

# DATA MANAGEMENT FOR NEUROSCIENCE: THE BIDS STANDARD AND DATALAD AN INTRODUCTION

Adina Wagner

 @AdinaKrik



---

Psychoinformatics lab,  
Institute of Neuroscience and Medicine, Brain & Behavior (INM-7)  
Research Center Jülich

Slides: <https://github.com/datalad-handbook/course/>

# **WHAT IS (RESEARCH) DATA MANAGEMENT?**



# **WHAT IS (RESEARCH) DATA MANAGEMENT?**

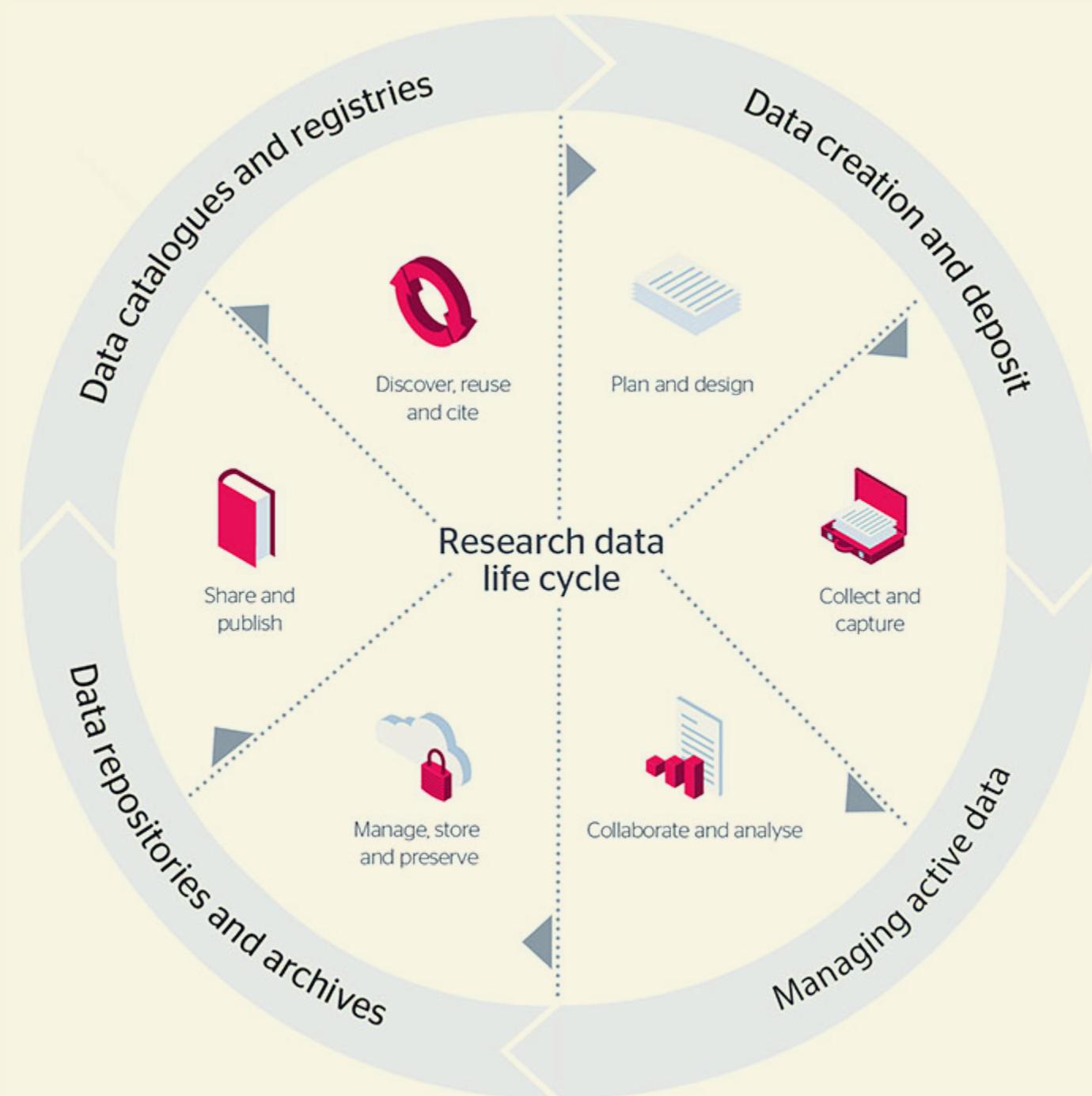
- (Research) Data = every digital object involved in your project: code, software/tools, raw data, processed data, results, manuscripts ...

# **WHAT IS (RESEARCH) DATA MANAGEMENT?**

- (Research) Data = every digital object involved in your project: code, software/tools, raw data, processed data, results, manuscripts ...
- ... needs to be properly managed - from its creation to its use, publication, sharing, archiving, re-use, or destruction... (keyword: **FAIR** data)

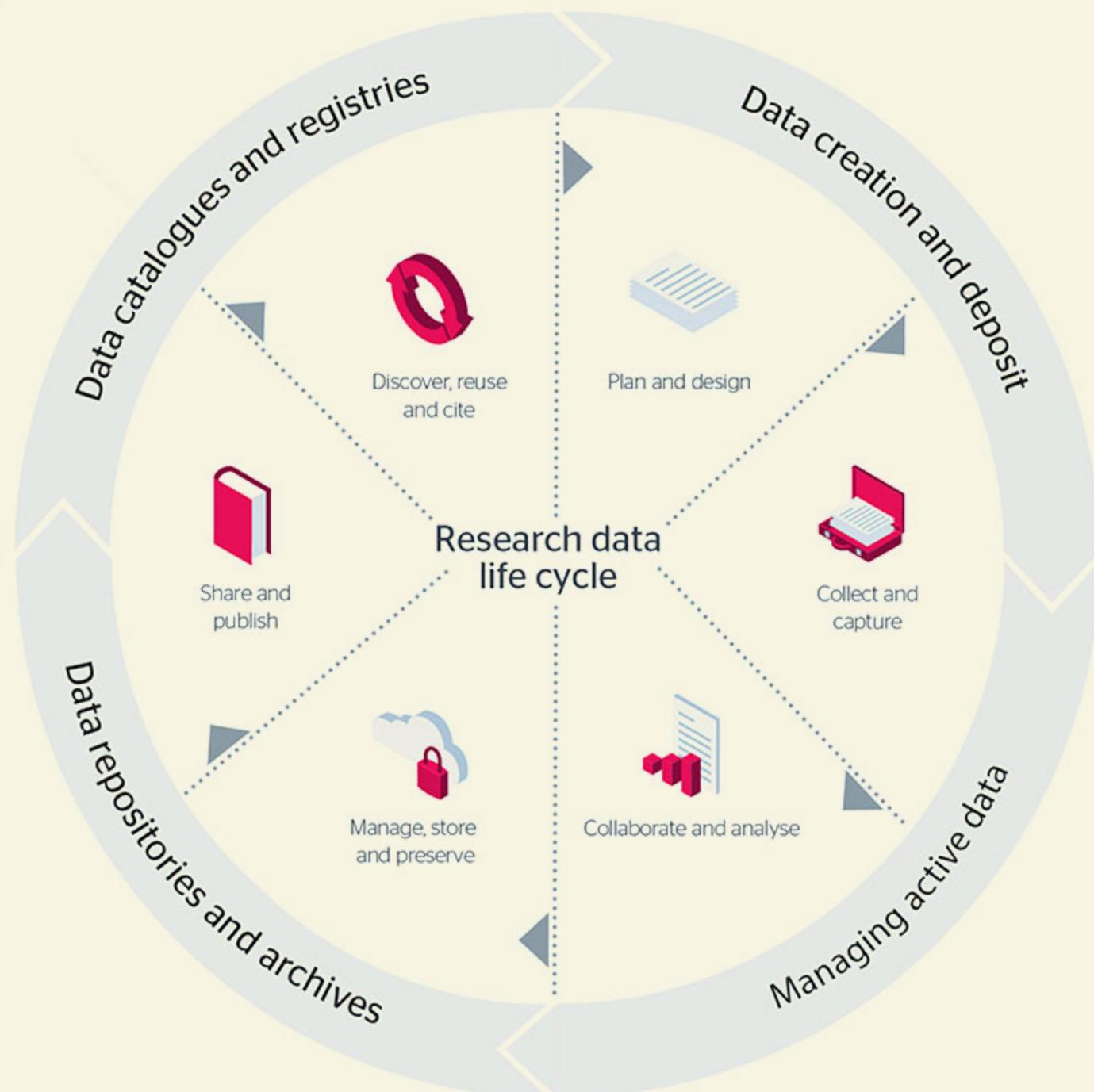
# **WHAT IS (RESEARCH) DATA MANAGEMENT?**

- (Research) Data = every digital object involved in your project: code, software/tools, raw data, processed data, results, manuscripts ...
- ... needs to be properly managed - from its creation to its use, publication, sharing, archiving, re-use, or destruction... (keyword: FAIR data)



