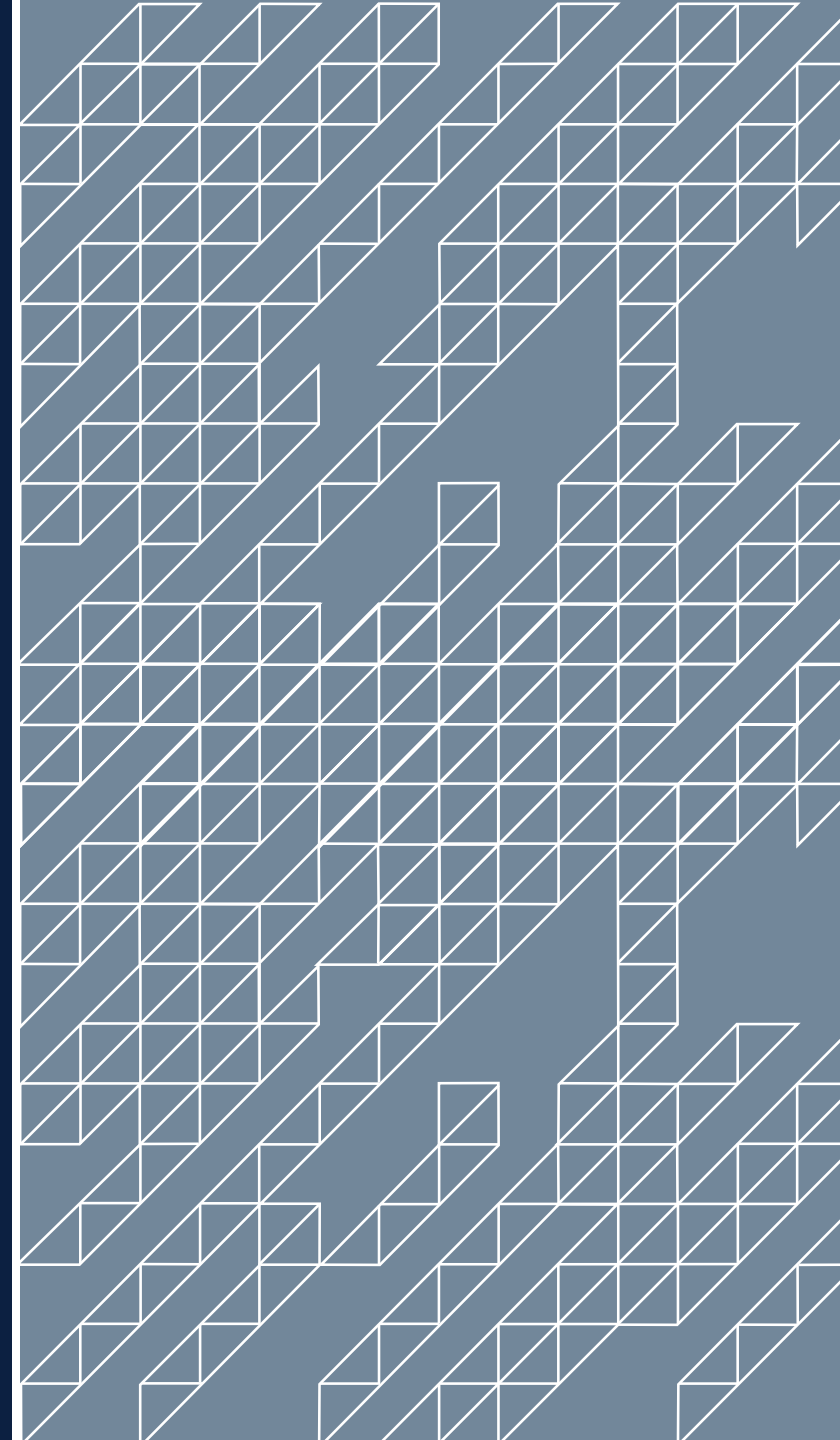


# Preparing Data for Reuse

---

Practical Approaches to Enhancing Data Reusability

Briana Wham,  
Data Learning Center Manager & Research Data Librarian



# Learning Objectives

- Understand why preparing data for reuse is essential
- Explore data repositories
- Learn key components of data documentation
- Investigate and determine documentation (README and data dictionary) needs

