

Ten Simple Rules for Data Storage

Edmund Hart ^{1,*}, Pauline Barmby ², Jeffrey Hollister ³, David LeBauer ⁴, Naupaka Zimmerman ⁵, Kara H. Woo ⁶, Sarah Mount ⁷, Timothée Poisot ⁸, François Michonneau ⁹,

1 Univeristy of Vermont, Department of Biology, Burlington

2 University of Western Ontario, Department of Physics and Astronomy

3 US Environmental Protection Agency, Atlantic Ecology Division

4 University of Illinois at Urbana-Champaign, National Center for Supercomputing Applications and Institute for Genomic Biology

5 University of Arizona, School of Plant Sciences

6 Washington State University, Center for Environmental Research, Education, and Outreach

7 University of Wolverhampton, School of Mathematics and Computer Science

8 Université de Montréal, Département de Sciences Biologiques

9 University of Florida, iDigBio, Florida Museum of Natural History, Gainesville, FL 32611-7800

* E-mail: edmund.m.hart@gmail.com

Introduction

Data is the central currency of science, but the nature of scientific data has changed dramatically with the rapid pace of technology. This rapid change has led to an increasing heterogeneity of data formats, dataset sizes, data complexity, and data use cases (including, e.g., the importance of data sharing). For example, improvements in high throughput DNA sequencing, sustained institutional support for large sensor networks [1], and sky surveys with large-format digital cameras [3] have led to the creation of massive quantities of data. At the same time, collaboration between researchers is becoming increasingly common [4], necessitating increased coordination between researchers collecting data [5]. These changes mean when modern scientists refer to data, they could be referring to petabytes of information stored in professionally-maintained databases, to Excel spreadsheets on a single computer, or to hand-written tables in lab notebooks on shelves. All of these forms of data remain important, but methods of data curation and storage must continue be updated in order to encompass the changes brought about by new forms of data collection and storage.

While much has been written about both the virtues of data sharing [6] and best practices to do so [8], how to store data has received much less attention. Proper storage of data is a prerequisite for any sharing, and indeed lack of proper storage may also contribute to the phenomenon of data decay or “data entropy”: as time passes data becomes less and less accessible (whether or not it has been publicly shared) [10]. Best practices for data storage often begin and end with “use a community standard repository” and this is, by all means, a great practice. However, data storage policies are highly variable between repositories [14], and best practices across all stages of the data life cycle will facilitate transition from local storage to repository. Good storage

practices are important even (or especially) in cases where data may not fit with an existing repository, or in the cases where only derived data products (versus raw data) are suitable for deposition, or in the case where an existing repository may have lax standards. Therefore, this manuscript describes 10 simple rules for data storage that grew out of a long discussion among instructors for the Software Carpentry initiative [15]. Software Carpentry instructors are scientists from diverse backgrounds who have encountered a variety of data storage challenges and are active in teaching other scientists best practices for scientific computing and data management. Thus, this paper represents a distillation of collective experience, and hopefully will be useful to scientists facing a variety of data storage challenges.

Rule 1: Know what to expect

Most of the troubles encountered during the analysis, management, and release of data can be avoided by having a clear roadmap of what to expect *before* data acquisition starts. For instance:

- How will the raw data be received?
- What is the format expected by the software used for analysis?
- Is there a community standard on the format for release?
- Does the study involve simulations, and what is the model output?

The answers to these questions can range from simple cases (e.g., sequencing data stored in the FASTA format, which can be used as is throughout the analysis), to experimental designs involving several instruments, each with its own output format. Knowing the state in which the data needs to be at each step can help (i) identify software tools to use in converting across data formats, (ii) orient technological choices about how and where the data should be stored, and (iii) rationalize the analysis pipeline, making it more amenable to re-use.

