

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

Abstract

Introduction

Data is the central currency of science. However the nature of scientific data has changed dramatically with the rapid pace of technology, leading to an increasing heterogeneity. Improvements in high throughput DNA sequencing and sustained institutional support for large sensor networks [1] have lead to the creation of massive quantities of data. At the same time collaboration between researchers is becoming increasingly common [3] with increased coordination between researchers collecting data [4]. These changes mean that data is more heterogenous than ever before in both size and complexity. It ranges from petabytes of information stored in hadoop clusters to Excel spreadsheets and lab notebooks on shelves.

While much has been written about both the virtues of data sharing [5] and best practices to do so [7] how to store data is a much less discussed topic. 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, as time passes data is less and less accessible (publicly shared or not) [9]. 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 [11] and data best storage practices will facilitate transition from local storage to repository. Furthermore your data may not fit with an existing repository, or only derived data products (vs raw data) are suitable, or an existing repository may have lax standards. What follows are 10 simple rules for data storage. These ideas grew out of a long discussion between Software Carpentry instructors, scientists from diverse backgrounds who have encountered 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* the data acquisition starts. How will the raw data be presented? In what format should they be for analysis? Does the study involves simulations, and what is the model output? Is there a community standard on the format for release? This can range from simple cases (sequencing data in the fasta format, that 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) create converters from these data, (ii) orient technological choices about how and where these data should be stored, and (ii) rationalizes the analysis pipeline, and make it more amenable to re-use.

Another side of preparedness is the ability to estimate the volume needed to store these data at each step. The strategy to apply is not the same when the total amount of data is in the order of a few Mb, than when it reaches the Gb or Tb sizes. Although (and we do not condone this practice) lighter datasets can be managed without much of a data management plan, larger ones require careful planning and preparation (see Rule 9).

Rule 2: Know your use case

Researchers should know their use case and store data appropriately. This involves answering the following questions. Should the raw data be archived (see rule 3)? Should the data used for analysis be prepared once, or re-generated from the raw data (and what difference does it means for storage and computing requirements)? Should the final data be released, and in what format? How do you track the changes made to the data, and where are they logged? Do you anticipate to make manual corrections, and why? Are there restrictions on the data, and how can you make them (*e.g.* for survey results) anonymous? Do you need validation from within your institution to release the data? Does your funding agency requires data deposition, and are there some specific platforms? Does the journal in which you plan to publish requires data deposition? 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 what *you* will do with the data, when, and how, will bring you close to a very detailed roadmap on how to handle these data from their acquisition to their publication.

Rule 3: Keep raw data raw

Since analytical and data processing procedures may improve or otherwise change over time, having access to the 'raw' or unprocessed data helps to facilitate 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 compare findings directly across studies.

Therefore, data should always be kept in a 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 subsequent data sets.

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 mailing the files on a physical drive. 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 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 the National Oceanic and Atmospheric Administration (NOAA) uses for satellite data sets.

Rule 4: Store data in open formats

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 simply unaffordable. Examples of open data formats include comma-separated values (CSV) for tabular data, hierarchical data format (HDF) for scientific data, portable network graphics (PNG) for images and extensible markup language (XML) for documents. Examples of closed formats include DWG (for AutoCAD drawings), Photoshop document (PSD), and Windows Media Audio (WMA) (need refs?). At a minimum, data being stored for archival purposes should be stored in open formats, even if day-to-day processing uses closed formats, for example due to software requirements.

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 01-02-2015 could be in January (US format) or February (UK format), but in

ISO 8601 format it has the canonical form 2015-02-01.

Semantic versioning, as described in [12], is a richer approach to solving the same problem. An example of this can be seen in the CellPack datasets [13].