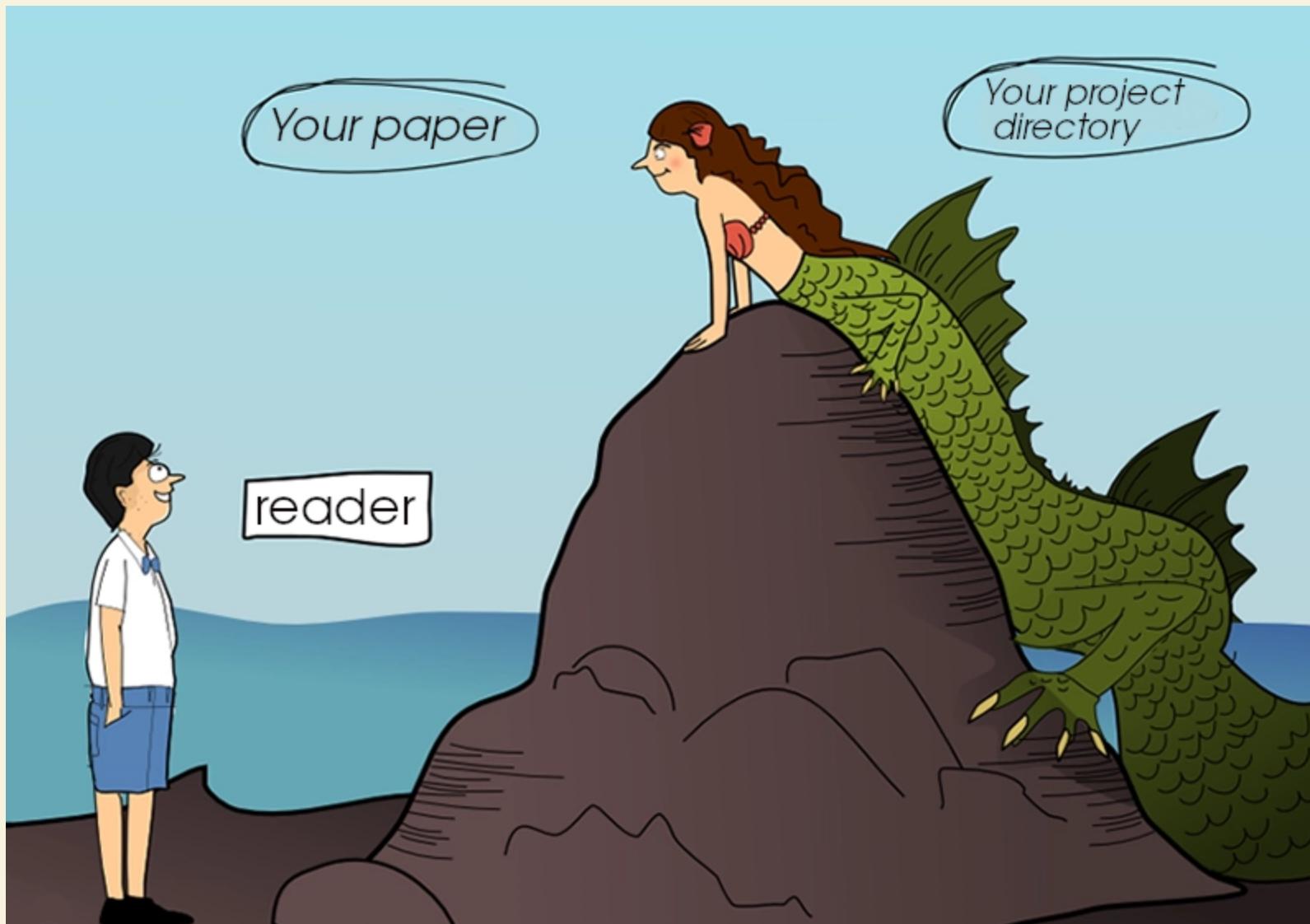
# **WHAT IS (RESEARCH) DATA MANAGEMENT?**

- (Research) Data = every digital object involved in your project: code, software/tools, raw data, processed data, results, manuscripts ...
- ... needs to be properly managed - from its creation to its use, publication, sharing, archiving, re-use, or destruction... (keyword: FAIR data)

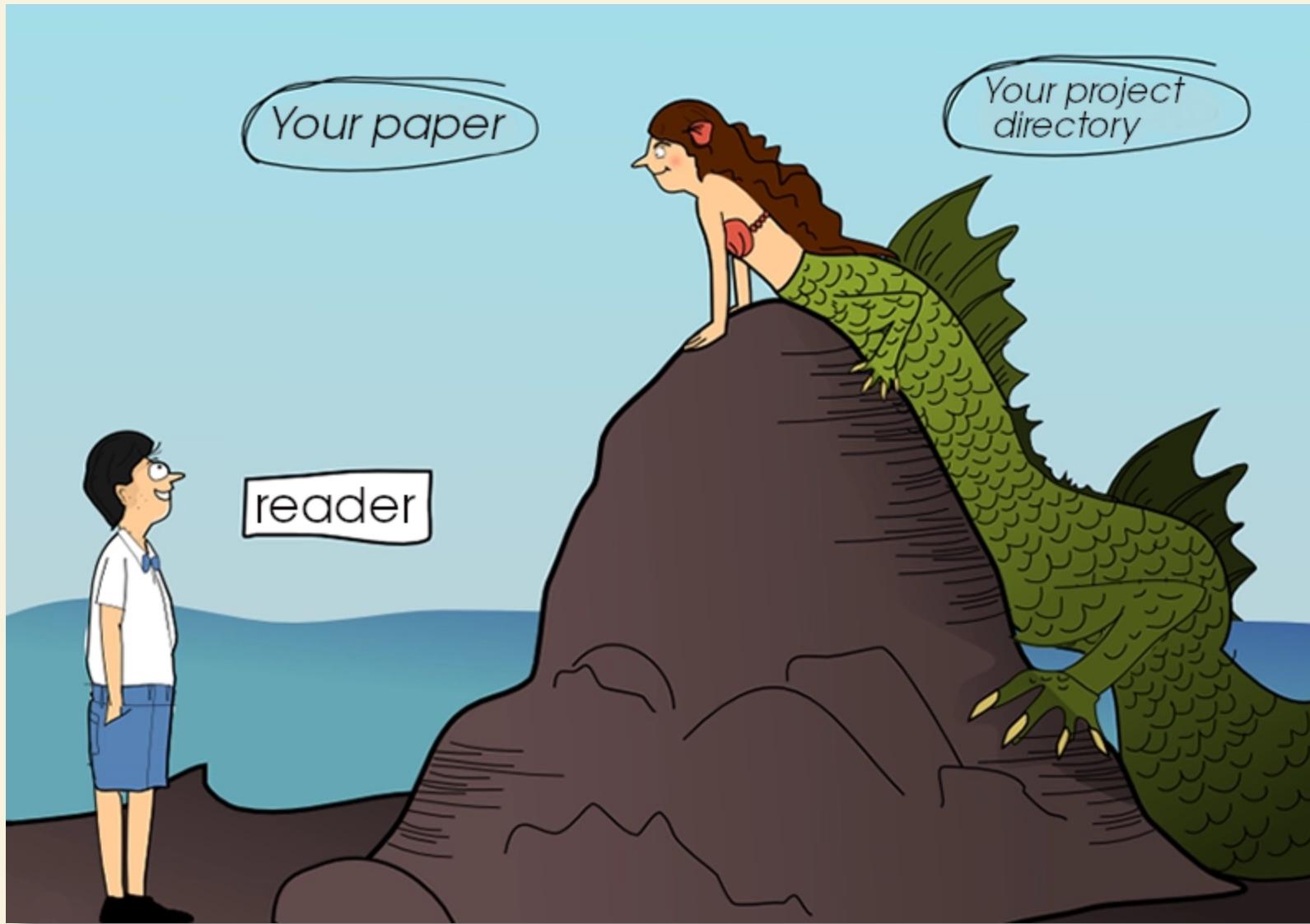


- Research data management is a key component for reproducibility, efficiency, and impact/reach of data analysis projects