# Data Should Be FAIR

## Findable

Persistent Identifiers (PIDs)

iD

Rich metadata



Indexed data repositories



PIDs in metadata



## Accessible

Standard communications protocol



Open, free protocol



Authentication, where necessary

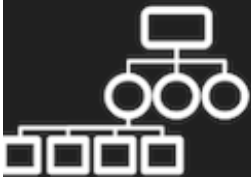


Metadata is always available



## Interoperable

Vocabularies



Vocabularies are FAIR



Linked metadata



## Reusable

Metadata have multiple attributes



Usage license



Provenance



Community standards

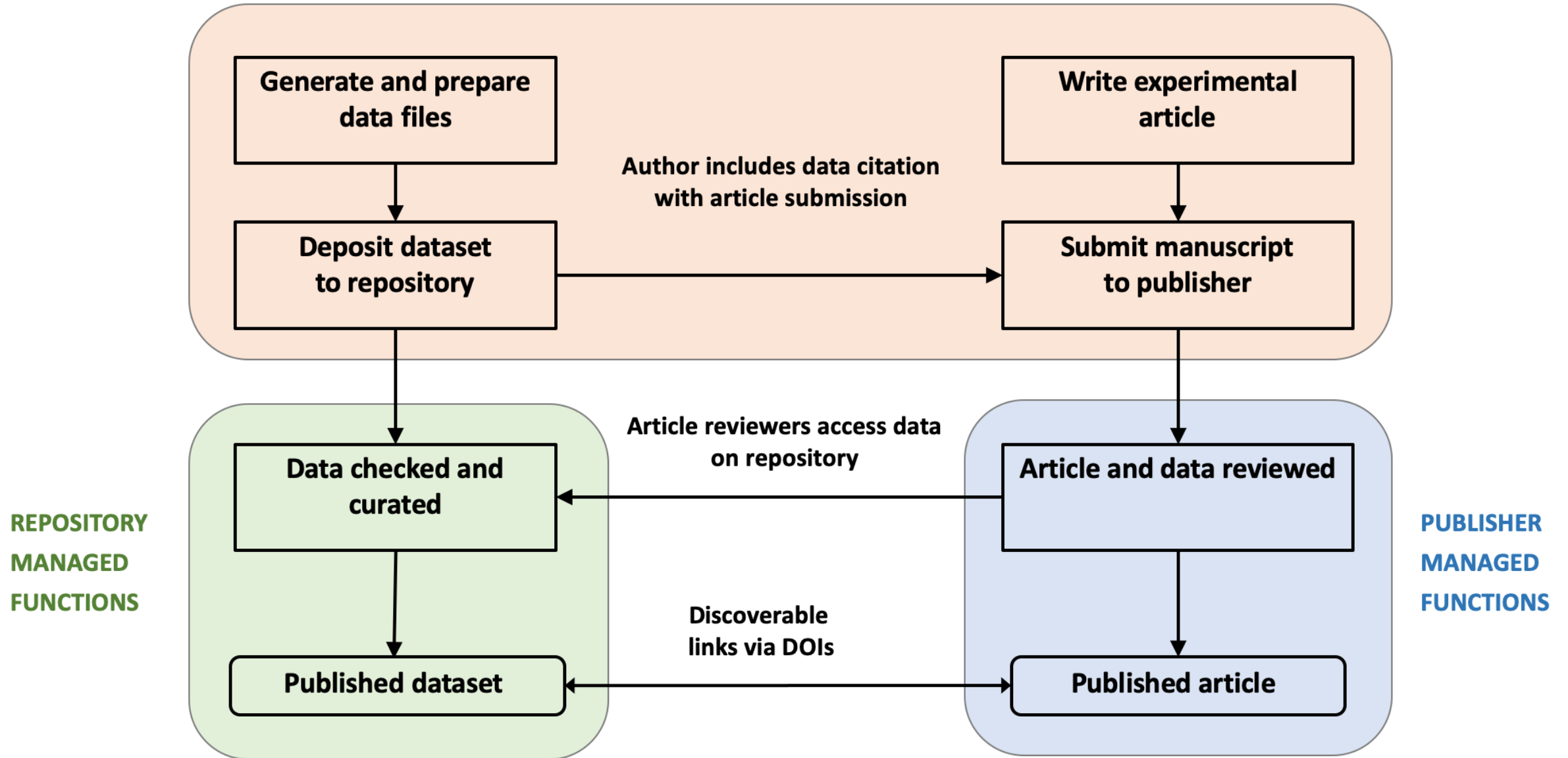


**What issues have you faced reusing data?**

## DATA FLOW

## AUTHOR MANAGED FUNCTIONS

## ARTICLE FLOW



# FAIR Checklist for Dataset/Files

- Is the dataset in a trusted data repository?
- Does the dataset have a registered DOI?
- Are data files in standard and/or commonly available open formats (as much as possible)?
- Are the data and/or metadata retrievable via an API and/or discoverable through an open search protocol?

# What do repositories offer?



**Preservation**

**Access**

**Curation**

# Types of Data Repositories

## Disciplinary

Set up to accommodate the data needs of a specific research community

## Institutional

Support and promote the research outputs of an institution and set up to accept various data types, formats, and disciplinary focuses.

## Generalist

Accept data regardless of type, format, content, or disciplinary focus.



