# Research Project Workshop 02: Research Data Management (RDM) Overview

March15, 2017
Instructor: Letisha Wyatt, PhD
**tinyurl.com/z8qn8o4**
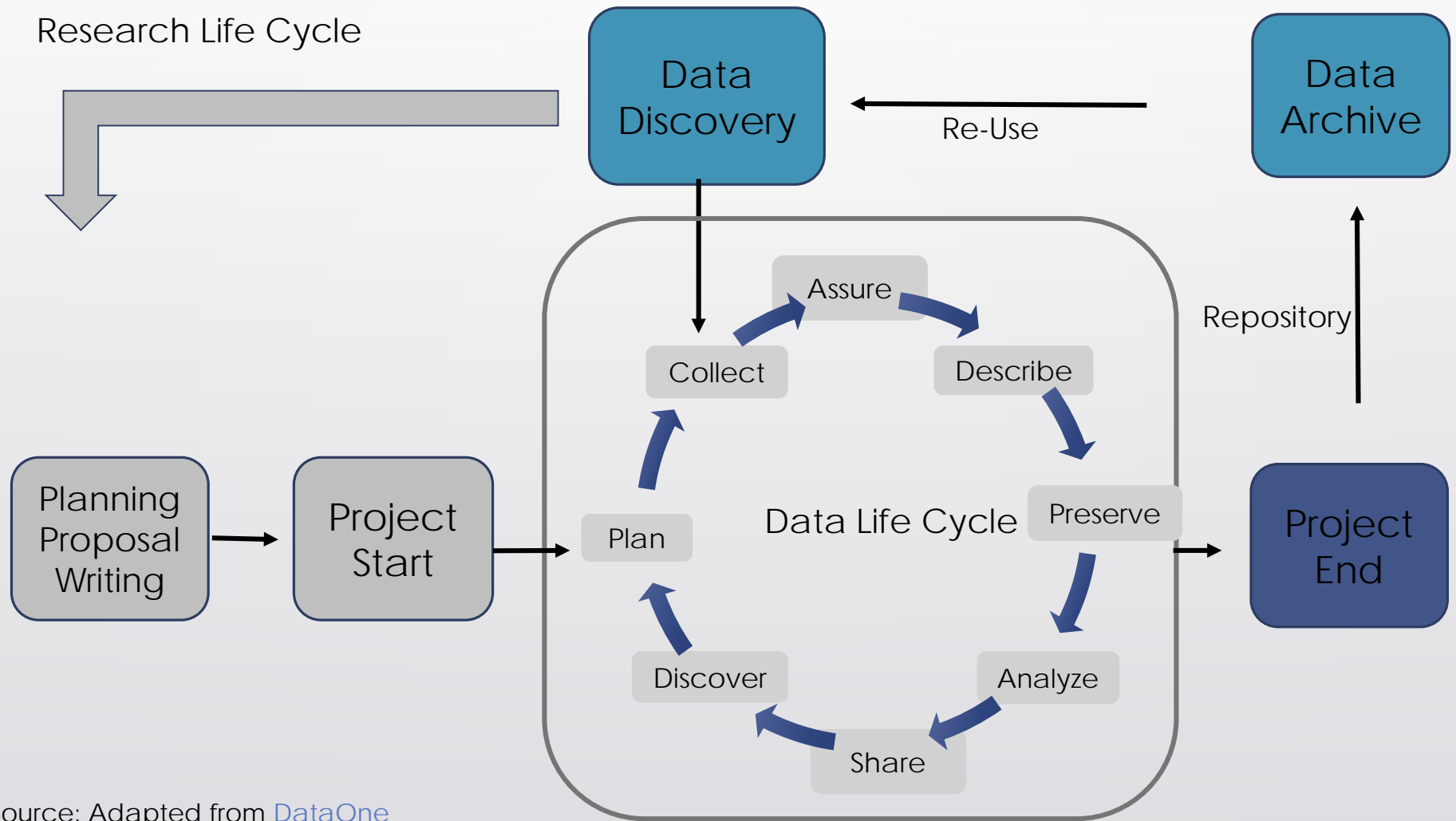
# agenda & learning outcomes

