# Research
# Data Management

Enhanced Research Skills Certificate Program

Week 3: Understanding Data

Winter 2016

# After this workshop you will be able to:

- Be able to explain what data is, and the importance of research data management

- Be able to use a file-naming system to name and organize data files

- Describe the value of open access to data, and what is required to facilitate discovery, access, and reuse of shared data

- Begin to practice good methods of data storage and backup

- Read and interpret a basic metadata record

- Begin recording metadata to accompany the data you collect and use for research

# Video:

[Data Sharing and Management Snafu in 3 Short Acts](#)

# Reflection questions

- What were the issues preventing the researcher from reusing the data?


- Why is it important for researchers to <u>document</u>, <u>preserve</u> and <u>share</u> their data?

But what do we mean by *Data*??

# Defining 'Data'

"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 Board. 2005. Long-Lived Data Collections: Enabling Research and Education in the 21st Century. Arlington, VA: National Science Foundation, p. 13.

# What is data?

## Any information you use in your research

# Can you find what you need, when you need it?

'What a mess' by .pst, via Flickr: http://www.flickr.com/photos/psteichen/3915657914/.

# Why Manage Data: Researcher Perspective

- Manage your data for yourself:
  - Keep yourself organized
  - Track your science processes for reproducibility
  - Better control versions of data
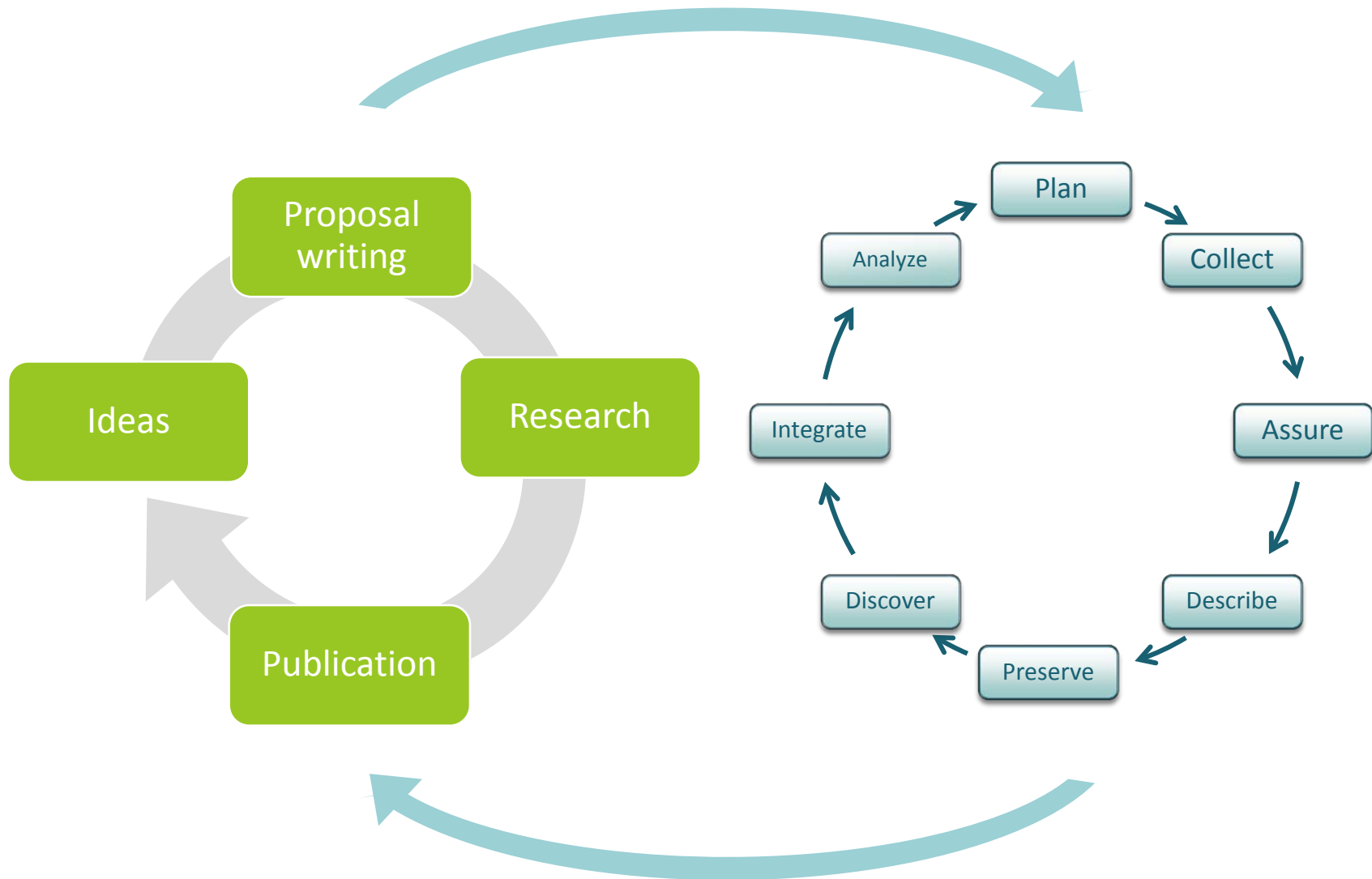  - Quality control your data more efficiently

