

The Research Data Management Workbook

Kristin Briney

2023-08-25

Contents

| | |
|--|-----------|
| About this Book | 5 |
| Description | 5 |
| Edition | 5 |
| License | 5 |
| The Author | 6 |
| 1 Introduction | 7 |
| 1.1 What is Research Data Management? | 7 |
| 1.2 Why do Research Data Management? | 7 |
| 1.3 How to Use this Book to Manage Your Research Data Better . . | 8 |
| 1.4 Further Resources on Research Data Management | 8 |
| 2 Documentation | 11 |
| 2.1 Evaluate a Laboratory Notebook | 11 |
| 2.2 Write a Project-Level README.txt | 12 |
| 2.3 Create a Data Dictionary | 14 |
| 3 File Organization and Naming | 15 |
| 3.1 Set Up a File Organization System | 15 |
| 3.2 Create a File Naming Convention | 16 |
| 4 Data Storage | 21 |
| 4.1 Store and Back Up the Data | 21 |
| 4.2 Test Your Backup | 24 |

| | | |
|----------|---|-----------|
| 5 | Data Management | 27 |
| 5.1 | Write a Living Data Management Plan (DMP) | 27 |
| 5.2 | Determine Data Stewardship | 29 |
| 6 | Data Sharing | 31 |
| 6.1 | Pick a Data Repository | 31 |
| 6.2 | Share Data | 33 |
| 7 | Project Wrap Up | 35 |
| 7.1 | Prepare Data for Future Use | 35 |
| 7.2 | Convert Data File Types | 36 |
| 7.3 | Create an Archive Folder | 37 |
| 7.4 | Separate from the Institution | 39 |

About this Book

Description

The Research Data Management Workbook is made up of a collection of exercises for researchers to improve their data management. The Workbook contains exercises across the data lifecycle, though the range of activities is not comprehensive. Instead, exercises focus on discrete practices within data management that are structured and can be reproduced by any researcher.

The book is divided into chapters, roughly by phases of the data lifecycle, with one or more exercises in each chapter. Every exercise comes with a description of its value within data management, instructions on how to do the exercise, original source of the exercise (when applicable), and the exercise itself.

The Workbook is intended as a supplement to existing data management education. If you would like to learn more about the principles of data management, please see the article “Foundational Practices of Research Data Management” [Briney et al., 2020] or read the book “Data Management for Researchers” [Briney, 2015].

Edition

The Research Data Management Workbook is currently in its beta edition, meaning I’m still tweaking the exercises and trying to catch stray typos. I’m more than happy to receive any feedback on the Workbook to improve it; you can message me at briney@caltech.edu.

License

This book is available under a Creative Commons Attribution-NonCommercial (CC BY-NC) 4.0 International license.

I encourage you to use and adapt all of the exercises in this book for educational and personal use. Just remember to cite me:

- Briney, K. A. (2023). *The Research Data Management Workbook*. Caltech Library.

The Author



Figure 1: Headshot of author, Kristin Briney

Kristin Briney is the Biology & Biological Engineering Librarian at the California Institute of Technology and author of the books “Data Management for Researchers” [Briney, 2015] and, with Becky Yoose, “Managing Data for Patron Privacy” [Briney and Yoose, 2022]. She has a PhD in chemistry and an MLIS, both from the University of Wisconsin-Madison. Her research focuses on research data management, institutional data policy, and patron privacy vis-a-vis library data handling. Kristin is an advocate for the adoption of the international date standard ISO 8601 (YYYY-MM-DD) and likes to spend her free time making data visualizations out of yarn and fabric.

Chapter 1

Introduction

1.1 What is Research Data Management?

Research data management is a set of collective practices and decisions that make it easier for you, your collaborators, and your future self to easily find, understand, and use your research data. These practices cover the entire lifecycle of research data, from its collection and analysis through sharing and reuse. There is no one magical data management practice to rule them all. Rather, data management consists of a number of small activities that make it easier to deal with your data. Research is hard enough as it is without having to fight with your files, so the goal of data management is for you to maximize your time doing research instead of wasting time with file handling.

1.2 Why do Research Data Management?

Most researchers have spent time, at some point in their careers, digging through their computer to find a specific file that can't be located. It's incredibly frustrating and a waste of time and resources, especially if you end up recollecting missing data. The good news is that it is possible to avoid this situation entirely by strategically managing your data better.

Done well, research data management means:

- always understanding what your data is and how you collected it even if the data is a year old
- always finding the file you need quickly
- never losing your data even if your hard drive crashes
- knowing what rights and responsibilities you have over your data

- knowing how and where to share your data to comply with your funder’s data sharing policy
- being able to pick up and easily reuse data from a past project

If all of those things sound like something you would like to implement in your research, you are reading the right book!

1.3 How to Use this Book to Manage Your Research Data Better

The “Research Data Management Workbook” is focused entirely on the “how” of data management. The Workbook consists of a series of worksheets, checklists, and procedures to set up new data management practices, check existing practices, and make good decisions about your data. You will find exercises covering all the ideals listed in the “Why Do Research Data Management?” section above, allowing you to streamline the your use of research data.

