# Ten Simple Rules for Data Storage

Edmund Hart [1,*], Pauline Barmby [2], Jeffrey Hollister [3], David LeBauer [4],

**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**

∗ E-mail: emh@emhart.info

## Abstract

## Introduction

1

Some example text with a citation [1]

2

## Rule 1: Rule

3

# 1   Rule 2: Rule {Know your use case}

4

Researchers should know their use case and store data appropriately. Is this data    5
collected and just being archived? Will it change regularly? How will those changes be    6
logged (e.g. provenance if any)? Will this be shared via an API? Linked to a paper?    7
What are the institutional restrictions? Can you use a commercial service like    8
Dropbox or use a personally maintained system? Knowing the reason why you're    9
sharing your data will constrain your choices here.    10

## Rule 3: Rule

11

## Rule 4: Rule

12

## Rule 5: Rule

13

## Rule 6: Rule

14

## Rule 7: Rule

15

## Rule 8: Rule

16

# 2   Rule 9: Rule {Data size matters / requires special considerations}

17

18

- #39 and related GH issues #16, #19, #25    19

- Size classes:  20
  - larger than RAM  21
  - larger than HD space  22
  - larger than data storage server  23

- Storage method depends on the size of data; storage costs, transfer time, and  24
  computing costs can become substantial.  25
  - data generated by simulation and derived data should consider cost of  26
    storage vs. the cost of re-generating output.  27
  - For analyses of large data sets, the speed of reading and writing data can  28
    limt the speed of computation.  29

- Larger data sets that are actively used in analysis should be stored on a disk  30
  that is attached to a computer rather than being moved around between analysis  31
  and storage.  32
  - inactive data can be put in longer-term storage; this is less expensive, but  33
    can be slow to retrieve. Archiving of 'stale' files can be automated (and is  34
    at HPC centers).  35

- Data that is larger than memory can handle,  36
  - can be handled by 'big memory' nodes.  37
  - Computing can also be done 'in the database'  38

- Don't move (large data) around more than you have to - it can become  39
  inefficient, and make storage slower than necessary.  40
  - New tools make it easier to find and download data combined with  41
    reproducible scripts can lead to excessive and careless abuse of resources.  42
  - subset and compute on the server, in the database where possible. The  43
    dplyr R package does lazy eval; SQL can perform a wide range of data  44
    summaries, by groups, etc. On the other hand, it may be quicker to transfer  45
    normalized (e.g. 'flattening' a relational database can increase the size of  46
    data by orders of magnitude)  47
  - Use tools to store local 'cached' copies, instead of writing scripts that always  48
    download archived data. Only update data if there are changes. * knitr has  49
    a cache argument that saves time in re-computing and in re-downloading.  50

- For data larger than a single hard drive disk, up to multiple servers  51
  - requires a meta-data server to allow fast access to distributed across many  52
    disks  53

- For very large data  54
  - it is not practical to store data  55
  - there are trade offs among cost, information content, and accessibility.  56

## Rule 10: Rule  57

# 3   Acknowledgements  58

## Figure Legends

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

## Tables

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

## References

1. 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. Public Library of Science; 2014;10: e1003542.