Another side of preparedness is the ability to estimate the volume of disk space needed to store the data at each step. The required strategy will differ for datasets of varying size. ‘Lighter’ datasets (e.g. datasets that are only a few megabytes in size) can be managed locally with a simpler data management plan, whereas larger datasets (e.g. gigabytes to terabytes and even petabytes) will in almost all cases require careful planning and preparation (Rule 9).

Rule 2: Know your use case

Researchers should know their use case and store data appropriately. This involves answering the following questions, ideally prior to beginning data collection:

- Should the raw data be archived (Rule 3)?
- Should the data used for analysis be prepared once, or re-generated from the raw data each time (and what difference would this choice make for storage, computing requirements, and reproducibility)?
- Will manual corrections be necessary (as opposed to programmatic approaches)?
- How will changes to the data be tracked, and where will these tracked changes be logged?
- Will the final data be released, and if so, in what format?
- Are there restrictions or privacy concerns associated with the data (e.g. for survey results)?
- Will institutional validation be required prior to releasing the data?

- Does the agency funding the research require data deposition, and if so, where? 71
- Does the journal in which you plan to publish require data deposition? 72

None of these questions have universal answers, nor are they the only questions one should ask before starting data acquisition. But similarly to Rule 1, knowing the what, when, and how of *your* use of the data will bring you close to a reliable roadmap on how to handle data from acquisition through publication to archival. 73 74 75 76

Rule 3: Keep raw data raw 77

Since analytical and data processing procedures may improve or otherwise change over time, having access to the ‘raw’ or unprocessed data can help facilitating future re-analysis and analytical reproducibility. As processing algorithms and computational power increase, new analyses will be enabled that were not possible at the time of the original work. If only derived data are stored, it can be difficult to impossible for other researchers to confirm analytical results, to assess the validity of statistical models, or to directly compare findings across studies. 78 79 80 81 82 83 84

Therefore, data should always be kept in raw format whenever possible (within the constraints of technical limitations). In addition to being the most appropriate way to ensure transparency in analysis, having the data stored and archived in their original state gives a common point of reference for derivative analyses. Despite the intuitive value of this approach, it is not always clear what constitutes sufficiently “raw” data (e.g., ohms off a temperature sensor or images off an Illumina sequencing flowcell are generally not archived after the initial processing). However, we focus here on the spirit of the rule. Data should be as “pure” as possible when they are stored. If derivations occur, they should be documented by also archiving relevant code and intermediate data sets. 85 86 87 88 89 90 91 92 93 94

The US National Ecological Observatory Network (NEON) handles this issue with a schema for various “levels” of data products that pertain to the amount of processing that has been performed on each (see here for a brief overview). In this case, raw data can include such products as voltage measurements or unprocessed LIDAR returns. These represent a tremendous amount of data; sharing this level of data often requires physically sending a hard drive through the postal service. NEON has handled this by writing detailed “Algorithm Theoretical Basis Documents” (ATBD’s) documenting the different processing “levels”. This approach is based around a similar one developed for the NASA EOSDIS program, and makes clear exactly what has been done to each dataset to derive it from its raw form. These levels, which start at 0 for raw data, and increase with the amount of derivation and processing, are also analogous to the levels used by the National Aeronautics and Space Administration (NASA) and National Oceanic and Atmospheric Administration (NOAA) uses for satellite data sets. 95 96 97 98 99 100 101 102 103 104 105 106 107

Rule 4: Store data in open formats 108

To maximize accessibility and long-term value, data should be stored in file formats whose specifications are freely available. The appropriate file type will depend on the type of data being stored (e.g. numeric measurements, text, images, video), but the key idea is that data should not require proprietary software or hardware, or a license, to be accessed. Proprietary formats can change, maintaining organizations can go out of business, and changes in license fees could make access to data in proprietary formats unaffordable for those who would make use of them. Examples of open data formats include comma-separated values (CSV) for tabular data, hierarchical data format (HDF) and NetCDF for hierarchically structured scientific data, portable 109 110 111 112 113 114 115 116 117

network graphics (PNG) for images, and extensible markup language (XML) for documents. Examples of closed formats include DWG (for AutoCAD drawings), Photoshop document (PSD, for bitmap images), Windows Media Audio (WMA), and Microsoft Excel (need refs?). Even if day-to-day processing uses closed formats, for example due to software requirements, data being stored for archival purposes should be stored in open formats. This is generally not prohibitive; most closed-sourced software enables users to export data to an open format.