The Workbook is not a complete set of exercises for everything under the umbrella of data management. Instead, the Workbook centers on activities that are structured, reproducible, and apply to many researchers. The strength and weakness of data management is that many of its practices are customizable to individual research workflows. Each exercise in the Workbook, therefore, is built on best principles while allowing for customization to suit local needs.

You may go through exercises in the Workbook collectively or individually as you chose. Do note that some exercises require completing one or more other exercises in the book, so it’s best to have the whole workbook on hand just in case. Finally, for the exercises that have been formatted as worksheets, I recommend printing them out and writing your answers in the space provided.

1.4 Further Resources on Research Data Management

“The Research Data Management Workbook” does not comprehensively explain data management and therefore works best for those with some foundational data management knowledge or in tandem with other educational resources on research data management.

My best recommendation is to use the Workbook as the exercise book for my first book, “Data Management for Researchers” [Briney, 2015]. The following table lays out how chapters in the Workbook match with chapters in “Data Management for Researchers,” allowing you to look up more information on any topic covered by the Workbook:

1.4. FURTHER RESOURCES ON RESEARCH DATA MANAGEMENT 9

| Research Data Management Workbook | Data Management for Researchers |
|-----------------------------------|---------------------------------|
| Chapter 2 | Chapter 4 |
| Chapter 3 | Chapter 5 |
| Chapter 4 | Chapter 8 |
| Chapter 5 | Chapter 3 |
| Chapter 6 | Chapter 10 |
| Chapter 7 | Chapter 9 |

Chapter 2

Documentation

Documentation has sometimes been called “a love letter to your future self” as it helps you remember important details about your research data. The great thing about research documentation is that it’s not limited to a laboratory or research notebook, though notebooks are still very important! This chapter introduces two types of useful documentation – a project-level README.txt and a data dictionary – and offers worksheets for writing both. The chapter also includes a worksheet to evaluate an older entry in your laboratory notebook to ensure your documentation is of sufficient quality.

2.1 Evaluate a Laboratory Notebook

Description: *The laboratory or research notebook is a fundamental documentation method for many researchers. But for how ubiquitous the lab notebook is, documentation can sometimes be lacking. The ideal laboratory notebook allows someone with similar training as you to be able to follow everything you did in your research. This exercise prompts you to review an old entry within your laboratory notebook to evaluate if your documentation is sufficient for reproducing your work.*

Instructions: *You will need a laboratory notebook entry from 6-12 months ago to do this exercise. Once you have the entry, read through it to try to understand what you did on that day. Answer the exercise questions to evaluate the entry and identify any note keeping improvements to make.*

Date of lab notebook entry being evaluated: _____

Read the entry and summarize the work you did on that date:

How easy was it to understand what you did from your notes? Could you reproduce your work solely from the information in your notes?

What worked well with your note keeping?

What would you improve about your note keeping?

List one change you plan to make to take better research notes:

2.2 Write a Project-Level README.txt

Description: Data files living on a computer often need extra documentation for someone to understand what research they correspond to. In particular, it is useful to record the most basic project information and store it in the top-level folder of each research project. This can be done with a README.txt. The name, “README”, indicates that the file conveys important information and the file type, TXT, can be opened by many different software programs, making the content maximally accessible. This exercise walks you through the key information needed in a project-level README.txt file. The same information can also be recorded at the front of a physical laboratory notebook.

Instructions: Pick a research project and answer the following questions. Copy all of the text into a TXT file and save it with the name “README.txt”. Store

this file in the top-level of the project folder on your computer, alongside the project files.

Source: *This exercise was adapted from the “Project Close Out Checklist” [Briney, 2020b].*

What is the title of the project?

What is the project description?

What is the time period the project was done over?

Who worked on the project?

Where are the files are stored?

How are files organized? Are any naming conventions used and, if so, what are they (see Chapter 3)?

2.3 Create a Data Dictionary

Description: Ideally, a spreadsheet is formatted with a row of variable names at the top, followed by rows of data going down. This makes easy for data to be used in any data analysis software (interoperability is a good thing) but makes it impossible to document a spreadsheet within the file itself. For this reason, it's useful to create a data dictionary to describe the spreadsheet so that others can interpret the data. This exercise walks you through the major information you should record for each variable in the spreadsheet, adding up to a complete dictionary to accompany the spreadsheet file.

Instructions: Fill out the information in each row for each variable in the spreadsheet; note that you will likely have more variables than columns in this table. Copy this information into a text document and save it next to the spreadsheet. It is useful to save the data dictionary with the same root name as its data file by appending “_dictionary” on the end of the file name; for example, the data dictionary for the file “myData.xlsx” would be “myData_dictionary.txt”.