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

You should make it almost impossible to separate your data from your metadata. The importance of metadata for context, reusability and discovery has been written about at length in many guides for data best practices [14]. It goes without saying these guidelines should be followed. Metadata should be as comprehensive as possible, use the relevant standards of your discipline, and be in a machine readable format (XML, JSON) That metadata should always accompany your data set wherever it is stored. How best to do this depends on the format of your archive. Formats such as NetCDF or HDF5 allow for embedded metadata so the data and metadata are always together. If you are 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 include metadata file(s). Whatever your archiving format, the goal should make the link between metadata and data as clear as possible. The best approach is dependent on your archiving plan, but even if your archive on just for yourself, metadata will provide future you with important context.

Rule 7: Rule

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 phase 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 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. Reasons they might not include:

- faulty backup software

- incorrect configuration (e.g., not backing up sub-directories) 163
- encryption (e.g., someone has encrypted the backups but lost the password) 164
- media errors 165

and many others. 166

Consider the backup plans of data repositories before publishing your data. Many 167
repositories mirror the data they host on multiple machines. If possible, find out 168
about the long-term storage plans of the repository. Are there plans in place to keep 169
data available if the organization that manages the repository dissolves? 170

Rule 9: Data size matters / requires special 171 considerations 172

- #39 and related GH issues #16, #19, #25 173
- Size classes: 174
 - larger than RAM 175
 - larger than HD space 176
 - larger than data storage server 177
- Storage method depends on the size of data; storage costs, transfer time, and 178
computing costs can become substantial. 179
 - data generated by simulation and derived data should consider cost of 180
storage vs. the cost of re-generating output. 181
 - For analyses of large data sets, the speed of reading and writing data can 182
limit the speed of computation. 183
- Larger data sets that are actively used in analysis should be stored on a disk 184
that is attached to a computer rather than being moved around between analysis 185
and storage. 186
 - inactive data can be put in longer-term storage; this is less expensive, but 187
can be slow to retrieve. Archiving of ‘stale’ files can be automated (and is 188
at HPC centers). 189
- Data that is larger than memory can handle, 190
 - can be handled by ‘big memory’ nodes. 191
 - Computing can also be done ‘in the database’ 192
- Don’t move (large data) around more than you have to - it can become 193
inefficient, and make storage slower than necessary. 194
 - New tools make it easier to find and download data combined with 195
reproducible scripts can lead to excessive and careless abuse of resources. 196
 - subset and compute on the server, in the database where possible. The 197
dplyr R package does lazy eval; SQL can perform a wide range of data 198
summaries, by groups, etc. On the other hand, it may be quicker to transfer 199
normalized (e.g. ‘flattening’ a relational database can increase the size of 200
data by orders of magnitude) 201
 - Use tools to store local ‘cached’ copies, instead of writing scripts that always 202
download archived data. Only update data if there are changes. * knitr has 203
a cache argument that saves time in re-computing and in re-downloading. 204
- For data larger than a single hard drive disk, up to multiple servers 205

- requires a meta-data server to allow fast access to distributed across many disks
- For very large data
 - it is not practical to store data
 - there are trade offs among cost, information content, and accessibility.

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

Not only data should be stored in an open format to ensure that data will be easily and widely accessible (see Rule #4), they should also be stored in a format that allows computers to make sense of it.

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.

Because a computer will be able to import your data directly (i.e., without the need for manual manipulation of your data), the script used to import and modify the data can be made available and the origin of the data used in the analysis will be evident. In turn, it will make the analysis more robust as there will be no opportunity to introduce mistakes in the data, and it will make the analysis reproducible.

To take full advantage of the data, it is important that it is structured such that the data can be manipulated and analyzed easily, in other words that the data is tidy (Wickham2014tidy). With tidy data, each variable is a column, each observation is a row, and each type of observational unit is a table. When data is organized like this, it reduces the duplication of information and it is easier to subset or summarize the dataset to include the variables or observations of interest.

To facilitate interoperability, it is best to use variable that can be 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 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.

1 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. 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
4. 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
5. Wolkovich EM, Regetz J, O'Connor ML. 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
6. 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
7. White E, Baldridge 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
8. 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
9. 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
10. 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
11. 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
12. Preston-Werner T. Semantic Versioning 2.0.0. <http://semver.org>;
13. 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
14. 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
15. 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