# Types of Data Repositories

## Disciplinary

Qualitative Data Repository

ICPSR

Vivli  
CENTER FOR GLOBAL CLINICAL RESEARCH DATA



Databrary



NIMH Data Archive

NDA | ABCD | CCF | NIAAA<sub>DA</sub> | OAI

## Institutional

ScholarSphere



datacommons@psu  
data: energy environment society technology

## Generalist



zenodo



# Types of Data Repositories

<b>Public-Public (Open access)</b>	<b>Public-Private (Mediated open access)</b>	<b>Private-Private (Closed access)</b>
Metadata is fully discoverable	Metadata is fully discoverable	Metadata is not publicly available
Data are accessible and immediately downloadable	Mediated access to data via data custodian	Data not discoverable or available to third parties
Preferred option for non-sensitive data from completed projects	Good option for sensitive or confidential data	Safest option for highly-sensitive data

Table originally from Curtin University, Research Data Team. <http://libguides.library.curtin.edu.au>

# Desirable Characteristics of Data Repositories

## Desirable Characteristics of Data Repositories for Federally Funded Research

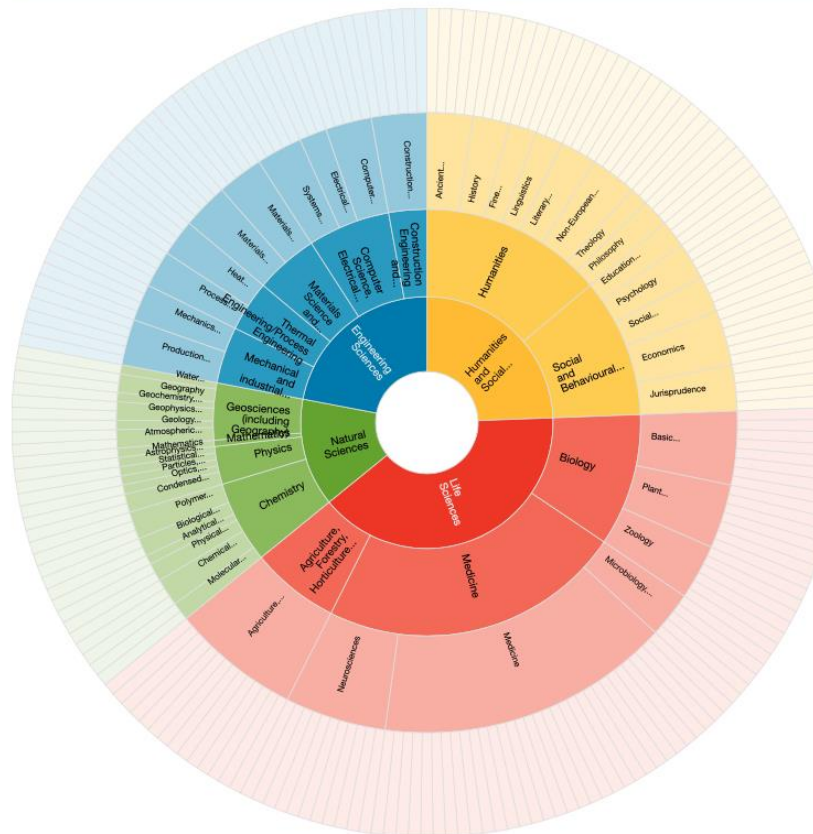
Unique Persistent Identifiers	Long-Term Sustainability	Metadata	Curation and Quality Assurance
Free and Easy Access	Broad and Measured Reuse	Clear Use Guidance	Security and Integrity
Risk Management/Confidentiality	Common Format	Provenance	Retention Policy