# WHY DATA MANAGEMENT?

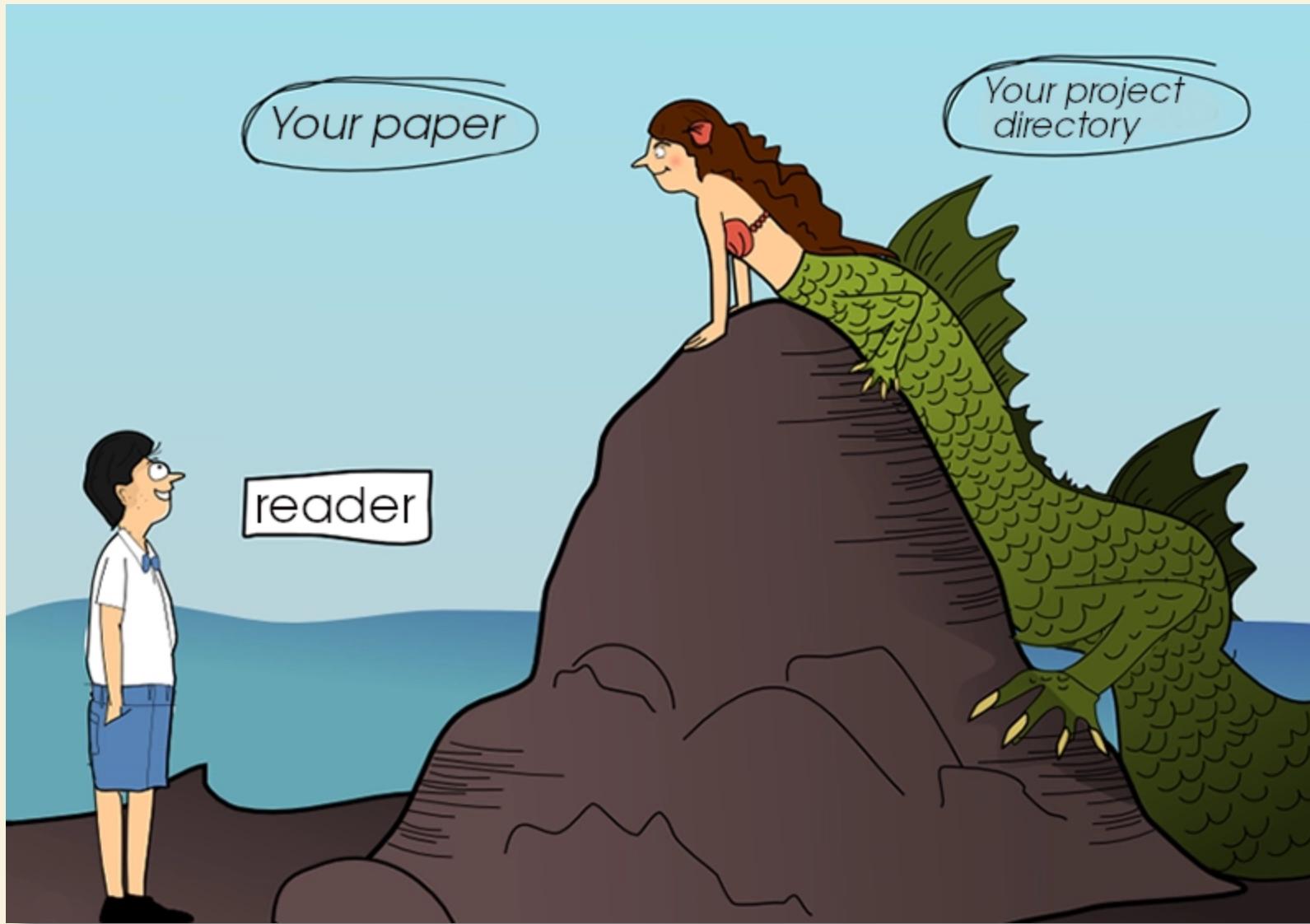


# WHY DATA MANAGEMENT?



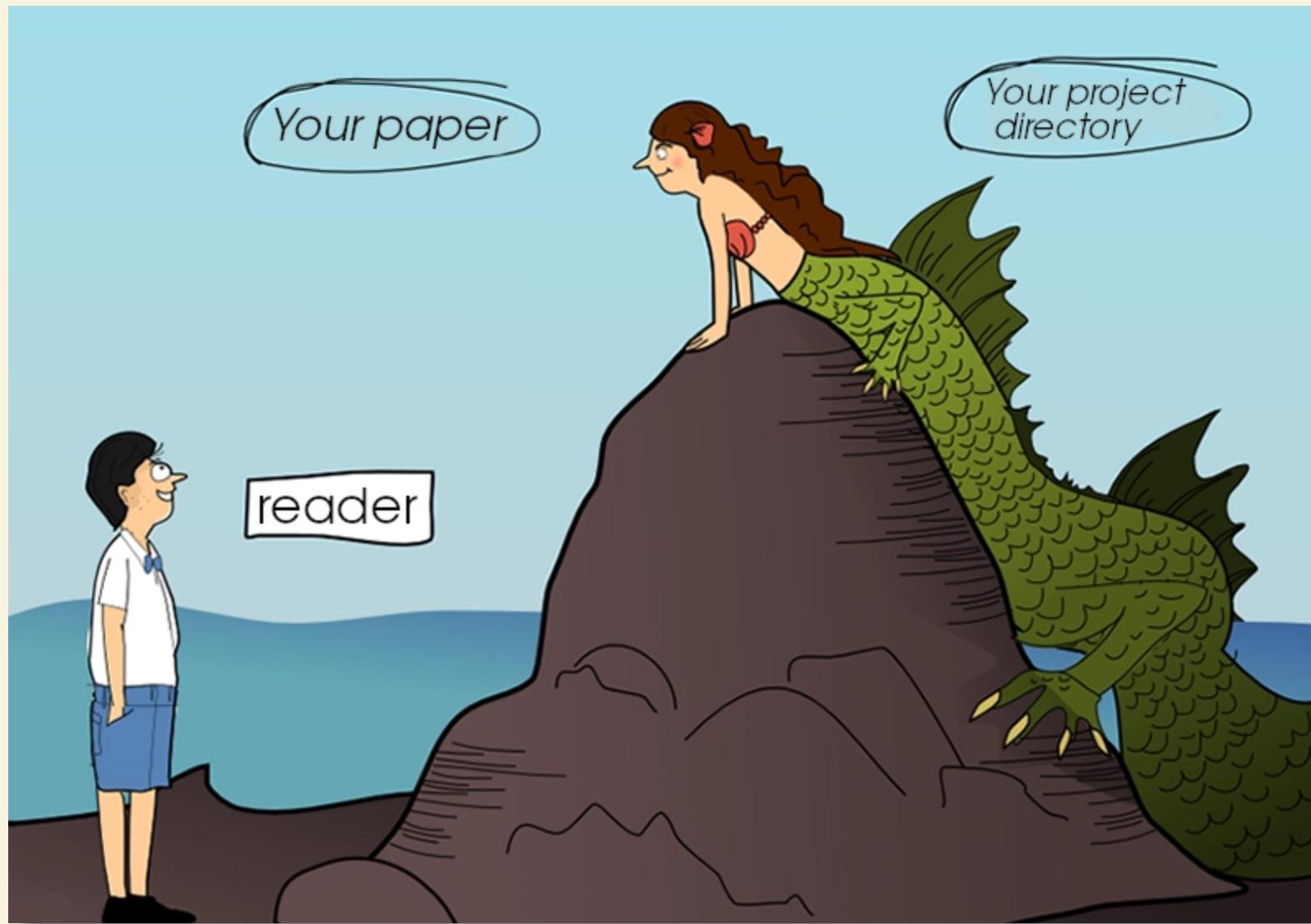
- Funders & publishers require it

# WHY DATA MANAGEMENT?



- Funders & publishers require it
- Scientific peers increasingly expect it

# WHY DATA MANAGEMENT?



- Funders & publishers require it
- Scientific peers increasingly expect it
- The quality and efficiency of your work improves