Source: This exercise was adapted from “Leveling Up Data Management” [Briney, 2023].

| Question | Variable 1 | Variable 2 | Variable 3 |
|---------------------------------------|------------|------------|------------|
| Variable name | | | |
| Variable description | | | |
| Variable units | | | |
| Relationship to other variables | | | |
| Variable coding values and meanings | | | |
| Known issues with the data | | | |
| Anything else to know about the data? | | | |

Chapter 3

File Organization and Naming

Good file organization and naming are foundational data management practices, as they help you find files quickly when you need them. To set up file organization and naming conventions, this chapter offers two exercises: a card-sorting process for brainstorming a file organization system; and a worksheet for creating a file naming convention for a group of files.

3.1 Set Up a File Organization System

Description: *Implementing a file organization system is the first step toward creating order for your research data. Well-organized files make it easier to find the data you need without spending lots of time searching your computer. Every researcher organizes their files slightly differently, but the actual organizational system is less important than having a place where all of your files should logically go. This exercise prompts you to brainstorm organizational groupings and hierarchies to come up with an order for managing your research data.*

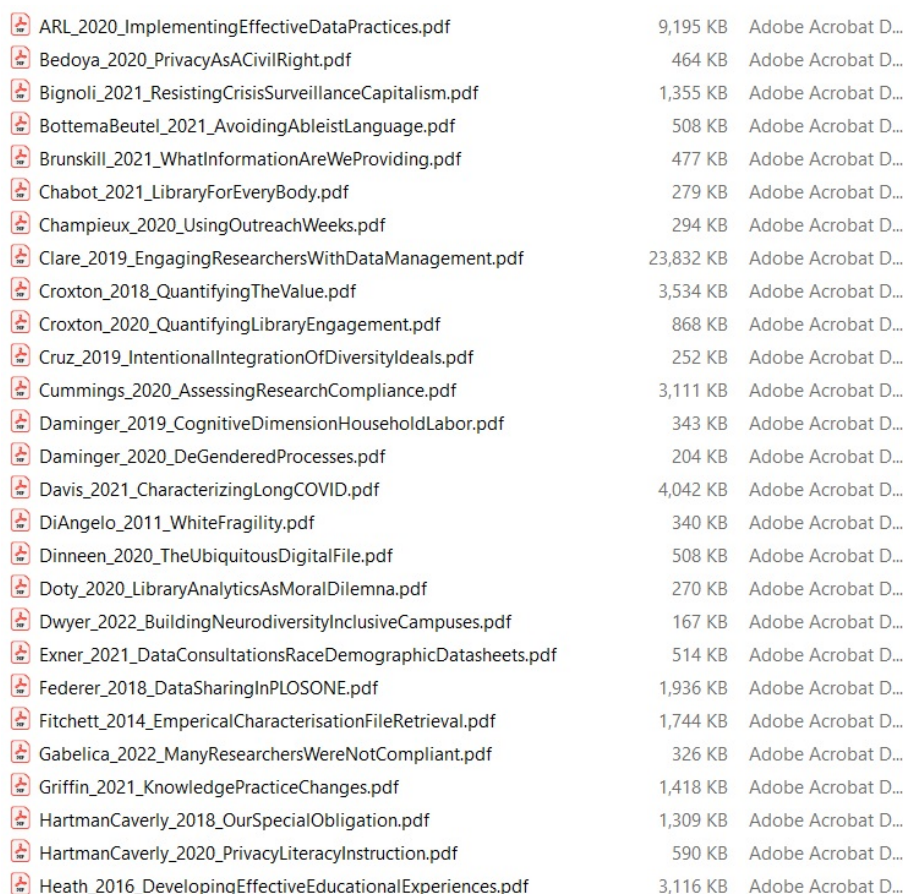
Instructions: *This is a card-sorting exercise, meaning you will need a stack of note cards or post-it notes to do this activity, ideally in three different colors. Follow the instructions to label cards and move them around until you develop your organizational system. There is no one correct way to do this so feel free to play around, add new cards, and move cards however you want! Once you put your new organizational system into place, be sure to always put your files where they're supposed to go.*

1. Take a stack of note cards or post-it notes in the first color and write the following labels on one card each, omitting any file types that you do not use in your research:
 - Raw data
 - Analyzed data
 - Code
 - Protocols
 - Article drafts
 - Figures
 - Your publications
 - Literature PDFs
 - Grant documents
 - Research notes
2. Move cards around and group together file-type cards that you want to store together. Files that will be stored near each other in a folder hierarchy, but not together, should be placed near each other while file types expected to be stored completely separately should be away from other cards.
3. Create hierarchies in file organization by adding new “folder” cards in one of two types; use different colored cards for each folder type:
 - Cards in the second color represent a single folder. These should be labeled with the folder name in quotations (e.g. “Literature” or “My publications”).
 - Cards in the third color represent a group of folders, such as for folders organized by date or by project. Use only one card to represent the organizational pattern that will be repeated. These cards should be labeled with the organizational system in square brackets (e.g. [By date] or [By project]).
4. Move existing file-type cards/groups of file-type cards underneath the new folder cards to show that a file type will be saved in a specific folder or group of folders. Organizational-group folders (cards in the third color) only need to be represented once in the card sorting, as they are assumed to represent multiple folders on a computer.
5. Make copies of any type of card and add folder levels, as needed. Adjust placement and hierarchies until you are happy with the organizational system you developed.
6. Record your organizational system in your lab notebook or in a README.txt.