For more see [NOT-OD-21-016](#) – Supplemental Information to the NIH Policy for Data Management and Sharing: Selecting a Repository for Data Resulting from NIH-Supported Research

## Browse by subject

Graphical **Text**

click to zoom into subjects or to select a bottommost subject in the hierarchy as filter for the re3data search page  
ctrl + click on a top subject to select it as filter



## Filter

- Subjects** ⊕
- Content Types** ⊕
- Countries** ⊕
- AID systems** ⊕
- API** ⊕
- Certificates** ⊕
- Data access** ⊕
- Data access restrictions** ⊕
- Database access** ⊕
- Database access restrictions** ⊕
- Database licenses** ⊕
- Data licenses** ⊕
- Data upload** ⊕
- Data upload restrictions** ⊕
- Enhanced publication** ⊕
- Institution responsibility type** ⊕
- Institution type** ⊕
- Keywords** ⊕
- Metadata standards** ⊕
- PID systems** ⊕
- Provider types** ⊕
- Quality management** ⊕
- Repository languages** ⊕
- Software** ⊕
- Syndications** ⊕
- Repository types** ⊕
- Versioning** ⊕

# Sustainable File Formats For Long-Term Storage & Preservation

TYPE OF DATA	RECOMMENDED FILE FORMATS FOR SHARING, RE-USE AND PRESERVATION
<b>Quantitative tabular data with extensive metadata</b> 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
<b>Quantitative tabular data with minimal metadata</b> a matrix of data with or without column headings or variable names, but no other metadata or labelling	comma-separated values (CSV) file (.csv)  tab-delimited file (.tab)  including delimited text of given character set with SQL data definition statements where appropriate
<b>Geospatial data</b> vector and raster data	ESRI Shapefile (essential: .shp, .shx, .dbf ; optional: .prj, .sbx, .sbn)  geo-referenced TIFF (.tif, .tiff)  CAD data (.dwg)  tabular GIS attribute data
<b>Qualitative data</b> textual	eXtensible Mark-up Language (XML) text according to an appropriate Document Type Definition (DTD) or schema (.xml)  Rich Text Format (.rtf)  plain text data, ASCII (.txt)
<b>Digital image data</b>	TIFF version 6 uncompressed (.tif)
<b>Digital audio data</b>	Free Lossless Audio Codec (FLAC) (.flac)
<b>Digital video data</b>	MPEG-4 (.mp4)  motion JPEG 2000 (.jp2)
<b>Documentation</b>	Rich Text Format (.rtf)  PDF/A or PDF (.pdf)  OpenDocument Text (.odt)

# FAIR Checklist for Data Documentation

- Are all associated data files unambiguously named and described including file types, software requirements, and/or conversion information in the metadata?
- Does the metadata include useful disciplinary notation and terminology?
- Does the metadata include machine-readable standards where available?
- Are related articles referenced and linked in the metadata?
- Is a citation format for the dataset provided?
- Are any license terms, attribution, or terms of use clearly indicated?
- Is the metadata exportable in a machine-readable structured text-based format?

# Data Documentation

- the context of data collection: project history, aim, objectives and hypotheses
- data collection methods: sampling, data collection process, instruments used, hardware and software used, scale and resolution, temporal and geographic coverage and secondary data sources used
- dataset structure of data files, study cases, relationships between files
- data validation, checking, proofing, cleaning and quality assurance procedures carried out
- changes made to data over time since their original creation and identification of different versions of data files
- information on access and use conditions or data confidentiality

# Documentation Formats

## Metadata

Highly-structured data laid out in fields, often with controlled vocabularies in each field.