1 Understand the elements of the data life cycle

2 Describe fundamental processes of research data management (RDM)

3 Identify the components of a data management plan and construct a draft

# why is data tlc important?
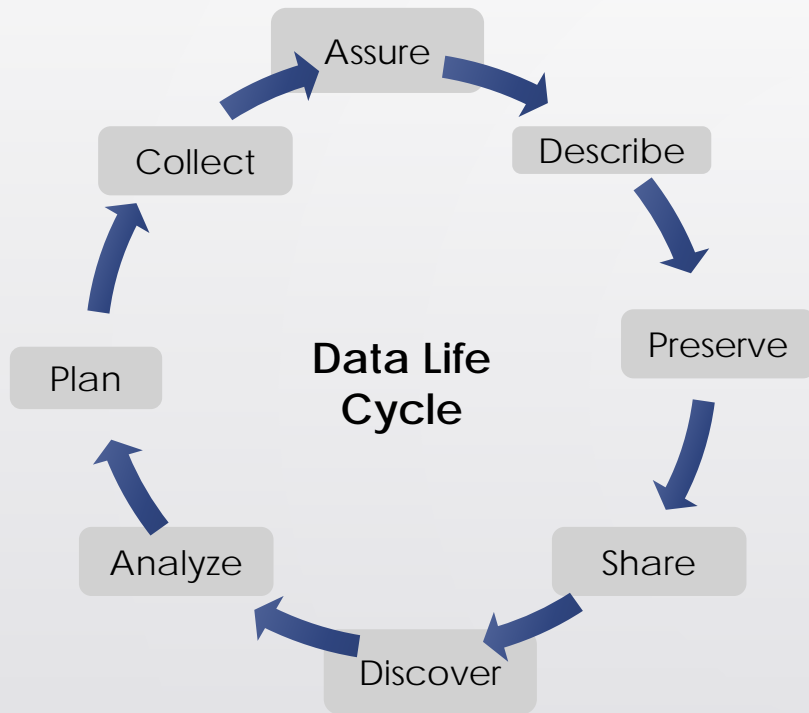
# research/data life cycles

Research Life Cycle

Data Discovery

Data Archive

Re-Use

Repository

Planning Proposal Writing → Project Start →

**Data Life Cycle**

- Assure
- Collect
- Describe
- Plan
- Preserve
- Discover
- Analyze
- Share

→ Project End

how about your
data-related activities?

before the data

# data management plan (dmp)

**Data Life Cycle**

- Assure
- Describe
- Preserve
- Share
- Discover
- Analyze
- Plan
- Collect

☐ Collection: ?
☐ Assurance: ?
☐ Description: ?
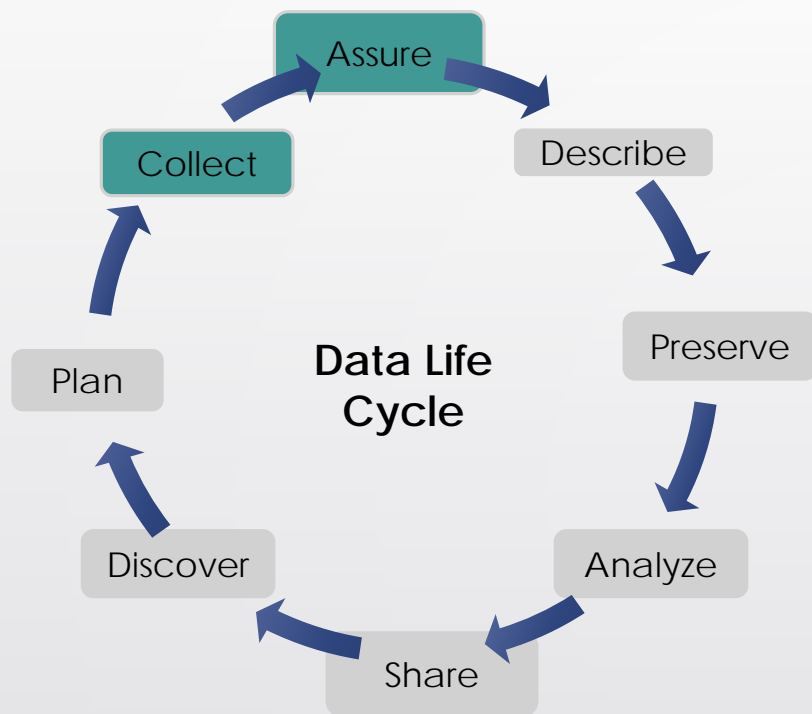☐ Preservation: ?
☐ Sharing/Discovery: ?
☐ Analysis: ?