# HOW DO YOU SPEND YOUR TIME?

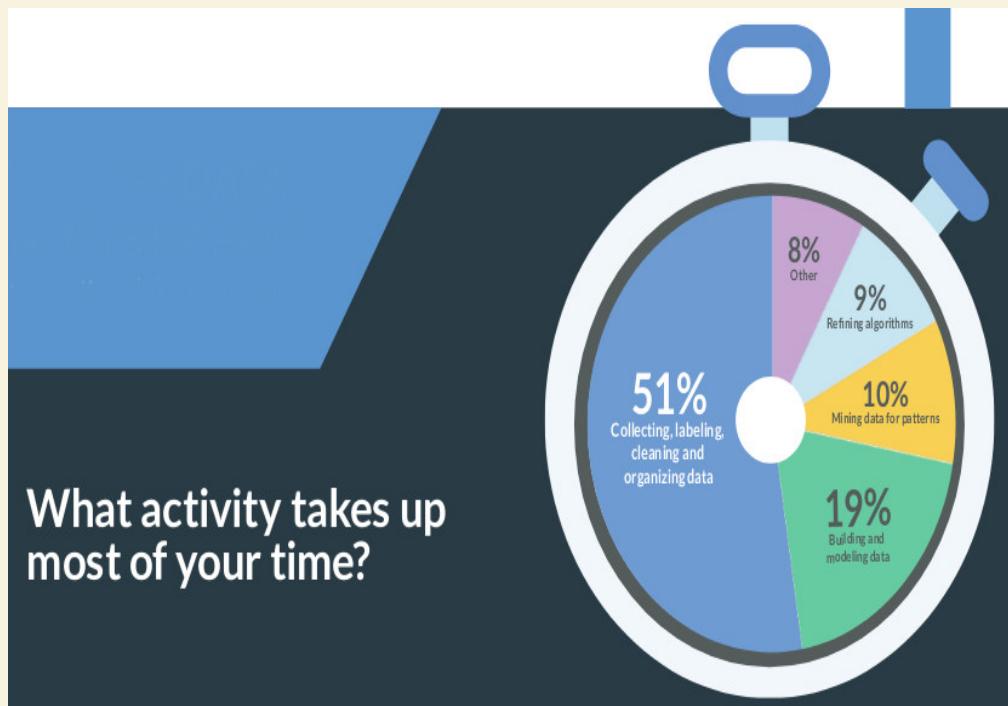
CrowdFlower  
DataScience Report 2017

---

# HOW DO YOU SPEND YOUR TIME?

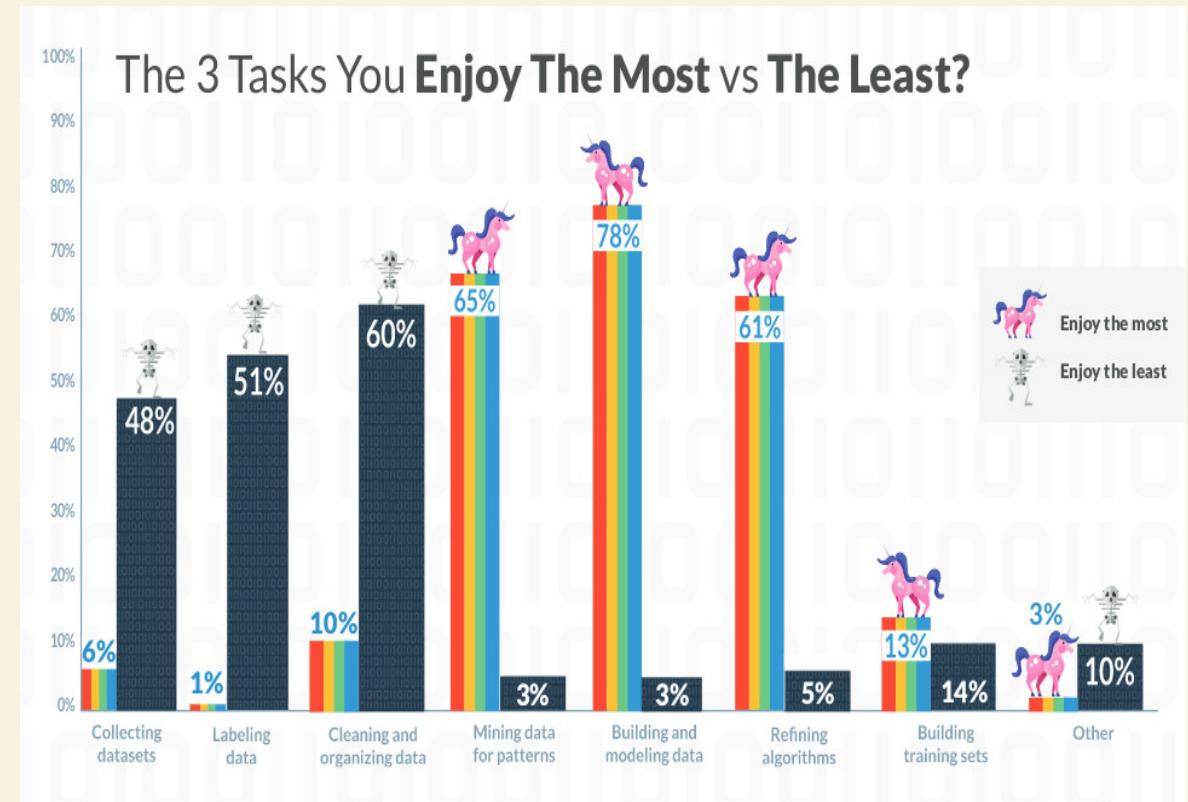
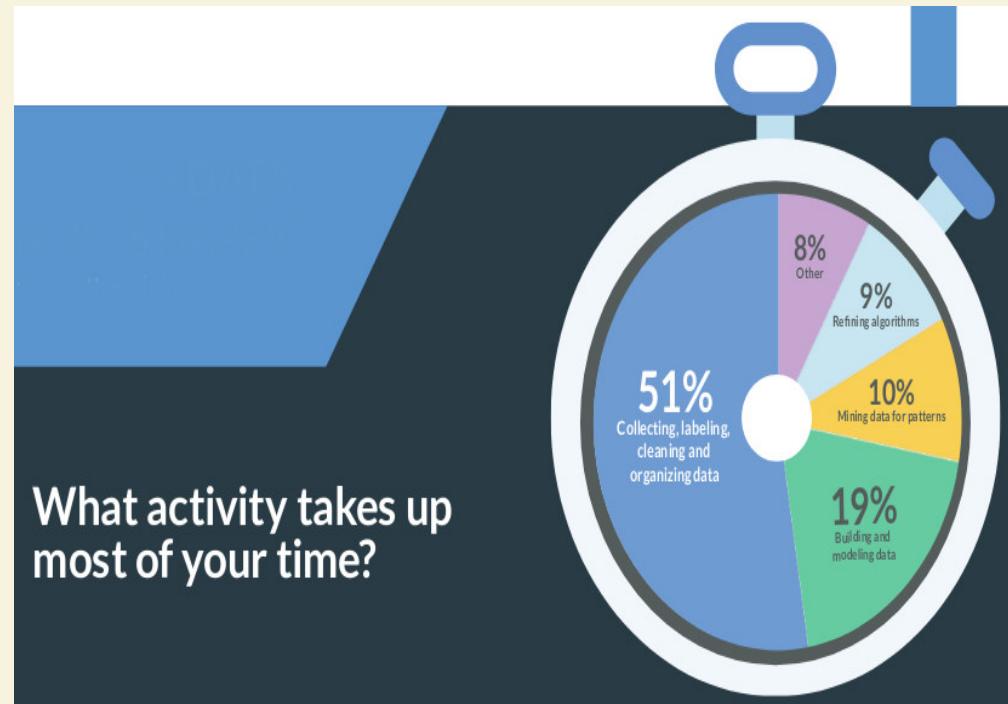
CrowdFlower  
DataScience Report 2017

---



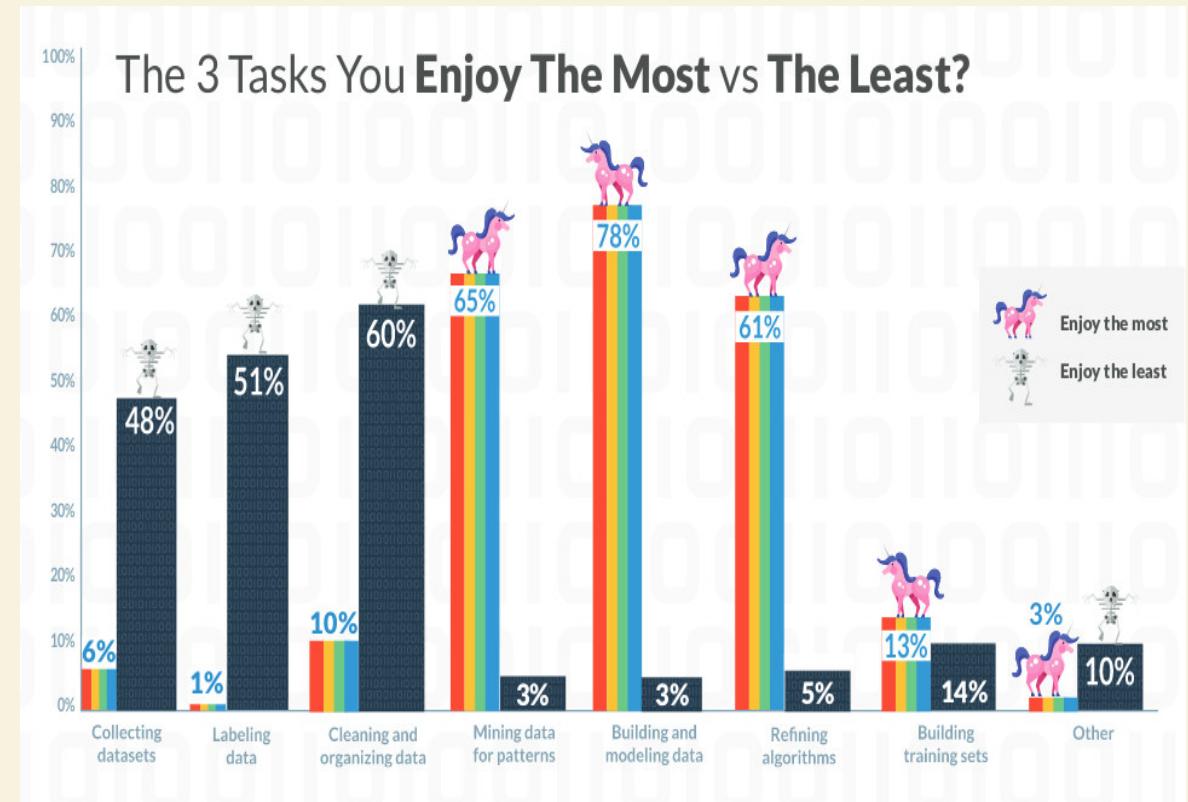
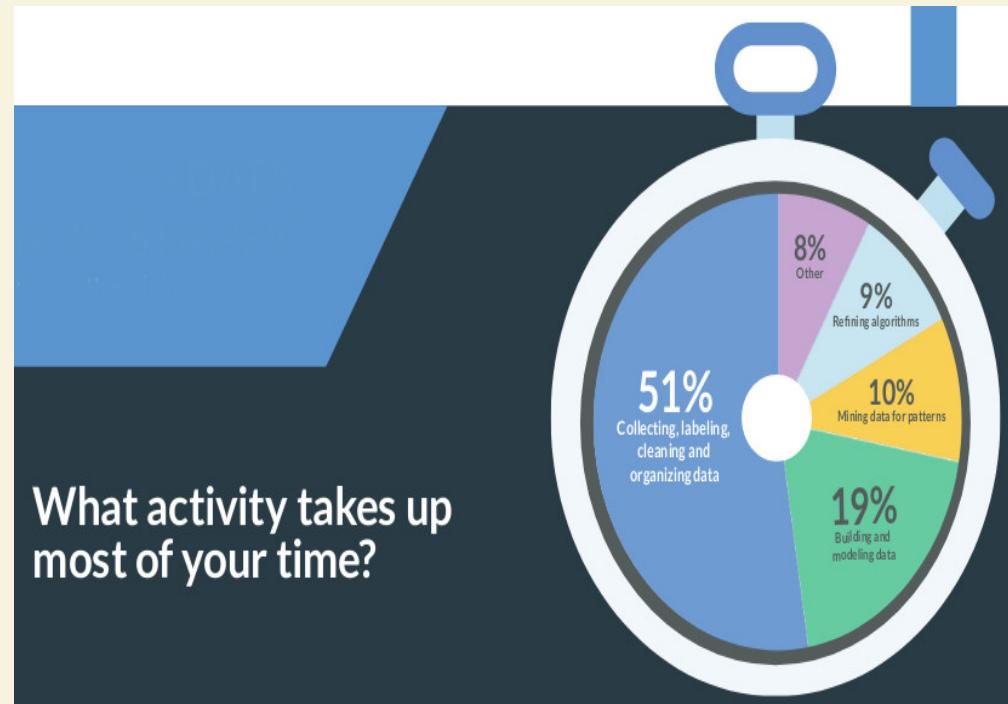
# HOW DO YOU SPEND YOUR TIME?