3.2 Create a File Naming Convention

Description: *File naming conventions are a simple way to add order to your files and help to find them later. Rich and descriptive file names make it easier*

to search for files, understand at a glance what they contain, and tell related files apart.



| | | |
|---|-----------|--------------------|
| ARL_2020_ImplementingEffectiveDataPractices.pdf | 9,195 KB | Adobe Acrobat D... |
| Bedoya_2020_PrivacyAsACivilRight.pdf | 464 KB | Adobe Acrobat D... |
| Bignoli_2021_ResistingCrisisSurveillanceCapitalism.pdf | 1,355 KB | Adobe Acrobat D... |
| BottemaBeutel_2021_AvoidingAbleistLanguage.pdf | 508 KB | Adobe Acrobat D... |
| Brunskill_2021_WhatInformationAreWeProviding.pdf | 477 KB | Adobe Acrobat D... |
| Chabot_2021_LibraryForEveryBody.pdf | 279 KB | Adobe Acrobat D... |
| Champieux_2020_UsingOutreachWeeks.pdf | 294 KB | Adobe Acrobat D... |
| Clare_2019_EngagingResearchersWithDataManagement.pdf | 23,832 KB | Adobe Acrobat D... |
| Croxton_2018_QuantifyingTheValue.pdf | 3,534 KB | Adobe Acrobat D... |
| Croxton_2020_QuantifyingLibraryEngagement.pdf | 868 KB | Adobe Acrobat D... |
| Cruz_2019_IntentionalIntegrationOfDiversityIdeals.pdf | 252 KB | Adobe Acrobat D... |
| Cummings_2020_AssessingResearchCompliance.pdf | 3,111 KB | Adobe Acrobat D... |
| Daminger_2019_CognitiveDimensionHouseholdLabor.pdf | 343 KB | Adobe Acrobat D... |
| Daminger_2020_DeGenderedProcesses.pdf | 204 KB | Adobe Acrobat D... |
| Davis_2021_CharacterizingLongCOVID.pdf | 4,042 KB | Adobe Acrobat D... |
| DiAngelo_2011_WhiteFragility.pdf | 340 KB | Adobe Acrobat D... |
| Dinneen_2020_TheUbiquitousDigitalFile.pdf | 508 KB | Adobe Acrobat D... |
| Doty_2020_LibraryAnalyticsAsMoralDilemma.pdf | 270 KB | Adobe Acrobat D... |
| Dwyer_2022_BuildingNeurodiversityInclusiveCampuses.pdf | 167 KB | Adobe Acrobat D... |
| Exner_2021_DataConsultationsRaceDemographicDatasheets.pdf | 514 KB | Adobe Acrobat D... |
| Federer_2018_DataSharingInPLOSONE.pdf | 1,936 KB | Adobe Acrobat D... |
| Fitchett_2014_EmpiricalCharacterisationFileRetrieval.pdf | 1,744 KB | Adobe Acrobat D... |
| Gabelica_2022_ManyResearchersWereNotCompliant.pdf | 326 KB | Adobe Acrobat D... |
| Griffin_2021_KnowledgePracticeChanges.pdf | 1,418 KB | Adobe Acrobat D... |
| HartmanCaverly_2018_OurSpecialObligation.pdf | 1,309 KB | Adobe Acrobat D... |
| HartmanCaverly_2020_PrivacyLiteracyInstruction.pdf | 590 KB | Adobe Acrobat D... |
| Heath_2016_DevelopingEffectiveEducationalExperiences.pdf | 3,116 KB | Adobe Acrobat D... |

Figure 3.1: Screenshot of pdf's with consistent file names using the convention FirstAuthorLastName_YEAR_ShortTitle.pdf

Instructions: *This exercise guides researchers through the process of creating a file naming convention for a group of related files. Fill in each section for a group of related files following the instructions; an example for microscopy files is provided. This exercise may be redone as needed, as different groups of files require different naming conventions.*

Source: *This exercise is based on the “File Naming Convention Worksheet” [Briney, 2020a].*

1. What group of files will this naming convention cover?

You can use different conventions for different file sets.

Example: This convention will apply to all of my microscopy files, from raw image through processed image.

2. What information (metadata) is important about these files and makes each file distinct?

Ideally, pick three pieces of metadata; use no more than five. This metadata should be enough for you to visually scan the file names and easily understand what's in each one.

Example: For my images, I want to know date, sample ID, and image number for that sample on that date.

- 1.
- 2.
- 3.
- 4.
- 5.

3. Do you need to abbreviate any of the metadata or encode it?

If any of the metadata from step 2 is described by lots of text, decide what shortened information to keep. If any of the metadata from step 2 has regular categories, standardize the categories and/or replace them with 2- or 3-letter codes; be sure to document these codes.