Rule 5: Data should be uniquely identifiable

To aid reproducibility, the data used in a scientific publication should be uniquely identifiable. Ideally, datasets should have a unique identifier such as a Document Object Identifier (DOI), Archival Resource Key (ARK), or a Persistent URL (PURL). An increasing number of online services, such as Figshare, Zenodo, or DataOne are able to provide these.

Datasets may evolve over time. In order to distinguish between different versions of the same data, each dataset should have a distinct name, which includes a version identifier. A simple way to do this is to use date stamps as part of the dataset name. To avoid regional ambiguities, it is wise to use the ISO 8601 standard, which mandates the date format YYYY-MM-DD (i.e. from largest time unit to smallest). For example, the date “February 1st, 2015” could be written as 01-02-2015 by someone located in the UK. However, someone in the US could interpret it as being “January 2nd, 2015”. In ISO 8601 format it has the unambiguous canonical form: 2015-02-01.

Semantic versioning, as described in [16], is a richer approach to solving the same problem. An example of this can be seen in the CellPack datasets [17]. A semantic version number takes the form: **Major.Minor.Patch**, e.g. 0.2.7. The **major version** numbers should be incremented (or *bumped*) when a dataset scheme has been updated, or some other change is made that is not compatible with previous versions of the data with the same major version number. This might mean that an experiment using version 1.0.0 of the dataset could not be run on version 2.0.0 without making some changes to the data analysis. The **minor version** should be bumped when a change has been made which is compatible with older versions of the data with the same Major version. This means that any analysis that can be performed on version 1.0.0 of the data should be repeatable with version 1.1.0 of the data. The **patch version** number should be bumped when typos or bugs have been fixed. For example version 1.0.1 of a dataset may fix a typo in version 1.0.0.

Rule 6: Link relevant metadata

It should be almost impossible to separate data from its associated metadata. The importance of metadata for context, reusability, and discovery has been written about at length in many guides for data best practices [12].

Metadata should be as comprehensive as possible, use the relevant standards of your discipline, and be machine-readable (e.g., XML, JSON). Metadata should always accompany a data set, wherever it is stored. How best to do this depends on the format of the data. Formats such as NetCDF or HDF5 allow for embedded metadata, so the data and metadata are always together. When using a database, metadata tables should be clearly labeled and linked to the relevant data. Ideally a schema will be provided that also shows the linkages between data tables and metadata tables. Another scenario is a set of flat text files—in this case a semantically versioned, compressed archive should be created that includes metadata file(s).

Whatever format is used for archiving, the goal should be to make the link between metadata and data as clear as possible. The best approach is dependent on the archiving plan used, but even if the data set is archived solely for personal use, metadata will provide crucial context for future reuse.

Rule 7: Adopt the proper privacy protocols

In data sets where privacy is important, be sure to have a plan in place to protect data confidentiality. You should consider the different data stakeholders when developing privacy protocols for your data storage. These stakeholders might include funding agencies, human subjects or entities, collaborators, and your own needs. Both the NSF and NIH have data sharing policies in their grant guidelines requiring that personally identifiable information not be shared and that human subject data must be anonymized. If your data set is small, with minimal personal information, a hashing scheme can be used to anonymize personal information. Make sure to not store the hashing scheme with the data to prevent inadvertent sharing and don't use a commonplace hashing technique. Famously, New York City officials shared what they thought was anonymized data on cab drivers and over 173 million cab rides. However, it was quickly recognized that the city anonymized the data with a simple MD5 hashing scheme and all 20 GB of data was de-anonymized in a matter of hours [19]. Sometimes, however, the data itself allows identifiability. This is the case with human genomic data, where the data itself can be used to identify a subject [20]. Therefore, the type of data you have needs to be carefully considered as well. Methods for dealing with these complex issues at the intersection of data storage and privacy are rapidly evolving. Ideas such as storing changes against a reference genome, which helps with privacy and data volume [21], or bringing computation to data storage facilities instead of vice versa are still being developed [23]. Having a plan for privacy before you store your data is important, because it can determine in part how best to (or how you must) store data in the first place.