## Data Dictionaries

Defines variables without cluttering datasets.

## README.txt files

Simple text files that provide support for navigating folders and understanding the contents of your files.



# Metadata

## Minimal

Enable basic discovery and access

- Creator
- Title
- Publisher
- Date
- Embargo term
- License
- Access terms and conditions

# Metadata

## Minimal

Enable basic discovery and access

- Creator
- Title
- Publisher
- Date
- Embargo term
- License
- Access terms and conditions

## General Contextual

May be of an administrative nature or relate to project methodologies

- Funder
- Grant #
- Project information
- Data generation process
- Geographical location of data collection
- Date range of data collection

# Metadata

## Minimal

Enable basic discovery and access

- Creator
- Title
- Publisher
- Date
- Embargo term
- License
- Access terms and conditions

## General Contextual

May be of an administrative nature or relate to project methodologies

- Funder
- Grant #
- Project information
- Data generation process
- Geographical location of data collection
- Date range of data collection

## Discipline/ Data Specific

Enables re-use

- Machine settings
- Experimental conditions

\*\*Can be difficult to capture in structured metadata fields unless depositing in a discipline specific repository or community standard exists  
-> information can be put in data documentation

# Metadata

<https://fairsharing.org/>

<http://www.dcc.ac.uk/drupal/resources/metadata-standards>

<http://rd-alliance.github.io/metadata-directory/>

## Minimal

Enable basic discovery and access

- Creator
- Title
- Publisher
- Date
- Embargo term
- License
- Access terms and conditions

## General Contextual

May be of an administrative nature or relate to project methodologies

- Funder
- Grant #
- Project information
- Data generation process
- Geographical location of data collection
- Date range of data collection

## Discipline Specific

Enables re-use

- Machine settings
- Experimental conditions

\*Example Metadata schema: Darwin Core (biodiversity taxa data), Ecological Metadata Language, ISO1911 & ISO 19139 (geospatial data), DDI (social science, survey data), MIxS (Genomic data)

\*Example standards: Biological classification system, NCBI Taxonomy  
Cox and Verbaan 2018

# Data Dictionary

## Common information included:

- Variable name
- Variable definition
- How the variable was measured
- Data units
- Data formats
- Minimum and maximum values
- Coded values and their meanings
- Representation of null vs. NA values
- Precision of measurement
- Known issues with the data (missing values, bias, etc.)
- Relationship to other variables
- Other important notes about the data

**Table 1.** List of taxa included in the data files, including taxonomic code used in all data files (Taxon), Latin name and common name of each taxon<sup>4,5</sup>.

count	Taxon	Latin_Name	Common_Name
1	CMED	<i>Cheirogaleus medius</i>	Fat-tailed dwarf lemur
2	DMAD	<i>Daubentonia madagascariensis</i>	Aye-aye
3	EALB	<i>Eulemur albifrons</i>	White-fronted brown lemur
4	ECOL	<i>Eulemur collaris</i>	Collared brown lemur
5	ECOR	<i>Eulemur coronatus</i>	Crowned lemur
6	EFLA	<i>Eulemur flavifrons</i>	Blue-eyed black lemur
7	EFUL	<i>Eulemur fulvus</i>	Common brown lemur
8	EMAC	<i>Eulemur macaco</i>	Black lemur
9	EMON	<i>Eulemur mongoz</i>	Mongoose lemur
10	ERUB	<i>Eulemur rubriventer</i>	Red-bellied lemur
11	ERUF	<i>Eulemur rufus</i>	Red-fronted brown lemur
12	ESAN	<i>Eulemur sanfordi</i>	Sanford's brown lemur
13	EUL	<i>Eulemur</i>	Eulemur hybrid
14	GMOH	<i>Galago moholi</i>	Mohol bushbaby
15	HGG	<i>Hapalemur griseus griseus</i>	Eastern lesser bamboo lemur
16	LCAT	<i>Lemur catta</i>	Ring-tailed lemur
17	LTAR	<i>Loris tardigradus</i>	Slender loris
18	MMUR	<i>Mircocebus murinus</i>	Gray mouse lemur
19	MZAZ	<i>Mirza coquereli</i>	Northern giant mouse lemur
20	NCOU	<i>Nycticebus coucang</i>	Slow loris
21	NPYG	<i>Nycticebus pygmaeus</i>	Pygmy slow loris
22	OGG	<i>Otolemur garnettii garnettii</i>	Northern greater galago
23	PCOQ	<i>Propithecus coquereli</i>	Coquerel's sifaka
24	PPOT	<i>Perodicticus potto</i>	Potto
25	VAR	<i>Varecia</i>	Varecia hybrid
26	VRUB	<i>Varecia rubra</i>	Red ruffed lemur
27	VVV	<i>Varecia variegata variegata</i>	Black-and-white ruffed lemur