**postdocs** when surveyed asking "to what extent have you **dealt with** NIH data sharing regulations or NSF **data management plans**?"



Legend: ■ Nationally ■ NYULMC

| Category | Nationally | NYULMC |
|---|---|---|
| Not aware of policies | 38% | 39% |
| Aware but no involvement | 48% | 48% |
| Had to write data plan | 12% | 11% |
| Had to Implement data plan | 8% | 10% |

Source: NYU Medical Library

"Every minute you spend in planning saves 10 minutes in execution; this gives you a 1,000 percent return on energy!"

(Brian Tracy)

data collection

**Data Life Cycle**

Cycle stages: Assure → Describe → Preserve → Analyze → Share → Discover → Plan → Collect → (Assure)

☐ Collection:
    Where is the raw data?
    Format and volume?
    Methods?
    Variables and units?

☐ Assurance:
    Validation or standards
    (instruments/data)
    Missing values?

**Source:** Digital Curation Centre

# data standards in practice

Data standards are the rules by which data are described and recorded. In order to share, exchange, and understand data, we must standardize the format as well as the meaning.

# different definitions for 'data standards'

"Standards are documented agreements containing technical specifications or other precise criteria to be used consistently as rules, guidelines, or definitions of characteristics to ensure that materials, products, processes, and services are fit for their purpose."

(International Organization for Standardization (ISO))

# ways to standardize I

❑ Data content standards

  ❑ Data content standards are a standardized, pre-defined terminology that is used consistently throughout a dataset.

  ❑ For example, when describing a variable, use the same terminology throughout, such as male/female for gender (as opposed to M/F, men/women, etc.)

# ways to standardize II

❑ Data dictionary

   ❑ The purpose of a data dictionary is to ensure consistent terminology throughout a dataset. It contains a list of descriptive data types that are being collected.

   ❑ Example:

| Data Dictionary - Canadian Men's Risk | | | | |
|---|---|---|---|---|
| **Name** | **Definition** | **Data Type** | **Format** | **Values** |
| Counter | identification number for each respondent | identifier | integer | [1, 390] |
| Diagnosis | | categorical | string | Cardiovascular Diabetes type 2 Prostate Cancer Osteoporosis Erectile Dysfunction Low Testosterone None |
| Sex | Sex of respondent | categorical | string | male |
| Age | Age range of respondent [Q1] | categorical | string | <45 45-64 >=65 |