Rule 8: Have a systematic backup scheme

Every storage medium can fail, and every failure can result in loss of data. Researchers should therefore ensure that data is backed up at all stages of the research process. Data stored on local computers or institutional servers during the collection and analysis phases should be backed up to other locations and formats to protect against data loss. No backup system is failsafe (see the stories of the Dedoose crash and the near deletion of Toy Story 2), so more than one backup system should be used. Kristin Briney advocates the "Rule of 3" for backing up data: two onsite copies (such as on a computer, an external hard drive, or a tape) and one offsite copy (e.g. in cloud storage). Keeping backups in multiple locations protects against data loss due to theft, natural disasters, etc.

Researchers should also test their backups regularly to ensure that they are functioning properly. Common reasons for backup failure include:

- faulty backup software
- incorrect configuration (e.g., not backing up sub-directories)
- encryption (e.g., someone has encrypted the backups but lost the password)
- media errors

Consider the backup plan of your selected data repository before publishing your data. Many repositories mirror the data they host on multiple machines. If possible,

find out about the long-term storage plans of the repository. Are there plans in place to keep data available if the organization that manages the repository dissolves?

Rule 9: The location and method of data storage depends on how much you have.

The storage method you should choose depends on the size and nature of your data; the cost of storage, the time it takes to transfer the data, how the data will be used and any privacy concerns. Data is increasingly generated in the range of many terabytes by environmental sensors, satellites, automated analytical tools, simulation models, and genomic sequencers. Even larger data generating machines like the Large Hadron Collider (LHC) and the Large Scale Synoptic Survey Telescope (LSST) generate many TBs per day, rapidly accumulating to PB scale over the course of any particular study. While the cost of storage continues to decrease, the volume of data to be stored impacts the choice of storage methods and locations: for large datasets it is necessary to balance the cost of storage with the time of access and costs of re-generating the data.

When data takes too long to transfer or is costly to store, it can become more efficient to use a computer that can directly access and use the data in place. Inactive data can be put in longer-term storage; this is less expensive, but can take longer to retrieve. Some storage systems automatically migrate ‘stale’ files to longer term storage. Alternatively, some computing can be done ‘in the database’ or ‘on disk’ via database query languages (e.g. SQL, MapReduce) that perform basic arithmetic or via the use of procedural languages (e.g. R, Python, C) embedded in the database server. Modern database technologies such as HDFS and Spark allow these computations to be done on data of almost any size.

If you regularly only need access to a small subset of your data or need to share it with many collaborators, a web based API might be a good solution. Using this method, many users can send requests via HTTP to a web service which can subset the data, perform in database computation, and return smaller volumes of data as specific slices. Tools based on web services make it easier to find and download data, and facilitate analysis via reproducible scripts, however they can lead to excessive and careless abuse of resources. The time required to re-download and recompute results can be reduced by ‘caching’. Caching stores copies of downloads and generated files that are recognized when the same script is run multiple times.

<!-- “When data is larger than RAM, it can be handled by a ‘big memory’ node - these are currently 1-4 TB. This allows the user to read in and use a large dataset without special tools.” I think this needs some explanation I can’t provide.-->

Rule 10: Data should be stored in a machine readable-format

Not only data should be stored in an open format, to ensure that it will be easily and widely accessible (Rule 4), but data should also be stored in a format that computers can easily make use of.