Example: Sample ID will use a code made up of: a 2-letter project abbreviation (project 1 = P1, project 2 = P2); a 3-letter species abbreviation (mouse = "MUS", fruit fly = "DRS"); and 3-digit sample ID (assigned in my notebook).

4. What is the order for the metadata in the file name?

Think about how you want to sort and search for your files to decide what metadata should appear at the beginning of the file name. If date is important, use ISO 8601-formatted dates (YYYYMMDD or YYYY-MM-DD) at the beginning of the file names so dates sort chronologically.

Example: My sample ID is most important so I will list it first, followed by date, then image number.

- 1.
- 2.
- 3.
- 4.
- 5.

5. What characters will you use to separate each piece of metadata in the file name?

Many computer systems cannot handle spaces in file names. To make file names both computer- and human-readable, use dashes (-), underscores (_), and/or capitalize the first letter of each word in the file names.

Example: I will use underscores to separate metadata and dashes between parts of my sample ID.

6. Will you need to track different versions of each file?

You can track versions of a file by appending version information to end of the file name. Consider using a version number (e.g. “v01”) or the version date (use ISO 8601 format: YYYYMMDD or YYYY-MM-DD).

Example: As each image goes through my analysis workflow, I will append the version type to the end of the file name (e.g. “_raw”, “_processed”, and “_composite”).

7. Write down your naming convention pattern.

Make sure the convention only uses alphanumeric characters, dashes, and underscores. Ideally, file names will be 32 characters or less.

Example: My file naming convention is “SA-MPL-EID_YYYYMMDD_###_status.tif”

Examples are “P1-MUS-023_20200229_051_raw.tif” and “P2-DRS-285_20191031_062_composite.tif”.

8. Document this convention in a README.txt (or save this worksheet) and keep it with your files.

Chapter 4

Data Storage

All research data needs to be stored and backed up, but it can be frustrating to pick these systems and ensure that they are working correctly. This chapter consists of two exercises: a worksheet to document available storage and backup options and decide between them; and a procedure for testing that a backup system is working.

4.1 Store and Back Up the Data

Description: *Research data needs to be stored and backed up reliably so that important data is not lost. But storage is commonly a challenge, as institutions don't always offer uniform options for storage and backup. This exercise prompts you to examine the storage and backup systems available to you, either free or paid, as well as options outside of your institution, before determining which is the best set of options for your data.*

Instructions: *Fill out the first three tables for storage and backup systems that are: 1) institutional and free; 2) institutional and paid; and 3) outside of your institution. Examine all of the options, evaluating them based on the criteria listed below the first three tables to pick your ideal storage and backup configuration. Note that this exercise works best for data under 1 TB in size that has no security restrictions.*

What storage or backup systems are already available to you, such as through your institution or workplace, at no cost?

| Question | System 1 | System 2 | System 3 |
|--|----------|----------|----------|
| System name | | | |
| Is it storage or backup? | | | |
| What is the hardware type? | | | |
| Is the system backed up? | | | |
| For backup systems, is backup automatic? | | | |
| Is the system local or remote? | | | |
| Is there a limit to storage capacity? | | | |
| Who manages the system? | | | |
| Is it easy or difficult to use? | | | |

What storage or backup systems are available to you, such as through your institution or workplace, at a cost?

| Question | System 1 | System 2 | System 3 |
|----------------------------|----------|----------|----------|
| System name | | | |
| Is it storage or backup? | | | |
| What is the cost? | | | |
| What is the hardware type? | | | |
| Is the system backed up? | | | |

| Question | System 1 | System 2 | System 3 |
|--|----------|----------|----------|
| For backup systems, is backup automatic? | | | |
| Is the system local or remote? | | | |
| Is there a limit to storage capacity? | | | |
| Who manages the system? | | | |
| Is it easy or difficult to use? | | | |

Are there other storage or backup systems that you can use?

| Question | System 1 | System 2 | System 3 |
|--|----------|----------|----------|
| System name | | | |
| Is it storage or backup? | | | |
| What is the cost? | | | |
| What is the hardware type? | | | |
| Is the system backed up? | | | |
| For backup systems, is backup automatic? | | | |
| Is the system local or remote? | | | |
| Is there a limit to storage capacity? | | | |

| Question | System 1 | System 2 | System 3 |
|---------------------------------|----------|----------|----------|
| Who manages the system? | | | |
| Is it easy or difficult to use? | | | |

Optimize your storage and backups on the following considerations:

1. You need one storage system that will hold all of your data files, is easy to use, and is managed by someone you trust.
2. You need at least one backup – preferably two – that will hold all of your data files, is managed by someone you trust, is easy to use, and backs up automatically.
3. At least one backup should be in a different location than your main storage system for disaster resiliency.
4. Balance cost with making sure that storage/backup systems are managed by trustworthy parties and are not difficult to use.

