolo, Geziane, Ladle, Richard, Rangel, Thiago, and Hortal, Joaquin. (2017). Temporal degradation of data limits biodiversity research. *Ecology and Evolution* 7:17, 6863-6870. doi: 10.1002/ece3.3259
- Urban, Tim. (2015). The Procrastination Matrix. *Wait But Why*. <https://waitbutwhy.com/2015/03/procrastination-matrix.html>
- DataONE Education Modules. <https://www.dataone.org/education-modules>
- Grolemund, Garrett, and Wickham, Hadley. (2017). Tidy Data. In: *R for Data Science*. <http://r4ds.had.co.nz/tidy-data.html>
- MIT Libraries Data Management Workshops. <https://libraries.mit.edu/data-management/services/workshops/>
- Knowledge Network for Biocomplexity – Tools. <https://knb.ecoinformatics.org/#tools>
- eMammal. Smithsonian Institute. <https://emammal.si.edu>
- Spencer, Linda A., Manning, Mary M., and Rickel, Bryce. (2013). Data Management, Storage, and Reporting. In: *A Technical Guide for Monitoring Wildlife Habitat* (tech. eds. Rowland, M.M.; Vojta, C.D.), Gen. Tech. Rep. WO-89. Washington, DC: U.S. Department of Agriculture, Forest Service. 16p. <https://www.fs.usda.gov/treearch/pubs/45226>