# Why Manage Data: Researcher Perspective

- Make backups to **avoid data loss**

- Format your data **for re-use (by yourself or others)**

- Be prepared: Document your data **for your own recollection, accountability, and re-use** (by yourself or others)

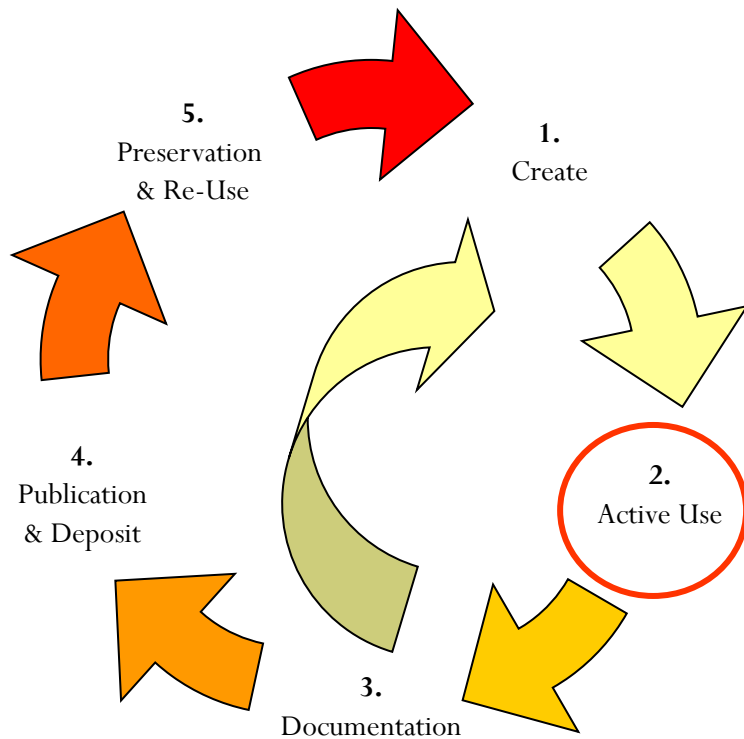- Prepare it **to share it – gain credibility and recognition** for your research efforts!

# Disciplines and Data Management

- Need for data management is universal
- Practices vary across disciplines
- Good data management extends knowledge creation across disciplines

# Research and Data Management Planning are Both Inter-related Iterative Processes

Proposal writing

Ideas

Research

Publication

Plan

Collect

Analyze

Assure

Integrate

Describe

Discover

Preserve

# How will you look after your data?

5.
Preservation
& Re-Use

1.
Create

4.
Publication
& Deposit

2.
Active Use

3.
Documentation

- Is your data safe?

- Is your data organised?

- Can you find your data?