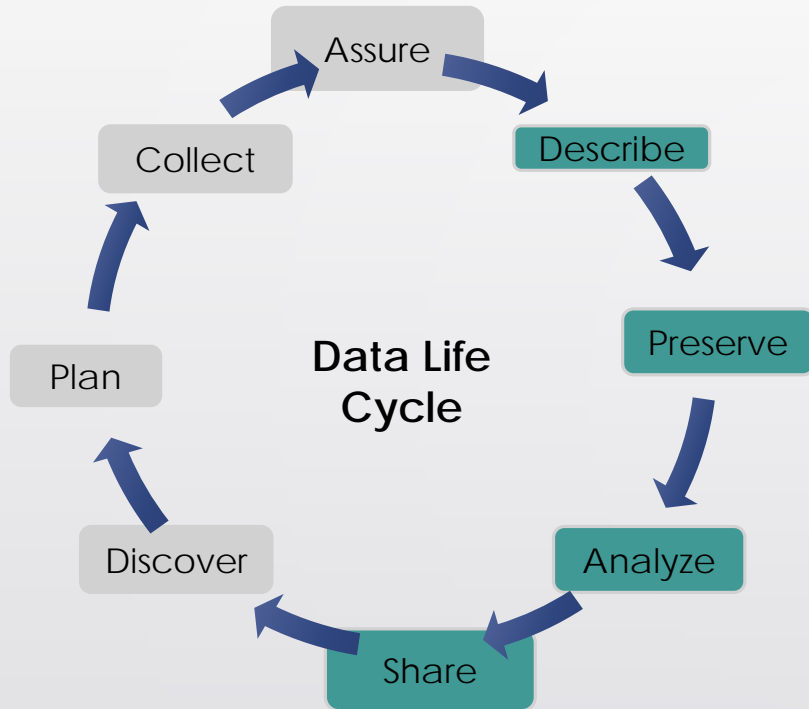Source: OHSU Open Educational Resources

# ways to standardize III

❑ Data modeling

    ❑ Data modeling is the a process that is typically performed before data collection. It involves identifying the user needs from the data set and determining the relationships among the data, and how the data should be best captured.

# ways to standardize IV

❑ Ontology

    ❑ A controlled vocabulary that defines the terms in a given domain or knowledgebase, and the relationships between those terms.

    ❑ *Example*:

       ❑ Gene Ontology

       ❑ Mammalian Phenotype Ontology

# ways to standardize V

❑ Metadata

> ❑ Metadata is often referred to as "data about data". Metadata is intended to describe other data, in order to clarify the meaning of a dataset. Metadata consists of both indexing terms (i.e., terms that can be used for search and retrieval of the data) and descriptive attributes.

> ❑ *Example*:

| Data Dictionary - Canadian Men's Risk | | | | |
|---|---|---|---|---|
| Name | Definition | Data Type | Format | Values |
| Counter | identification number for each respondent | identifier | integer | [1, 390] |
| Diagnosis | | categorical | string | Cardiovascular<br>Diabetes type 2<br>Prostate Cancer<br>Osteoporosis<br>Erectile Dysfunction<br>Low Testosterone<br>None |
| Sex | Sex of respondent | categorical | string | male |
| Age | Age range of respondent [Q1] | categorical | string | <45<br>45-64<br>>=65 |

Source:

after the data

# Information entropy

Time of data development

Specific details about problems with individual items or specific dates are lost relatively rapidly

General details about datasets are lost through time

Retirement or career change makes access to "mental storage" difficult or unlikely

Accident or technology change may make data unusable

Loss of data developer leads to loss of remaining information

DATA DETAILS

TIME

From Michener et al 1997

"Your closest collaborator is you six months ago, but you don't reply to emails…"

# data management plan (dmp)



Data Life Cycle

Assure → Describe → Preserve → Analyze → Share → Discover → Plan → Collect

- ❑ Description:
  - ❑ Descriptive details about the project & experiment

- ❑ Preservation:
  Back up schedule?
  3 Back up locations?
  How long to save it?

- ❑ Sharing:
  Public repository?
  Sensitive data?

- ❑ Analysis:
  Tools (software)?

# how do you describe your data?

Why are you doing this study?
Who is involved?
Where was the data collected?
What were the data collection processes
What do the values in the table mean?
What software is needed to view the data?
How should the data be cited?

README.txt
README.md

# stable file formats

| this | not that |
|------|----------|
| .csv | .xls |
| .txt .html .rtf | .doc |
| .tiff .png | .jpg .gif |

# how will you store & backup the data?

**Computer/Laptop**

**Network Drive**

**External Devices**

**Remote/Cloud**

**Physical Storage (e.g. notebooks)**

# how will you store & backup the data?

| Storage option | The good | The bad |
|---|---|---|
| Computer/Laptop | Convenient for active data | Easily lost/stolen<br>Fail<br>Manual backup |
| Network | Automatic backup and security | Access<br>Capacity limitations |
| External devices | Low cost<br>Portable<br>Easy to use | Easily lost/stolen<br>Fail |
| Remote/Cloud | Global access<br>Collaboration | Security/privacy limitations |
| Physical storage (e.g. notebooks) | Convenient<br>Tangible | Manual backup |

# how long will you keep the data?

# preservation & sharing



Data Availability

**The following policy applies to all PLOS journals, unless otherwise noted.**

PLOS journals require authors to make all data underlying the findings described in their manuscript fully available without restriction, with rare exception.