As datasets become increasingly larger, it is crucial that they can be parsed efficiently. This is best achieved by using standard data formats that have clear specifications (e.g., CSV, XML, JSON, HDF5). Such data formats can be handled by a variety of programming languages, as efficient and well-tested libraries for parsing them are typically available. These standard data formats also ensure interoperability,

facilitate re-use, and reduce the chances of data loss or mistakes being introduced during conversion between formats.

When data can be easily imported into familiar software, whether it be a scripting language, a spreadsheet, or any other computer program that can import these common files, data become easier to re-use. Computer source code: the human readable software code that uses data, provide meta-data as well, making the analysis more transparent, such that all assumptions are implicitly stated in a human readable script. This also enables extraction of the analyses performed, their reproduction, and their modification. This principle is exemplified by the scripts used to import allometric measurements into the BAAAD database.

To take full advantage of data, it can be useful for it to be structured in a way that make manipulation and analysis easy. One such structure for data has been named *tidy* data [24]. Technically known as the ‘third normal form’, each variable is a column, each observation is a row, and each type of observational unit is a table. When data is organized in this way, the duplication of information is reduced and it is easier to subset or summarize the dataset to include the variables or observations of interest.

Interoperability is facilitated when variable names are mapped to existing data standards. For instance, for biodiversity data, the Darwin Core Standard provides a set of terms that describe observations, specimens, samples, and related information for a taxa. Because each term is clearly defined and documented, each dataset can use the terms consistently, facilitating data sharing across institutions, applications, and disciplines.

With machine-readable data, it is also easier to build an Application Programming Interface (API) to query the dataset to retrieve a subset of interest.

11: Ask a librarian!

Academic librarians are increasingly assisting researchers in the annotation, storage, and identification of data.

Conclusions

Acknowledgements

National Center for Supercomputing Applications. Software Carpentry Foundation. iDigBio/NSF.

Figure Legends

Figures here: Will need to figure out numbering...

Tables

Tables here: Will need to figure out numbering...

References

1. Reid JG, Carroll A, Veeraraghavan N, Dahdouli M, Sundquist A, English A, et al. Launching genomics into the cloud: deployment of Mercury, a next generation

sequence analysis pipeline. *BMC bioinformatics*. 2014;15: 30.
doi:10.1186/1471-2105-15-30

2. Hampton SE, Strasser C a, Tewksbury JJ, Gram WK, Budden AE, Batcheller AL, et al. Big data and the future of ecology. *Frontiers in Ecology and the Environment*. 2013; 130312142848005. doi:10.1890/120103

3. Eisenstein DJ, others. SDSS-III: Massive Spectroscopic Surveys of the Distant Universe, the Milky Way, and Extra-Solar Planetary Systems. *The Astronomical Journal*. 2011;142: 72.

4. Adams J. Collaborations: The rise of research networks. *Nature*. Nature Publishing Group, a division of Macmillan Publishers Limited. All Rights Reserved. 2012;490: 335–6. doi:10.1038/490335a

5. Fraser LH, Henry HA, Carlyle CN, White SR, Beierkuhnlein C, Cahill JF, et al. Coordinated distributed experiments: an emerging tool for testing global hypotheses in ecology and environmental science. *Frontiers in Ecology and the Environment*. Ecological Society of America; 2013;11: 147–155. doi:10.1890/110279

6. Wolkovich EM, Regetz J, O'Connor MI. Advances in global change research require open science by individual researchers. *Global Change Biology*. 2012;18: 2102–2110. doi:10.1111/j.1365-2486.2012.02693.x

7. Roche DG, Lanfear R, Binning SA, Haff TM, Schwanz LE, Cain KE, et al. Troubleshooting public data archiving: suggestions to increase participation. *PLoS biology*. 2014;12: e1001779. doi:10.1371/journal.pbio.1001779