Slide adapted from the
PrePARe Project

# Storage and Security



*Photo credits: Harvey Rutt*
*http://www.ecs.soton.ac.uk/regenesis/pictures/*

- **3… 2… 1…  Backup!**
  - at least **3 copies** of a file
  - on at least **2 different media**
  - with at least **1 offsite**
  - Do not use Dropbox and be very cautious of cloud storage for sensitive data.

- **Access**
  - Protect your hardware
  - For sensitive data, use file encryption (Talk to your advisor)
  - Keep passwords safe (e.g. Keepass)
  - At least **2 people** should have access to your data
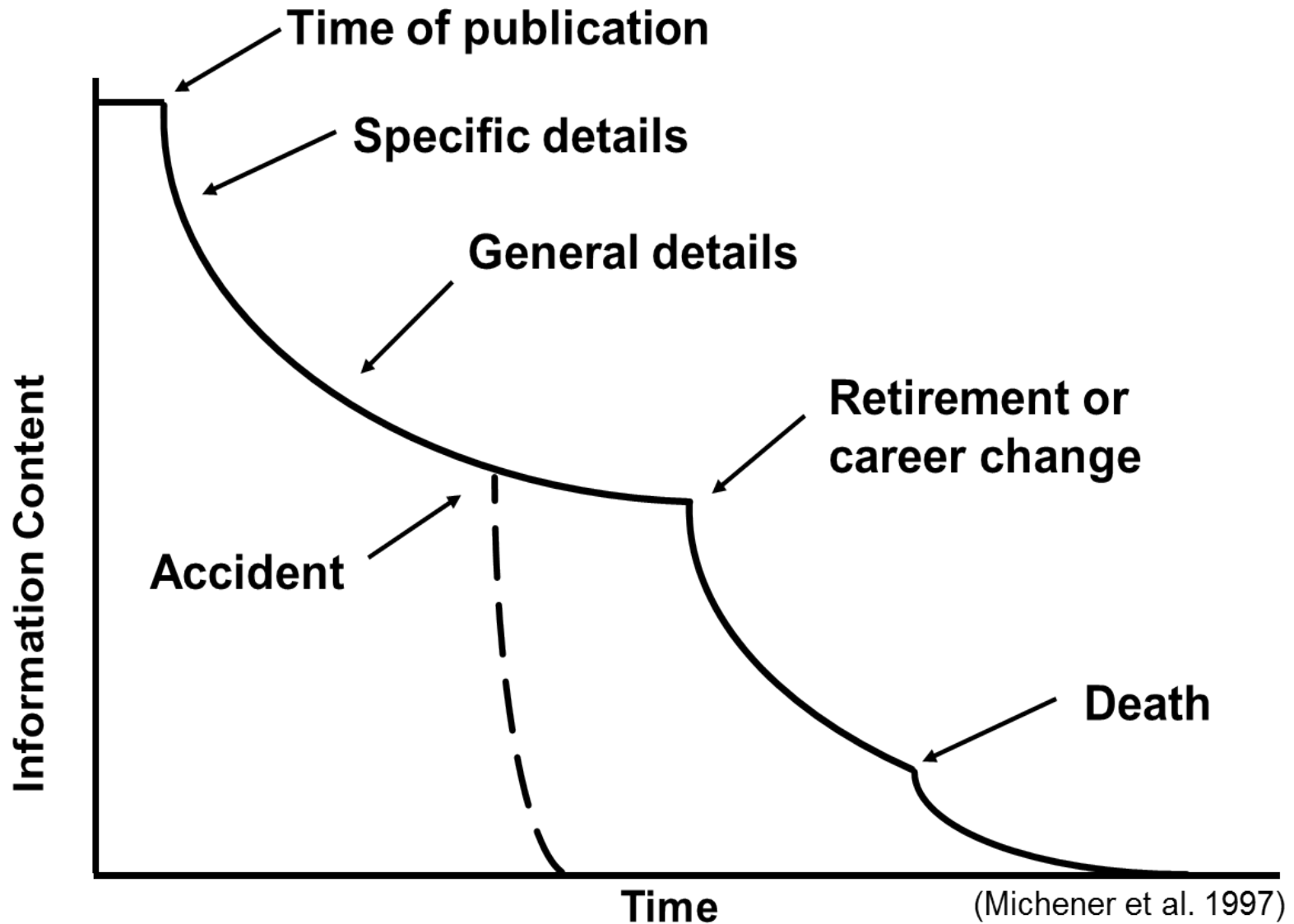
# Data Loss
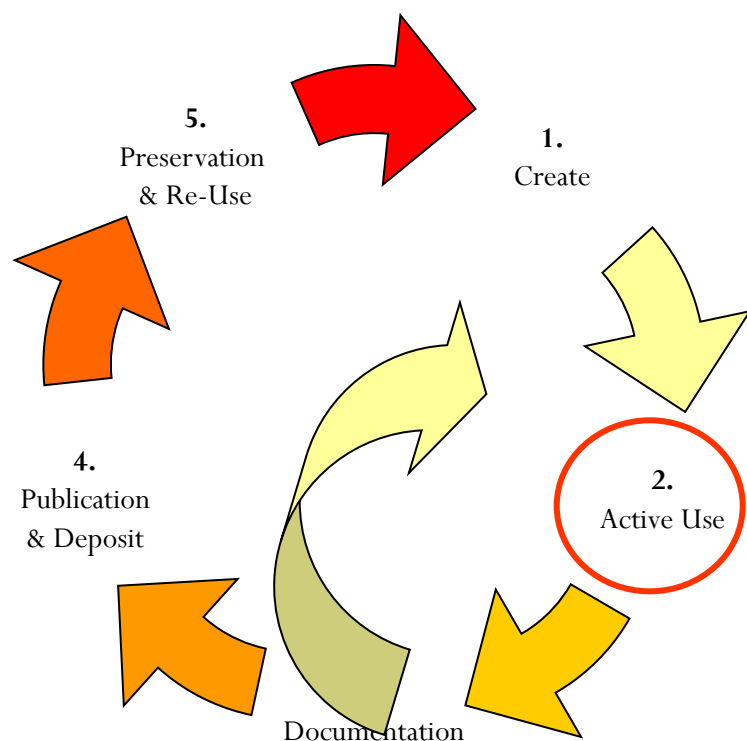