CrowdFlower  
DataScience Report 2017



# HOW DO YOU SPEND YOUR TIME?

CrowdFlower  
DataScience Report 2017

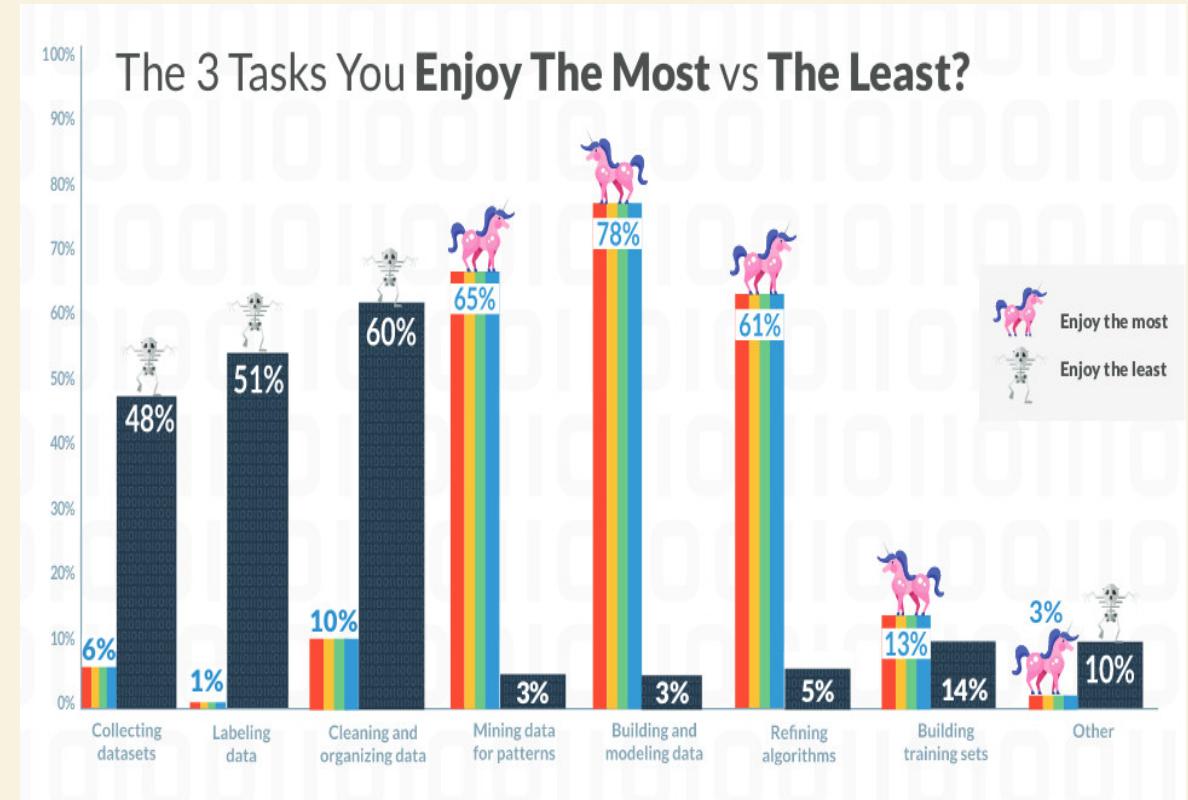
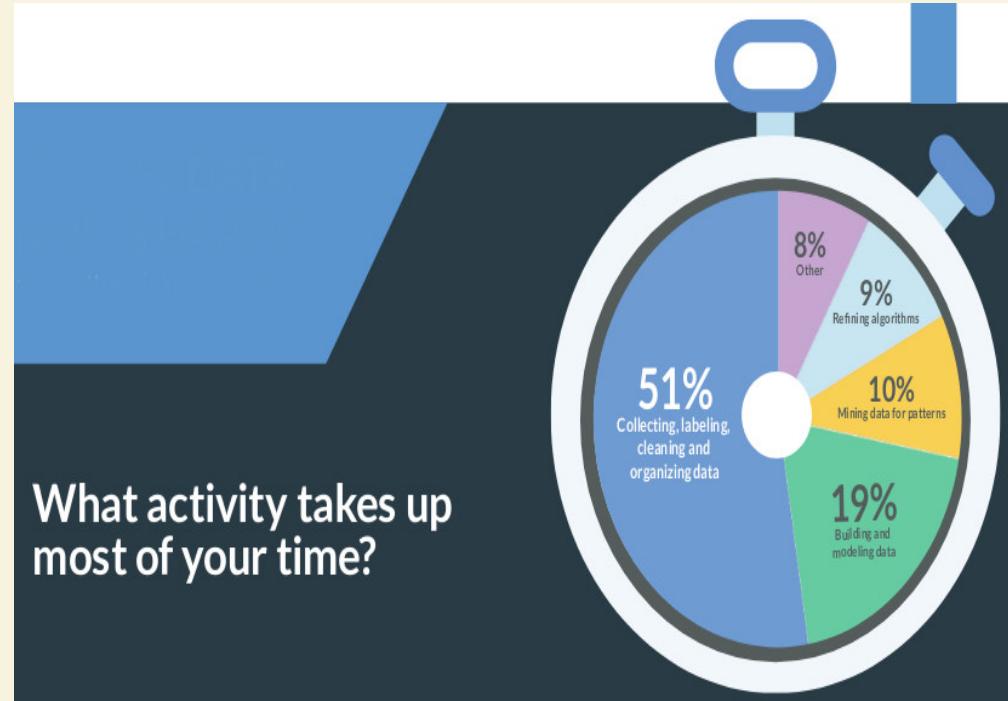


*Collaborative work and re-use of data are hampered by the effort it takes to access and understand the data.*

Thomas Wachtler

# HOW DO YOU SPEND YOUR TIME?

CrowdFlower  
DataScience Report 2017



*Collaborative work and re-use of data are hampered by the effort it takes to access and understand the data.*

Thomas Wachtler

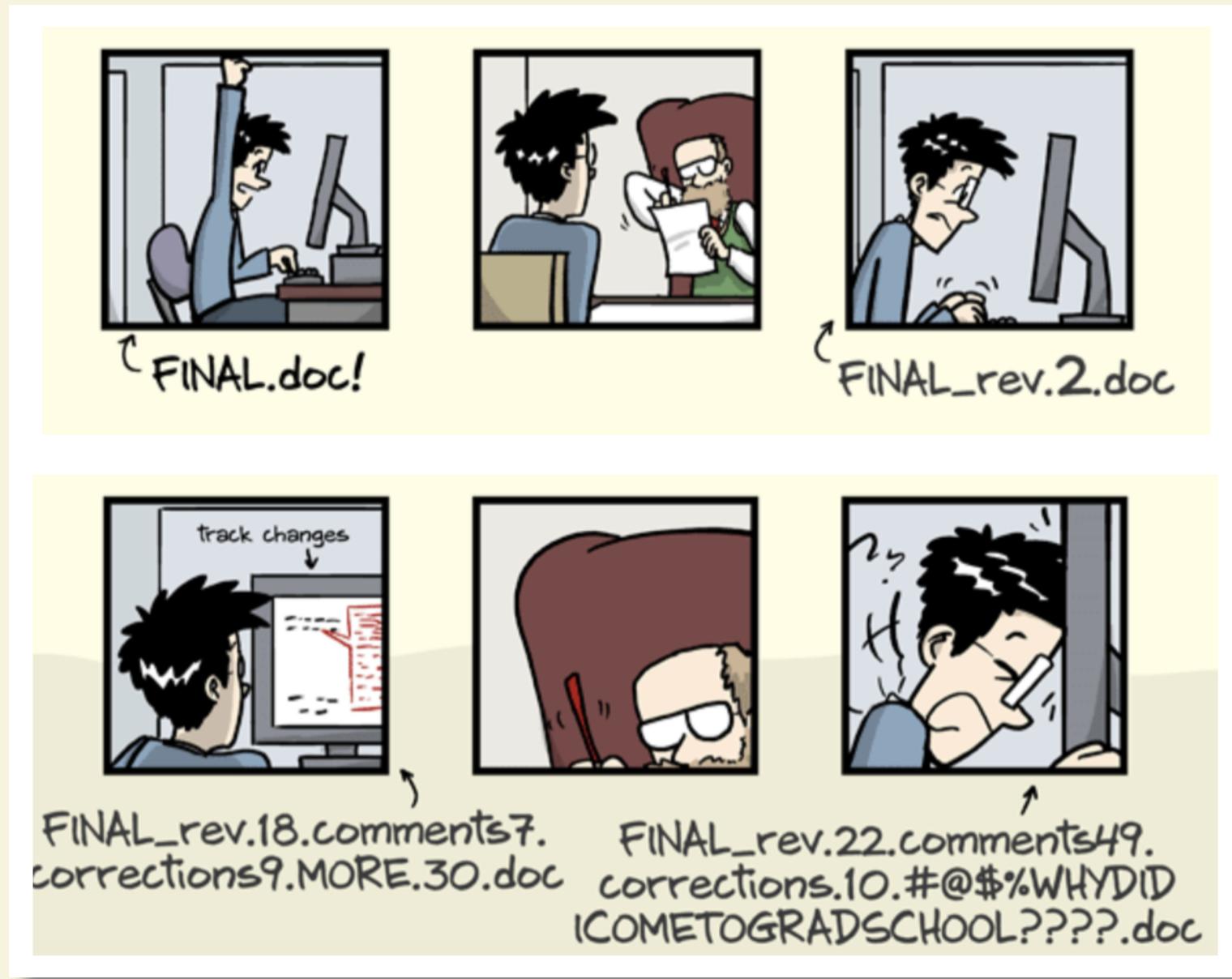
- Good data management can make your and others work & life much easier!