**Table 5:** DLC Weight File variable descriptions

count	Weight File Variable Name	Weight File Variable Definition
1	Taxon	Taxonomic code: In most cases, comprised of the first letter of the genus and the first three letters of the species; if taxonomic designation is a subspecies, comprised of the first letter of genus, species, and subspecies, and hybrids are indicated by the first three letters of the genus. See table 1 for details.
2	Hybrid	Hybrid status: N=not a hybrid. S=species hybrid. B=subspecies hybrid. If sire is one of multiple possible and animal could be a hybrid, it is designated a hybrid.
3	DLC_ID	Specimen ID: Unique DLC number assigned at accession of animal
4	Sex	Sex: M=male. F=Female. ND=Not determined

## Data Dictionary

Variable	Description	Data Type	Prescribed Values/ Format	Unit of Measurement	Sample Data
<b>phase</b>	important stages of the experiment when plot size, crop rotation or treatments changed	categorical	1-5	-	3
<b>year</b>	year planted and harvested	date	YYYY	-	1955
<b>plot</b>	label indicating plot number, North/South location, and A/B/C/D subplot (see Data Sources section above for an important note about plot names)	categorical	3NA; 3NB; 3NC; 3ND; 3SA; 3SB; 3SC; 3SD; 4NA; 4NB; 4NC; 4ND; 4SA; 4SB; 4SC; 4SD; 5NA; 5NB; 5NC; 5ND; 5SA; 5SB; 5SC; 5SD	-	4SD
<b>plot_num</b>	plot number only	categorical	3-5	-	4
<b>plot_dir</b>	N/S and E/W location	categorical	NE; NW; SE; SW	-	SE
<b>rotation</b>	number of years in the crop rotation schedule for this plot	categorical	1-3	-	2
<b>corn</b>	flag making it easy to group corn in rotation and continuous corn	T/F	TRUE; FALSE	-	TRUE
<b>crop</b>	crop planted this year with separate values for corn in rotation (C) and continuous corn (CC)	categorical	A [alfalfa]; C [corn]; CC [continuous corn]; H [hay]; O [oats]; S [soybeans]	-	C
<b>variety</b>	crop variety name	text	-	-	Illinois 1570
<b>all_corn</b>	whether this was a year when corn was planted in all plots	T/F	TRUE; FALSE	-	TRUE
<b>yield_bush</b>	yield for all crops except hay	numerical	-	bushels/acre	95.6
<b>yield_ton</b>	yield for hay	numerical	-	tons/acre	3.48
<b>treated</b>	whether this plot was treated that year	T/F	TRUE; FALSE	-	TRUE

# Demo - Creating a Data Dictionary in R

# Writing a README – Best Practices

- Create readme files for logical “clusters” of related files/data
- Name the README so that it is easily associated with the data file(s) it describes
- Write your README document as plain text file, often formatted with markdown (.md)
- Format multiple README files identically
- Use standardized date formats
- Follow the scientific conventions for your discipline for taxonomic, geospatial, and geologic names and keywords



# Writing a README – Recommended Content

## General Information

- **Provide a title for the dataset**
- **Name/institution/address/email information for**
  - **Principal investigator or person responsible for collecting the data**
  - Associate or co-investigators
  - Contact person for questions
- **Date of data collection**
- **Information about geographic location of data collection**
- Keywords
- Language
- Information about funding sources

# Writing a README – Recommended Content

## Data and file overview

- **For each filename, a short description of what data it contains**
- Format of the file if not obvious from the file name
- The relationship between files if multiple files relate to one another
- **Date that the file was created**
- Date(s) that the file(s) was updated (versioned) and the nature of the update(s), if applicable
- Information about related data collection that is not in the described dataset