CC image by Sharyn Morrow on Flickr



CC image by momboleum on Flickr

- Natural disaster
- Facilities infrastructure failure
- Storage failure
- Server hardware/software failure
- Application software failure
- Format obsolescence
- Legal encumbrance
- Human error
- Malicious attack by human or automated agents
- Loss of staffing competencies
- Loss of institutional commitment
- Loss of financial stability
- Changes in user expectations and requirements

# Loss does need not be catastrophic, it can be 'entropic'



(Michener et al. 1997)

# How will you look after your data?

**5.**
Preservation
& Re-Use

**1.**
Create

**2.**
Active Use

**4.**
Publication
& Deposit

Documentation

- Is your data safe?

- Is your data organised?

- Can you find your data?

library.uoregon.edu/datamanagement/

# The Benefits of Organizing Files and Folders

- Future-proof your data!

- Meet funder obligations

- Save time

- Easier collaboration; ensure that you work on correct version

- Good research practice

- Archive correct version

# File Naming

**File Naming Conventions**

- Describes essential facts

- Separates files by key data points

- Useful sorting/browsing arrangement

Example:

*20120613_southbend_s23.txt*

[date]_[location]_[instrument].[extension]

# File Naming: Best Practice

- Brief and relevant
- No special characters, dots or spaces
- For separation use underscores _
- Name independent of location
- Date:YYYY_MM_DD

# File Naming

File name = **principal identifier** of file

Easy to: identify, locate, retrieve, access

Provides context e.g.:

✓ version number e.g. *FoodInterview_v1*

   *FoodInterview_v2*

✓ date e.g. *HealthTest_2011_04_06*

✓ content description e.g. *BGHSurveyProcedures*

✓ creator name e.g. *CommsPlanHLJ*

# File Naming: Have a System!

- **Consistent** and logical naming system

- Develop a system with colleagues for shared data

- Be selective about what you save!