# HOW IS (RESEARCH) DATA MANAGEMENT POSSIBLE?

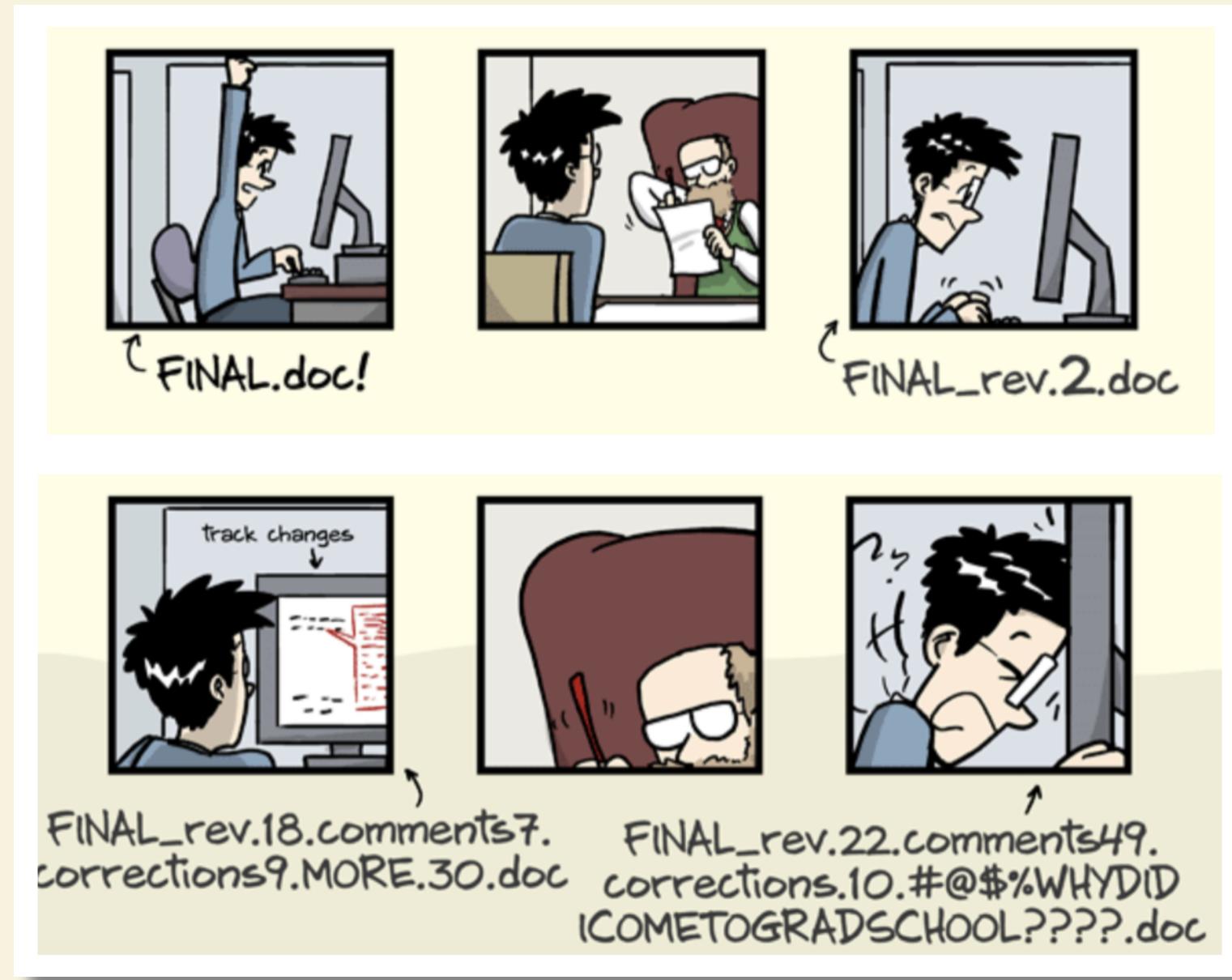
There are tools and concepts that can help:

- **Version control** your data
- **Standardize** file names and organization
- **Document everything**, ideally automatically

# WHY VERSION CONTROL?

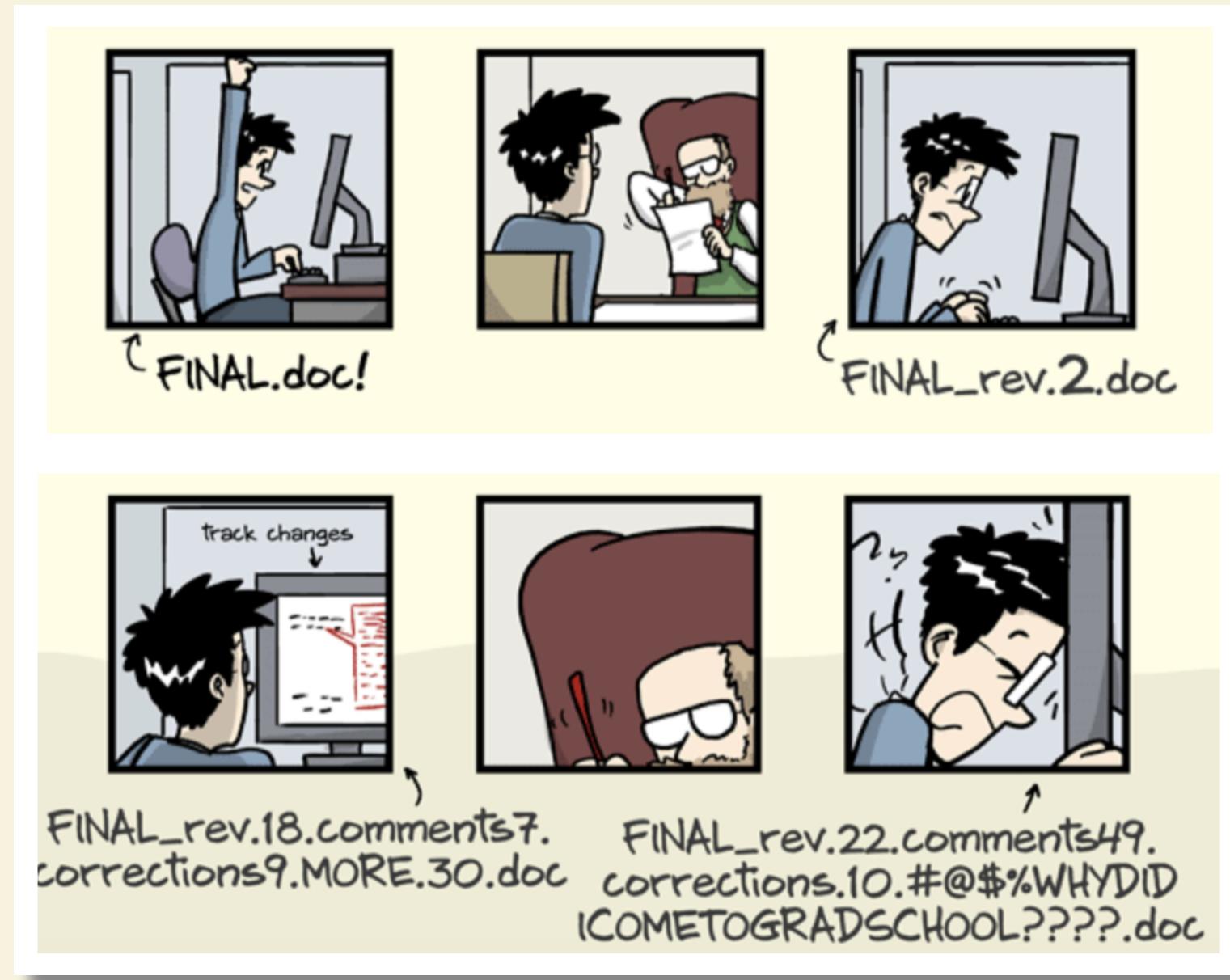


# WHY VERSION CONTROL?



- keep things organized

# WHY VERSION CONTROL?



- keep things organized
- keep track of changes

# WHY STANDARDS?

I SENT YOU THE DATA.

THANKS!