8. White E, Baldrige E, Brym Z, Locey K, McGlinn D, Supp S. Nine simple ways to make it easier to (re)use your data. *Ideas in Ecology and Evolution*. 2013;6: 1–10. doi:10.4033/iee.2013.6b.6.f

9. Goodman A, Pepe A, Blocker AW, Borgman CL, Cranmer K, Crosas M, et al. Ten simple rules for the care and feeding of scientific data. *PLoS computational biology*. 2014;10: e1003542. doi:10.1371/journal.pcbi.1003542

10. Pepe A, Goodman A, Muench A, Crosas M, Erdmann C. How Do Astronomers Share Data? Reliability and Persistence of Datasets Linked in AAS Publications and a Qualitative Study of Data Practices among US Astronomers. Golden AA-J, editor. *PLoS ONE*. 2014;9: e104798. doi:10.1371/journal.pone.0104798

11. Vines TH, Albert AYK, Andrew RL, Débarre F, Bock DG, Franklin MT, et al. The availability of research data declines rapidly with article age. *Current biology : CB*. Elsevier; 2014;24: 94–7. doi:10.1016/j.cub.2013.11.014

12. Michener WK, Jones MB. Ecoinformatics: supporting ecology as a data-intensive science. *Trends in ecology & evolution*. 2012;27: 85–93. doi:10.1016/j.tree.2011.11.016

13. Michener WK, Brunt JW, Helly JJ, Kirchner TB, Stafford SG. Nongeospatial metadata for the ecological sciences. *Ecological Applications*. Eco Soc America; 1997;7: 330–342.

14. Marcial LH, Hemminger BM. Scientific data repositories on the Web: An initial survey. *Journal of the American Society for Information Science and Technology*. 2010;61: 2029–2048. doi:10.1002/asi.21339

15. Wilson G. Software Carpentry: lessons learned. *F1000Research*. 2014;3: 62. doi:10.12688/f1000research.3-62.v1

16. Preston-Werner T. Semantic Versioning 2.0.0. <http://semver.org>; 2014.

17. Johnson GT, Goodsell DS, Autin L, Forli S, Sanner MF, Olson AJ. 3D molecular models of whole HIV-1 virions generated with cellPACK. *Faraday Discuss*. The Royal Society of Chemistry; 2014;169: 23–44. doi:10.1039/C4FD00017J

18. Strasser C, Cook R, Michener W, Budden A. Primer on Data Management : What you always wanted to know [Internet]. California Digital Libraries; 2012. doi:10.5060/D2251G48

19. Goodin D. Poorly anonymized logs reveal NYC cab drivers' detailed whereabouts [Internet]. 2015. Available: [http://arstechnica.com/tech-policy/2014/06/poorly-anonymized-logs-reveal-nyc-ca](http://arstechnica.com/tech-policy/2014/06/poorly-anonymized-logs-reveal-nyc-cab-drivers-whereabouts/)
20. Homer N, Szelinger S, Redman M, Duggan D, Tembe W, Muehling J, et al. Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays. *PLoS genetics*. Public Library of Science; 2008;4: e1000167. doi:10.1371/journal.pgen.1000167
21. Kahn SD. On the future of genomic data. *Science* (New York, NY). 2011;331: 728–9. doi:10.1126/science.1197891
22. Wandelt S, Bux M, Leser U. Trends in genome compression. *Current Bioinformatics*. Bentham Science Publishers; 2014;9: 315–326.
23. Gaye A, Marcon Y, Isaeva J, LaFlamme P, Turner A, Jones EM, et al. DataSHIELD: taking the analysis to the data, not the data to the analysis. *International journal of epidemiology*. 2014;43: 1929–44. doi:10.1093/ije/dyu188
24. Wickham H. Tidy Data. *Journal of Statistical Software*. 2014;59: 1–23. Available: <http://www.jstatsoft.org/v59/i10>