- Prepare it for your future self

# File Naming Exercise

Which of the following filenaming conventions would you use, and why?
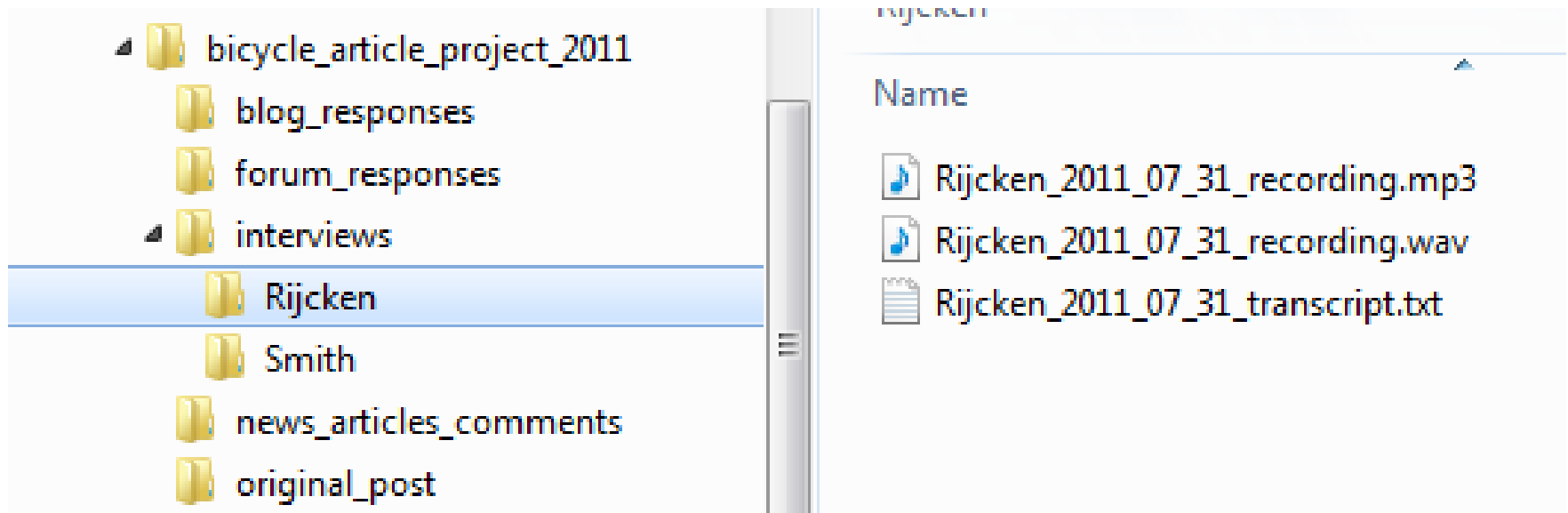
# File naming strategies – examples

- In retrospect I am not very happy with the method I used for naming files. The biggest problem was with the newspaper articles I downloaded… I named the files only based on the topic of the article, without mentioning the name of the periodical and the year of publication, which would have been very useful later, when I began writing the thesis.

    – Doctoral student researching communication history

# Folder Organization

## Folder Arrangement Conventions

- Organize project content into discrete units
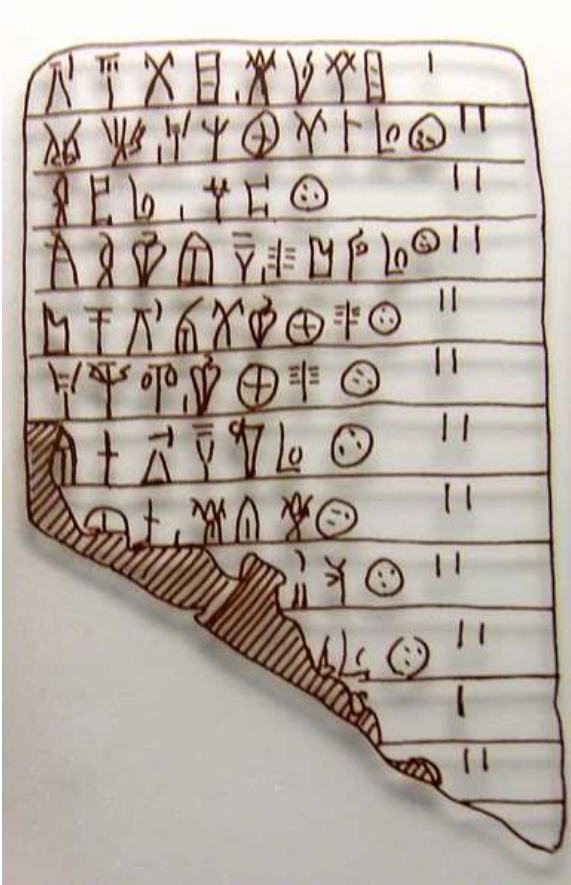- Standardize subfolder arrangement