When submitting a manuscript online, authors must provide a *Data Availability Statement* describing compliance with PLOS's policy. If the article is accepted for publication, the data availability statement will be published as part of the final article.

Refusal to share data and related metadata and methods in accordance with this policy will be grounds for rejection. PLOS journal editors encourage researchers to contact them if they encounter difficulties in obtaining data from articles published in PLOS journals. If restrictions on access to data come to light after publication, we reserve the right to post a correction, to contact the authors' institutions and funders, or in extreme cases to retract the publication.

Methods acceptable to PLOS journals with respect to data sharing are listed below, accompanied by guidance for authors as to what must be indicated in their data availability statement and how to follow best practices in reporting. If authors did not collect data themselves but used another source, this source must be credited as appropriate. Authors who have questions or difficulties with the policy, or readers who have difficulty accessing data, are encouraged to contact the relevant journal office or data@plos.org.

The data policy was implemented on March 3, 2014. Any paper submitted before that date will not have a data availability statement. However for all manuscripts submitted or published before this date, data must be available upon reasonable request.

Download the full text of the older policy (PDF).

*Navigation sidebar:*
- Acceptable Data-Sharing Methods
- Unacceptable Data Access Restrictions
- Explanatory Notes and Guidance
- Recommended Repositories
- FAQs for Data Policy

Source: PLOS Author Guidelines

# 'FAIR' data principles

❑ What is 'FAIR' data?

**F**INDABLE: persistent identifier, registered, rich metadata

**A**CCESSIBLE: retrievable by ID, metadata accessible

**I**NTEROPERABLE: uses shared language, vocab is 'FAIR', references other data

**R**E-USABLE: accurate attributes, released with data use license, standards are met, provenance available

❑ Don't forget about 'FAIR' metadata!

❑ Keep in mind for how you might craft a better plan for your managing your data…

# data plans - why should we care?

# incentives

Maximize your research dollars by sharing and re-using data

Maximize your research impact by sharing data

Treat yo'self: general organization

Source:

# funding agency mandates

❑NIH

    ❑Data sharing plan for new proposals

    ❑Specific institutes and Centers have specific data sharing requirements

❑NSF

    ❑Requires a maximum 2 page data plan for data management and sharing of the products of research for all proposals

Source: OHSU Open Educational Resources

# funding agency mandates

☐NIH/Genomic data

# journal reporting requirements

Reporting requirements vary

Methods section is meant to contain enough information so someone could independently replicate your experiment

Some journals are starting to increase their standards

# hang your checklist!

1. A data management plan has been created ☐

2. Data and files are organised logically ☐

3. Consistent wording and labels have been used ☐

4. Object and methodological metadata have been recorded ☐

5. Files are in a durable, common format, where possible (e.g. CSV, not XLS) ☐

6. Data is stored in a safe location ☐

7. Back-ups have been made ☐

8. A retention period has been applied ☐

9. The owner of the data is clearly identified ☐

10. Ethics guidelines have been complied with ☐

11. A collection-level metadata record has been published to a relevant discovery portal ☐

12. An appropriate licence has been applied ☐

13. A DOI has been minted to allow citation of the dataset ☐

# dmp tool

your turn…

# recap

1 Understand the elements of the data life cycle

2 Describe fundamental processes of research data management (RDM)

3 Identify the components of a data management plan and construct a draft

# key resources

❑ OHSU Big Data to Knowledge (BD2k) [Open Educational Resources (OERs)](#)

❑ OHSU Library [Data Webpage](#)

❑ OHSU RDM Librarian (1:1 consulting)

❑ OHSU Library Workshops: [Contact Us](#)

# questions?

Letisha R. Wyatt, PhD

wyattl@ohsu.edu
BICC 341
503.494.8627