# Writing a README – Recommended Content

## Sharing and access information

- **Licenses or restrictions placed on the data**
- Links to publications that cite or use the data
- Recommended citation for the data

# Writing a README – Recommended Content

## Methodological information

- **Description of methods for data collection or generation** (include links or references to publications or other documentation containing experimental design or protocols used)
- **Description of methods used for data processing**
- Any software or instrument-specific information needed to understand or interpret the data, including software version numbers
- Describe any quality-assurance procedures performed on the data
- Definitions of codes or symbols
- People involved with sample collection, processing, analysis, and/or submission

# Writing a README – Recommended Content

## Data-specific information

\*Repeat this section as needed for each dataset (or file, as appropriate)

- Count of number of variables and rows
- **Variable list, including full names and definitions of column headings for tabular data**
- **Units of measurement**
- **Definitions for codes or symbols used to record missing data**

# Demo – README Template

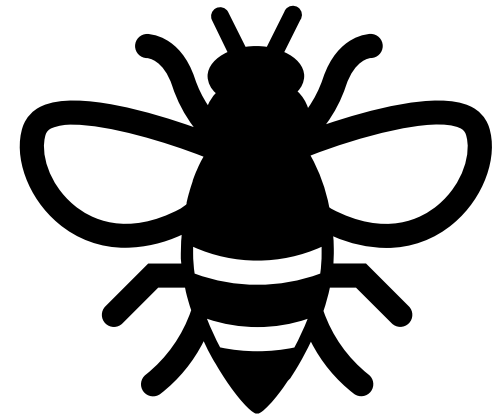
# Activity – Can you Reuse This Data?

**Imagine you're a researcher from another institution. You found this dataset and want to use it. Can you understand it well enough?**

## **Dataset Description:**

This dataset was deposited by an entomologist to their local data repository; it has been slightly simplified for the purposes of this activity.

It contains experimental measurements documenting how honey bee colonies exposed to high temperatures maintain and meet the water needs of the colony to survive.

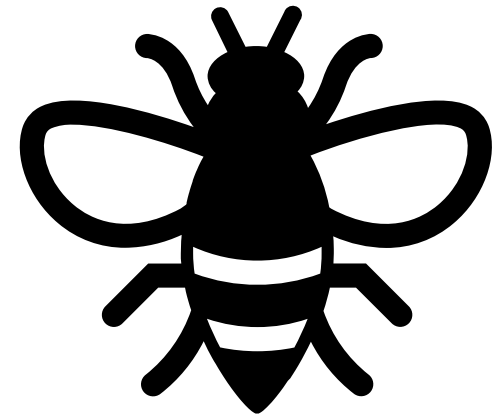


**Access Activity Materials:** <https://tinyurl.com/2e5834as>

# Activity – Can you Reuse This Data?

10 min

1. Review the dataset and the bare-bones README
2. Consider:
  - i. What's missing or unclear?
  - ii. What assumptions did you have to make?
  - iii. What information would help you trust and reuse this data?



Access Activity Materials: <https://tinyurl.com/2e5834as>

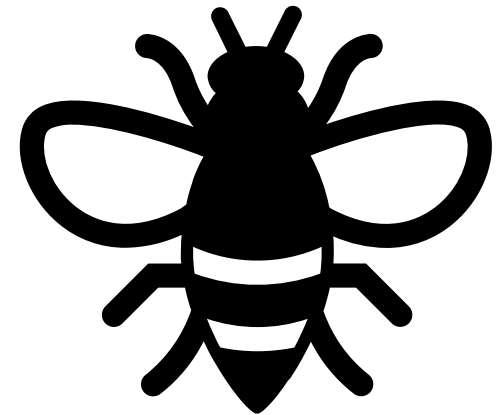


# Activity – Can you Reuse This Data?

Access Updated data & README: <https://tinyurl.com/5fpnkmd2>

10 min

1. What's different?
2. What information was added?
3. How does this improved README support reuse?



# Activity – Finding the Right Home for Your Data

**Imagine you are preparing to share your own dataset. How would you choose a repository and documentation format?**