# What Format Should I Store my Data in? (Open Formats)

- What is an Open Format?
  - The specification has been published
  - The format can be used and implemented by anyone
- Advantages of Open Formats
  - Not limited to one piece of software
  - More chance of being able to use the format in the future
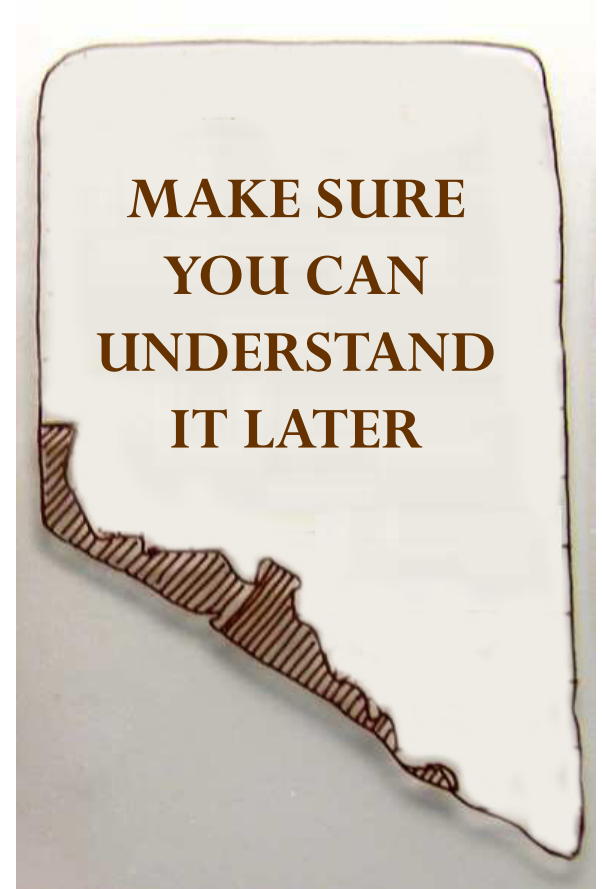- [Examples from UO Data Management pages](#)

# In small groups...

- What data management challenges have you encountered?

- What strategies have you personally found useful?

- Be ready to feed back to the group

- Use the "Nine simple ways…" article for ideas.

# Make material understandable

What's obvious now might not be in a few months, years, decades…

**MAKE SURE YOU CAN UNDERSTAND IT LATER**

Slide adapted from the
PrePARe Project

# Make material verifiable



- Detailing your methods helps people understand what you did

- And helps make your work reproducible

- Conclusions can be verified

# Understanding your data

- Will you be able to write up your methods at the end of your studies?

    - Will you be able to respond to reviewers' comments?

    - Will you be able to find the information you need for final project reports?

- Can you reproduce your work if you need to?

- What information would someone else need to replicate your work?

# Understanding your data

- Do you know how you generated your data?

  - Equipment or software used
  - Experimental protocol
  - Other things included in (e.g.) a lab notebook
  - Can reference a published article, if it covers everything

- Are you able to give credit to external sources of data?

  - Include details of where the data are held, identified & accessed
  - Cite a publication describing the data
  - Cite the data itself e.g.