Pick your storage and backup systems:

| Storage | Backup 1 | Backup 2 |
|---------|----------|----------|
| | | |

4.2 Test Your Backup

Description: Backups are super important for your data, so it's always good to test that your backups are still working. Nothing is worse than losing your data from your primary storage and then realizing that your backup isn't working either. Beyond checking that your backup is working, it's also good to know how to recover your files so that you don't have to learn this for the first time while panicking about lost data. This short exercise walks you through getting a file off your backup to test that it is working and to learn how the data-recovery process works.

Instructions: Pick a backup system and a file to recover and work through the steps. The hard part of this exercise is finding instructions for file recovery and recovering the file, which vary by backup system.

1. Identify where your data is backed up.
2. Find instructions for recovering data from your backup system.
3. Pick a data file from your computer.
4. Follow the instructions from step 2 to get a copy of the data file from step 3 out of your backup system.
5. If this process didn't work, fix your backup system. If this process did work, congrats your backup is working and you know how to recover your files!

Chapter 5

Data Management

While this entire workbook covers data management activities, it's often useful to take a step back and document the data management decisions that have been made. This chapter provides exercises in documenting data management in two areas: a worksheet for writing a living data management plan (which builds on exercises from previous chapters); and a worksheet for discussing roles and responsibilities around data management with your research collaborators.

5.1 Write a Living Data Management Plan (DMP)

Description: *Many researchers are aware of the two-page data management plan (DMP) for a grant application, but you may not be aware of the more useful type of DMP: a living DMP. This document describes how data will be actively managed during a project and may be updated whenever necessary to reflect current data practices. A living DMP is a useful touchstone for understanding where data lives, how it's labelled, how it moves through the research process, and who will oversee the data management. This exercise guides you through the process of creating a living DMP for your research.*

Instructions: *Pick a project and answer the following questions to build your living DMP. This DMP may be changed at any time to improve practices. If you are doing collaborative research, work through this exercise with your collaborators to agree on shared conventions.*

Write a short summary of what project this DMP is for:

How will your data be organized? (See Exercise 3.1: Set Up a File Organization System.)

What naming convention(s) will you use for your data? (See Exercise 3.2: Create a File Naming Convention.)

Where will data be stored? How will data be backed up? (See Exercise 4.1: Store and Back Up the Data.)

How will you document your research? Where will your research notes be stored?

How will data move through the collection and analysis pipelines?

Do you need to do any version control on your files? How will that be done?

As necessary, record any project roles and responsibilities around data management:

5.2 Determine Data Stewardship

Description: It is often helpful to be up front about requirements and permissions around research data. This exercise encourages you to discuss these issues with supervisors and peers to make sure that there are no misunderstandings about who has what rights to use, retain, and share data.

Instructions: Determine which research data should be discussed. Bring together the Principle Investigator, the researcher collecting the data, and anyone else who works with that data. As a group, answer the questions in the exercise, making sure that everyone agrees on the final decisions. Record the results of the discussion and save them with the project files.

Source: This exercise was adapted from the “Project Close Out Checklist” [Briney, 2020b].

What data is being discussed?

Are there security or intellectual property restrictions on the data and, if so, what are they?

Are there any requirements to publicly share the data and, if so, what are they?

Who will store the copy of record of the data and for how long?

Who is allowed to keep a copy of the data after the project ends?
Which data?

Who is allowed to reuse the data after the project ends? Are there
any requirements for reuse, such as co-authorship?

Who keeps any physical research notebooks after the project ends?

Chapter 6

Data Sharing

Sharing data that underlies research has become a common expectation within scholarly research. However, the landscape of data repositories is still uneven and many researchers are still learning best practices for data sharing. To help in this area, this chapter offers of two exercises: a decision tree-inspired worksheet for picking the best data repository for your data; and checklist for working through the process of sharing data in a data repository.

6.1 Pick a Data Repository

Description: *It can be difficult to know where to share research data as so many sharing platforms are available. Current guidance is to deposit data in data repository that will give you a DOI. This exercise guides you through the process of picking a data repository, starting with repositories for very specific types of data and defaulting to generalist data repositories. Note that some repositories charge fees for deposit, most often for large data (500 GB or larger).*

Instructions: *Identify the data that needs to be shared and work through repository selection from discipline-specific data repositories to more general data repositories. Once you have identified a repository for all of your data, the exercise is over and you do not need to answer any further questions. Note that, depending on data types, you may need to deposit your data into multiple repositories (for example, a discipline-specific repository for one type of data and an institutional data repository for the rest of the data).*

Identify all of the data that needs to be shared.

Is there a known disciplinary data repository for some or all of the data? For example, is there a data repository used by everyone in your research area or required for your data type by your funder?

If so, deposit some or all of your data there. End the exercise if the repository will accept all of your data or keep going if there is still some data left to deposit.

Review the list of recommended data repositories from PLOS. Is there a logical disciplinary data repository for some or all of your data?

If so, deposit some or all of your data there. End the exercise if you have shared all of your data or keep going if there is still some data left to deposit.

Does your institution have a data repository?

If so, deposit your data there and end the exercise.

Do you have a preferred generalist data repository?

If so, deposit your data there and end the exercise.

Pick a generalist data repository and deposit your data.