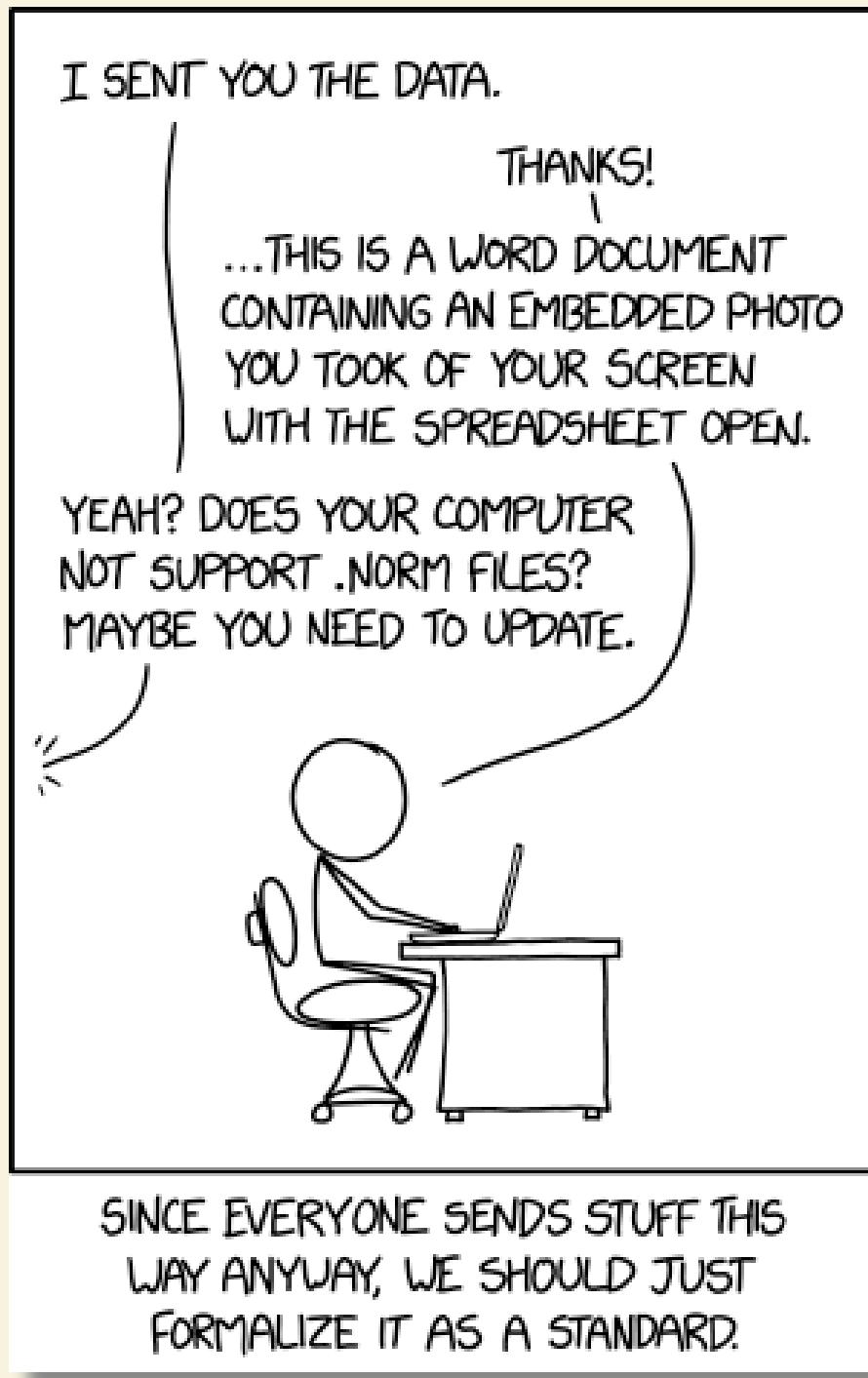
...THIS IS A WORD DOCUMENT  
CONTAINING AN EMBEDDED PHOTO  
YOU TOOK OF YOUR SCREEN  
WITH THE SPREADSHEET OPEN.

YEAH? DOES YOUR COMPUTER  
NOT SUPPORT .NORM FILES?  
MAYBE YOU NEED TO UPDATE.



SINCE EVERYONE SENDS STUFF THIS  
WAY ANYWAY, WE SHOULD JUST  
FORMALIZE IT AS A STANDARD.

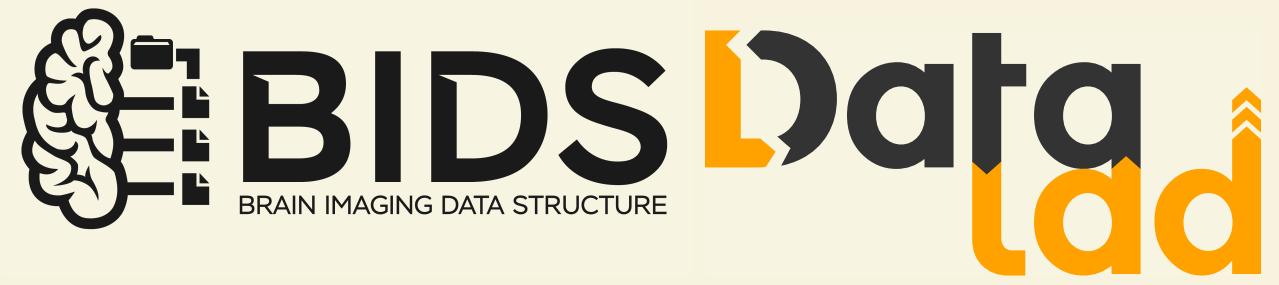
# WHY STANDARDS?



- reduce misunderstanding and rewriting/rearranging efforts



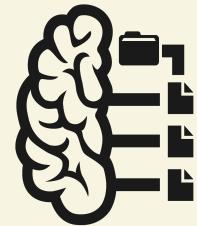




## TOOLS THAT CAN HELP WITH DATA MANAGEMENT



## TOOLS THAT CAN HELP WITH DATA MANAGEMENT



**BIDS**

BRAIN IMAGING DATA STRUCTURE



*What is it?*

## TOOLS THAT CAN HELP WITH DATA MANAGEMENT



**BIDS**

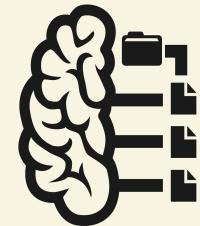
BRAIN IMAGING DATA STRUCTURE

Data  
Lad

*What is it?*

*Why should I use it?*

## TOOLS THAT CAN HELP WITH DATA MANAGEMENT



**BIDS**

BRAIN IMAGING DATA STRUCTURE



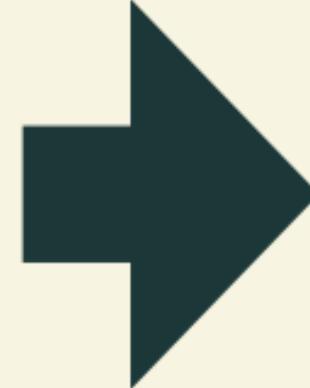
*What is it?*

*Why should I use it?*

*How can I use it?*



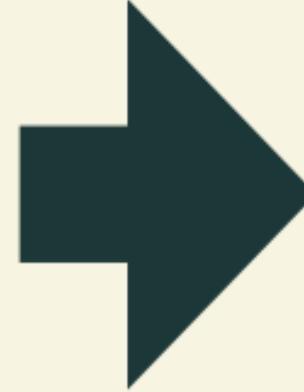
dicomdir/  
  1208200617178\_22/  
    1208200617178\_22\_8973.dcm  
    1208200617178\_22\_8943.dcm  
    1208200617178\_22\_2973.dcm  
    1208200617178\_22\_8923.dcm  
    1208200617178\_22\_4473.dcm  
    1208200617178\_22\_8783.dcm  
    1208200617178\_22\_7328.dcm  
    1208200617178\_22\_9264.dcm  
    1208200617178\_22\_9967.dcm  
    1208200617178\_22\_3894.dcm  
    1208200617178\_22\_3899.dcm  
  1208200617178\_23/  
  1208200617178\_24/  
  1208200617178\_25/



my\_dataset/  
  participants.tsv  
  sub-01/  
    anat/  
      sub-01\_T1w.nii.gz  
    func/  
      sub-01\_task-rest\_bold.nii.gz  
      sub-01\_task-rest\_bold.json  
    dwi/  
      sub-01\_dwi.nii.gz  
      sub-01\_dwi.json  
      sub-01\_dwi.bval  
      sub-01\_dwi.bvec  
  sub-02/  
  sub-03/  
  sub-04/



```
dicomdir/
  1208200617178_22/
    1208200617178_22_8973.dcm
    1208200617178_22_8943.dcm
    1208200617178_22_2973.dcm
    1208200617178_22_8923.dcm
    1208200617178_22_4473.dcm
    1208200617178_22_8783.dcm
    1208200617178_22_7328.dcm
    1208200617178_22_9264.dcm
    1208200617178_22_9967.dcm
    1208200617178_22_3894.dcm
    1208200617178_22_3899.dcm
  1208200617178_23/
  1208200617178_24/
  1208200617178_25/
```



```
my_dataset/
  participants.tsv
  sub-01/
    anat/
      sub-01_T1w.nii.gz
    func/
      sub-01_task-rest_bold.nii.gz
      sub-01_task-rest_bold.json
    dwi/
      sub-01_dwi.nii.gz
      sub-01_dwi.json
      sub-01_dwi.bval
      sub-01_dwi.bvec
  sub-02/
  sub-03/
  sub-04/
```

BIDS is a standard for a multitude of neuroimaging data (MRI, EEG, ...). It defines a data organization, naming schemes for files, and meta data descriptors.