Gary King; Langche Zeng, 2006, "Replication Data Set for 'When Can History be Our Guide? The Pitfalls of Counterfactual Inference'" hdl:1902.1/DXRXCFAWPK UNF:3:DaYlT6QSX9r0D50ye+tXpA== Murray Research Archive [distributor]

# What is Metadata?

**Metadata is: Data 'reporting'**

- **WHO** created the data?
- **WHAT** is the content of the data?
- **WHEN** were the data created?
- **WHERE** is it geographically?
- **HOW** were the data developed?
- **WHY** were the data developed?

DataONE

# Metadata Exercise

**Look at <u>one</u> of these data records with your group:**

Ecology: http://bit.ly/1ToAb9d

Psychology: http://openpsychologydata.metajnl.com/articles/10.5334/jopd.e2/

**Please answer the following questions:**

a. Is there a license or data use requirement for you to use this dataset?

b. Are the data collection and processing methods described?

c. How would you cite this dataset?

# Backups vs. Preservation

What's the difference?

- Backup $=$ Workflow Recovery

- Preservation $=$ Access (long term)

- **Backups**
- periodic snapshots in case current work is destroyed or lost

# **Backups vs. Preservation**

- What's the difference?

  - **Backups**
  - periodic snapshots in case current work is destroyed or lost

  - **Preservation**
  - long-term retention for future discovery and use

# Why Share Data?

- Share with your future self – avoid repeating research!

- To fulfill funder requirements. For instance, NIH has outlined [key elements to consider in preparing a data sharing plan](#).

- Some journals and societies require data archiving, i.e., the [Dryad partner journals](#) have a [Joint Data Archiving Policy](#) (JDAP).

- To raise interest in your research. [Sharing detailed research data is associated with increased citation rate](#).

- To speed research, particularly in complex research fields

- To establish priority and a public record.

# How and where do you share* data?

- Data repositories – more next week!



**\* and find data, too.**

# Take-aways

- Manage data well throughout the lifecycle

- Good practices at the start of research make it easier to follow through at the end of research

- Data sharing improves science, increases citation of your work, and is emphasized by research funders

# What steps will you take?

Share with your future self!

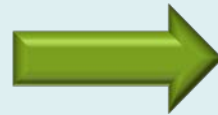| Will you be able to write your thesis/ final report in 3 years time? | → | Use better file names and capture contextual information |

# What steps will you take?

Share with your future self!

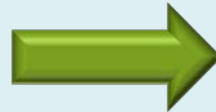What would happen if your computer failed? → Make use of centrally managed research data storage

# What steps will you take?

Share with your future self!

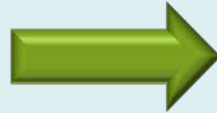| Your digital data is safe but what about your lab/research notebook? | → | Daily scans of your notebook to your personal storage area |

# What steps will you take?

Share with your future self!

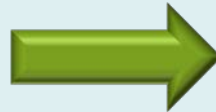| Can you or anyone else understand your older research data? | → | Capture methods and abbreviations in metadata or 'readme' files |

# What steps will you take?

Share with your future self!

| Does your funder require you to make your data accessible? | → | Publish your data in a data repository or data journal |

# References and Sources

- UO Data Management Web site: https://library.uoregon.edu/datamanagement/

- King, Gary. 2011. "Ensuring the Data-Rich Future of the Social Sciences." *Science (New York, N.Y.)* 331 (6018) (February 11): 719–21. http://apps.webofknowledge.com/full_record.do?product=WOS&search_mode=GeneralSearch&qid=9&SID=3DMGm2g4b4nJISv4h6x&page=1&doc=1.

- Piwowar, Heather A, Roger S Day, and Douglas B Fridsma. 2007. "Sharing Detailed Research Data Is Associated with Increased Citation Rate." *PloS One* 2 (3): e308. http://www.plosone.org/article/fetchObject.action?uri=info:doi/10.1371/journal.pone.0000308&representation=PDF.

Some slides and materials in this presentation are from:

- Managing Your Research Data, Catherine Pink and Jez Cope, University of Bath, UK

- DaMaRO Project Introduction to Research Data Management - PowerPoint presentation with presenter's notes - June 2013

# Questions?