6.2 Share Data

Description: Data sharing is becoming common and expected by funding agencies and journals. While the process of depositing data into a data repository will vary between repositories, there are some common actions that should be taken to prepare data for sharing. This exercise walks you through these standard requirements for sharing data.

Instructions: This checklist enumerates the necessary steps and decisions to deposit data supporting a research article into a data repository. Identify the data to be shared and work through the list. Note that, if data will be shared as multiple deposits or in multiple repositories, the checklist should be worked separately for each group of data.

Data Selection

- ___ Select data that reproduces published results.
- ___ Perform quality control on the data files.
- ___ Convert data in proprietary file types to open file types, as appropriate; see Exercise 7.2: Convert Data File Types.
- ___ Determine if data will be shared in one group or as several deposits. If there will be multiple deposits in one repository or data will be divided across more than one data repository, work through the remainder of the checklist separately for each group of shared files.

Data Documentation

- ___ Document any spreadsheet data with a data dictionary (see Exercise 2.3: Create a Data Dictionary). The data dictionary should be shared with the other files.
- ___ Write a brief description of each data file, including any data dictionaries, and what it contains. Save this information in a README.txt file and share it with the other files.
- ___ Write a brief description of the overall data to be used during deposit process for the dataset Description/Abstract.

Sharing Information (Metadata)

- ___ Give the dataset a title. Default is “Data from: [name of the article]”.
- ___ Determine who will be listed as authors of the data and in what order; this may be different than the authors of the article.
- ___ Chose a license for reuse. Default is CC0 [Creative Commons Wiki, 2014].

Deposit Data

___ Pick a data repository using Exercise 6.1: Pick a Data Repository

___ Deposit the data and documentation files, and fill in metadata as determined above.

Share Data

___ Share data with its DOI.

___ Link the publication to its data, either in a Data Availability Statement or as a citation.

Chapter 7

Project Wrap Up

The end of a project is a key time to perform data management activities in order to set yourself up for future data reuse. This is because you still remember all of the important details about your data and can make good decisions about preparing it for the future. This chapter has three exercises to work through for project wrap up: a worksheet on converting data to more open file types; a checklist for populating a project Archive folder; and a checklist for preparing data for reuse, which leverages the previous two exercises.

This chapter also covers project wrap up in the form of separating from your institution. This checklist exercise for the departing researcher is important to work through so that critical data does not get lost in the transition. A fuller version of this checklist, intended for both the departing personnel and a project administrator to work through together, is also available [Goben and Briney, 2023].

7.1 Prepare Data for Future Use

Description: *The end of a project is a good time to prepare data for potential future reuse, as you still know the important details about the data to record and have access to any software used to create the data. This checklist exercise walks you through steps to gather your data into a central place and document the project. Working through the checklist results in project data being in one central location, well documented, and organized and formatted in a way to make future reuse easier.*

Instructions: *Gather all of the data from a project and work through the checklist to organize and document the data for future reuse. This exercise refers to several other exercises in the Workbook that should be completed during this process, if they have not been already*

Source: This exercise was adapted from the “Project Close Out Checklist” [Briney, 2020b].

Prepare Data

___ Move all data into one central project folder; this folder may have sub-folders and should be organized however makes sense for your data.

___ As necessary, work through Exercise 7.2: Convert Data File Types to copy data into more open/common file formats.

Back Up Your Research Notes

___ If your notes are electronic, save a copy in the project folder

___ If your notes are physical, scan them and save a copy in the project folder.

Create a Project Archive Folder

___ Work through Exercise 7.3: Create an Archive Folder.

___ Put the Archive folder in the project folder.

Create a Project-Level README File

___ If you haven’t done so already, work through Exercise 2.2: Project-Level README.txt.

___ Store a copy of the README file with the data.

Save Files in a Stable Location

___ Save the project folder on a storage system that you will have access to for the next several years.

7.2 Convert Data File Types

Description: Data is often stored in a file type that can only be opened by specific, costly software – this is referred to as a “proprietary file type.” You can tell that you have data in a proprietary file type if you lose access to the data when you lose access to the software. When data is in a proprietary file type, it’s always a good idea to copy the data into a more common, open file type as a backup; you may lose a bit of functionality, but it’s better to have a backup than to not have your data at all! This exercise works through identifying possible alternative file types for the data’s proprietary file type before instructing you to make a copy of the data in the new file type.

Instructions: For any data in a proprietary file type, identify the data and answer the following questions. Once you have picked a more open, common

file type, make a copy of the data in that file type but do not delete the original data.

Is your data stored in a proprietary file type? What file type and how does this limit future data reuse?

Is it possible to convert your data to other file types? If so, list the possible types:

Which of the possible file types are in common use? Which of the possible file types can be opened by multiple software programs?

Of the possible options above, do you have a preference for a specific file type?

Pick one of the more open or common file types and copy your important data files into that file type. Do not delete the original files.

7.3 Create an Archive Folder

Description: *To save your future-self time spent digging through all of your research files, set aside the most important files into a separate “Archive” folder.*

Do this at the end of the project while you still remember which files are important and where they are located. The Archive folder should only contain a small subset of the most important documents that are likely to be reused; you may still need to go through all of your files but, in the majority of instances, you will save time by easily finding what you need in the Archive folder.

Instructions: *This exercise consists of a checklist of the key documents that are likely to be most useful in a research project archive. Create a separate folder within the larger project folder (or in a highly visible place within the storage system) labelled “Archive”. Copy – do not move – the files on this checklist into the Archive folder. Add copies other important research documents, as needed. Remember, the Archive folder does not need to be comprehensive, so focus on the subset of files that are most likely to be reused or referenced in the future.*

Source: *This exercise was adapted from the “Project Close Out Checklist” [Briney, 2020b].*

Project Documentation

___ README file of project information

Data Snapshots

___ Important raw data

___ Key data analyses

___ Final data

Code

___ Analysis code

___ Record software version, as appropriate

Other Research Documents

___ Protocols

___ Survey instruments

Research Notes

___ Scan of research notebook

___ Digital notes

Images

___ Flat files of figures (e.g. .JPG or .TIFF)

___ Editable image files (e.g. Photoshop)

Publications

- ☐ Published article in .PDF format
- ☐ Final version of the article in editable document format (e.g. .DOCX)
- ☐ Posters

Administrative Documents

- ☐ Grant proposals
- ☐ Grant progress reports and final report

7.4 Separate from the Institution

Description: Researchers regularly leave institutions in order to take new jobs. For how common this occurrence is, it represents a critical transition during which data may be lost. This checklist enumerates a number of important steps that researchers can take to ensure that they retain the appropriate data yet leave behind what belongs to the institution.

Instructions: The researcher leaving the institution should work through this checklist to ensure they keep the proper information while returning what does not belong to them. The researcher and project administrator may also jointly work through the extended version of this exercise, the Data Departure Checklist [Goben and Briney, 2023].

Source: This exercise was adapted from the “Data Departure Checklist” [Goben and Briney, 2023].

Return Lab Hardware

- ☐ Keys
- ☐ Key card
- ☐ Individual computer / workstation
- ☐ Tablet(s)
- ☐ Phone
- ☐ Peripherals (e.g. keyboard, mouse, monitor)
- ☐ Headsets, webcams
- ☐ External drives
- ☐ Other lab equipment (e.g. cameras, recording devices)

Delete Personal Information and Remove Personal Devices

- ___ Remove personal information from lab devices
- ___ Remove personal devices from lab
- ___ Remove personal access to shared accounts (e.g. lab Github, lab repository page, lab website, mailing lists, or social media)

Retain Copies of Data that You Have Permission to Keep

- ___ If you have not done so already, work through Exercise 5.2: Determine Data Stewardship to determine what data you may retain
- ___ Identify and keep pertinent research data from personal devices
- ___ Identify and keep pertinent research data from storage systems (e.g. AWS/Azure, Box, campus HPC, Dropbox, Electronic Lab Notebook, Globus, Google Drive, lab/department/college servers, Microsoft OneDrive, Microsoft Sharepoint, or shared collaborator drives)
- ___ When appropriate, make a copy of research notes

Update Research Administration Documents, As Necessary

- ___ Update/transfer Institutional Review Board
- ___ Update/transfer IACUC
- ___ Update/transfer Data Use Agreements (DUA)
- ___ Update/transfer Material Transfer Agreements (MTA)
- ___ Update/transfer research grants

Handle Email

- ___ Set out of office, providing forwarding information
- ___ Forward/backup important emails
- ___ Check with University Archivist or Records Manager for retention policies (depends on rank)

Bibliography

Kristin Briney and Becky Yoose. *Managing Data for Patron Privacy: Comprehensive Strategies for Libraries*. ALA Editions, Chicago, IL, 2022. URL <https://www.alastore.ala.org/mdpp>.

Kristin Briney, Heather Coates, and Abigail Goben. Foundational Practices of Research Data Management. *Research Ideas and Outcomes*, 6:e56508, July 2020. ISSN 2367-7163. doi: 10.3897/rio.6.e56508. URL <https://riojournal.com/article/56508/>.

Kristin A. Briney. *Data Management for Researchers : Organize, Maintain and Share Your Data for Research Success*. Pelagic Publishing, 2015. URL <https://pelagicpublishing.com/products/data-management-for-researchers-briney>.

Kristin A. Briney. File Naming Convention Worksheet, June 2020a. URL <https://doi.org/10.7907/894q-zr22>.

Kristin A. Briney. Project Close-Out Checklist for Research Data, May 2020b. URL <https://doi.org/10.7907/yjph-sa32>.

Kristin A. Briney. Leveling Up Data Management, June 2023. URL <https://doi.org/10.7907/syk7-3z92>.

Creative Commons Wiki. CC0 use for data, 2014. URL https://wiki.creativecommons.org/wiki/CC0_use_for_data.

Abigail Goben and Kristin A. Briney. Data Departure Checklist, August 2023. URL <https://doi.org/10.7907/h314-4x51>.