**BIDS is a structure**



**BIDS is a structure**

BIDS is not a new file format



**BIDS is a structure**

BIDS is not a new file format

**BIDS is a standard**



## **BIDS is a structure**

BIDS is not a new file format

## **BIDS is a standard**

BIDS is not a software tool. However, there is a large and growing amount of tools that are compatible with it that ease data archiving, data discovery/search, and analysis



## **BIDS is a structure**

BIDS is not a new file format

## **BIDS is a standard**

BIDS is not a software tool. However, there is a large and growing amount of tools that are compatible with it that ease data archiving, data discovery/search, and analysis

## **There is no "one" BIDS**



## **BIDS is a structure**

BIDS is not a new file format

## **BIDS is a standard**

BIDS is not a software tool. However, there is a large and growing amount of tools that are compatible with it that ease data archiving, data discovery/search, and analysis

## **There is no "one" BIDS**

BIDS exist for (a growing amount of) different modalities. There is constant (open!) development of all of them





**BIDS helps you and others to intuitively understand your data**



**BIDS helps you and others to intuitively understand your data**

A consensus on data organization spares you and others time to dive in to or rearrange data or scripts



**BIDS helps you and others to intuitively understand your data**

A consensus on data organization spares you and others time to dive in to or rearrange data or scripts

**BIDS opens up a large range of tools for you:**



**BIDS helps you and others to intuitively understand your data**

A consensus on data organization spares you and others time to dive in to or rearrange data or scripts

**BIDS opens up a large range of tools for you:**

➲ BIDS Apps, e.g., fmriprep or MRIQC



## **BIDS helps you and others to intuitively understand your data**

A consensus on data organization spares you and others time to dive in to or rearrange data or scripts

## **BIDS opens up a large range of tools for you:**

- ➲ BIDS Apps, e.g., fmriprep or MRIQC
- ➲ BIDS-aware tooling, e.g., PyBIDS



## **BIDS helps you and others to intuitively understand your data**

A consensus on data organization spares you and others time to dive in to or rearrange data or scripts

## **BIDS opens up a large range of tools for you:**

- ⌚ BIDS Apps, e.g., fmriprep or MRIQC
- ⌚ BIDS-aware tooling, e.g., PyBIDS
- ⌚ brainlife.io



## BIDS helps you and others to intuitively understand your data

A consensus on data organization spares you and others time to dive in to or rearrange data or scripts

## BIDS opens up a large range of tools for you:

- ⌚ BIDS Apps, e.g., fmriprep or MRIQC
- ⌚ BIDS-aware tooling, e.g., PyBIDS
- ⌚ brainlife.io
- ⌚ Upload your (or download others) BIDS-compliant datasets from OpenNeuro



## BIDS helps you and others to intuitively understand your data

A consensus on data organization spares you and others time to dive in to or rearrange data or scripts

## BIDS opens up a large range of tools for you:

- ⌚ BIDS Apps, e.g., fmriprep or MRIQC
- ⌚ BIDS-aware tooling, e.g., PyBIDS
- ⌚ brainlife.io
- ⌚ Upload your (or download others) BIDS-compliant datasets from OpenNeuro



## Useful links and pointers

- ➲ Read the [paper](#)
- ➲ Get started with the [BIDS starter-kit](#)
- ➲ Work through the Stanford Center for Reproducible Neuroscience [BIDS Tutorial Series](#)
- ➲ Use the [BIDS validator](#) to check your datasets
- ➲ Get involved on [Github](#) to help shape BIDS to your needs. You can also checkout the [Google discussion group](#)
- ➲ Follow [@BIDS-standard](#)

# Data



DataLad is a data management multitool.



Let's see it in action

# **REPRODUCIBLE PAPER - A MAGIC TRICK?**

If curious, you can read up all the details and a step-by-step instruction [here](#).

# Data Lad IN BRIEF

- A command-line tool, available for all major operating systems (Linux, macOS/OSX, Windows)
- Build on top of Git and Git-annex
- Allows...
  - ... version-controlling arbitrarily large content,
  - ... easily sharing and obtaining data (note: no data hosting!),
  - ... (computationally) reproducible data analysis,
  - ... and *much* more
- Completely domain-agnostic

# Data Lad IN BRIEF

- A command-line tool, available for all major operating systems (Linux, macOS/OSX, Windows)
- Build on top of Git and Git-annex
- Allows...
  - ... version-controlling arbitrarily large content,
  - ... easily sharing and obtaining data (note: no data hosting!),
  - ... (computationally) reproducible data analysis,
  - ... and *much* more
- Completely domain-agnostic

Today: Basic concepts and commands.

# DataLad IN BRIEF

- A command-line tool, available for all major operating systems (Linux, macOS/OSX, Windows)
- Build on top of Git and Git-annex
- Allows...
  - ... version-controlling arbitrarily large content,
  - ... easily sharing and obtaining data (note: no data hosting!),
  - ... (computationally) reproducible data analysis,
  - ... and *much* more
- Completely domain-agnostic

Today: Basic concepts and commands.

☞ For more: Read the DataLad Handbook



# DATALAD DATASETS

- DataLad's core data structure

# DATALAD DATASETS

- DataLad's core data structure
  - Dataset = A directory managed by DataLad

# DATALAD DATASETS

- DataLad's core data structure
  - Dataset = A directory managed by DataLad
  - Any directory of your computer can be managed by DataLad.

# DATALAD DATASETS

- DataLad's core data structure
  - Dataset = A directory managed by DataLad
  - Any directory of your computer can be managed by DataLad.
  - Datasets can be *created* (from scratch) or *installed*

# DATALAD DATASETS

- DataLad's core data structure
  - Dataset = A directory managed by DataLad
  - Any directory of your computer can be managed by DataLad.
  - Datasets can be *created* (from scratch) or *installed*
  - Datasets can be nested: *linked subdirectories*

# **EXPERIENCE A DATALAD DATASET**

Code to follow along:

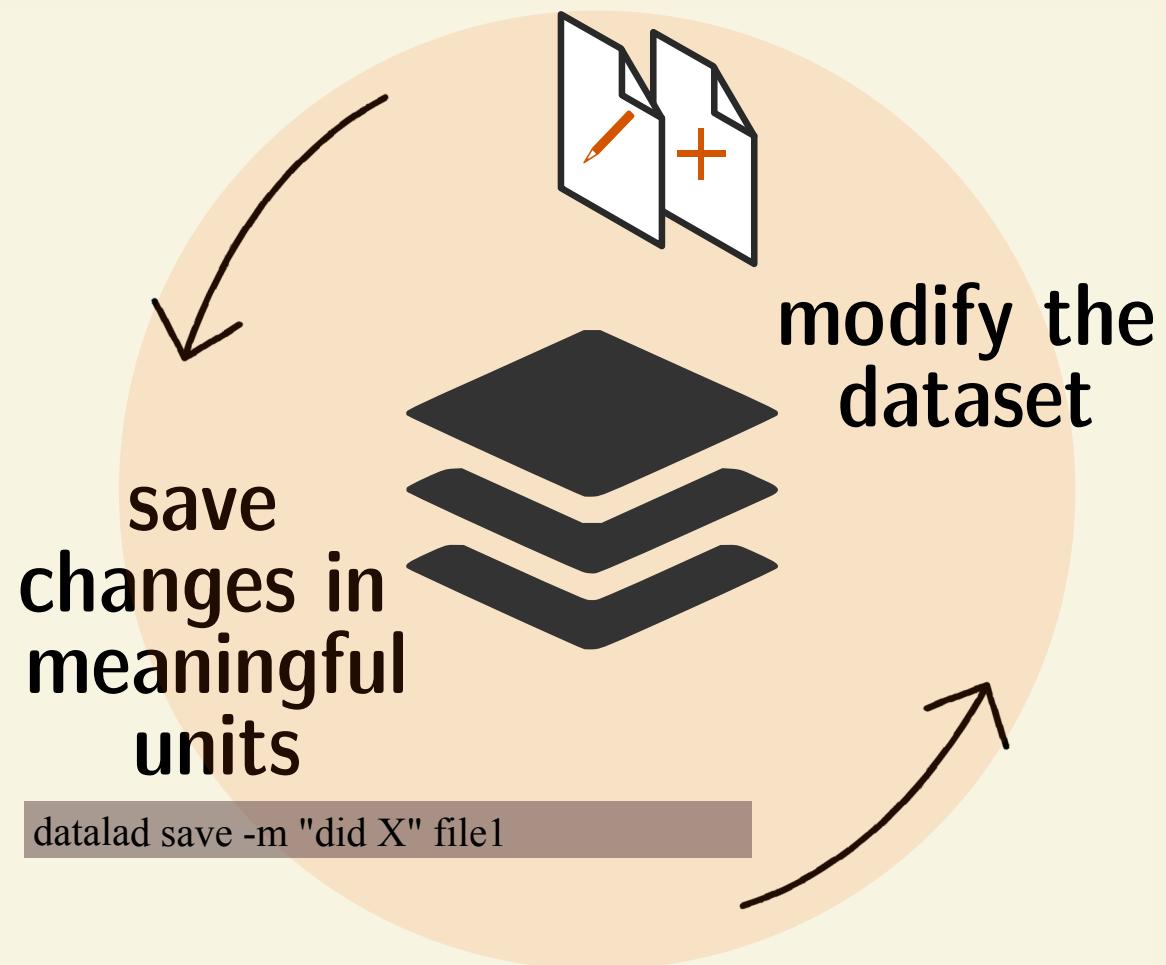
[http://handbook.datalad.org/en/latest/code\\_from\\_chapters/01\\_dataset\\_basics\\_code.html](http://handbook.datalad.org/en/latest/code_from_chapters/01_dataset_basics_code.html)

# **LOCAL VERSION CONTROL**

Procedurally, version control is easy with DataLad!

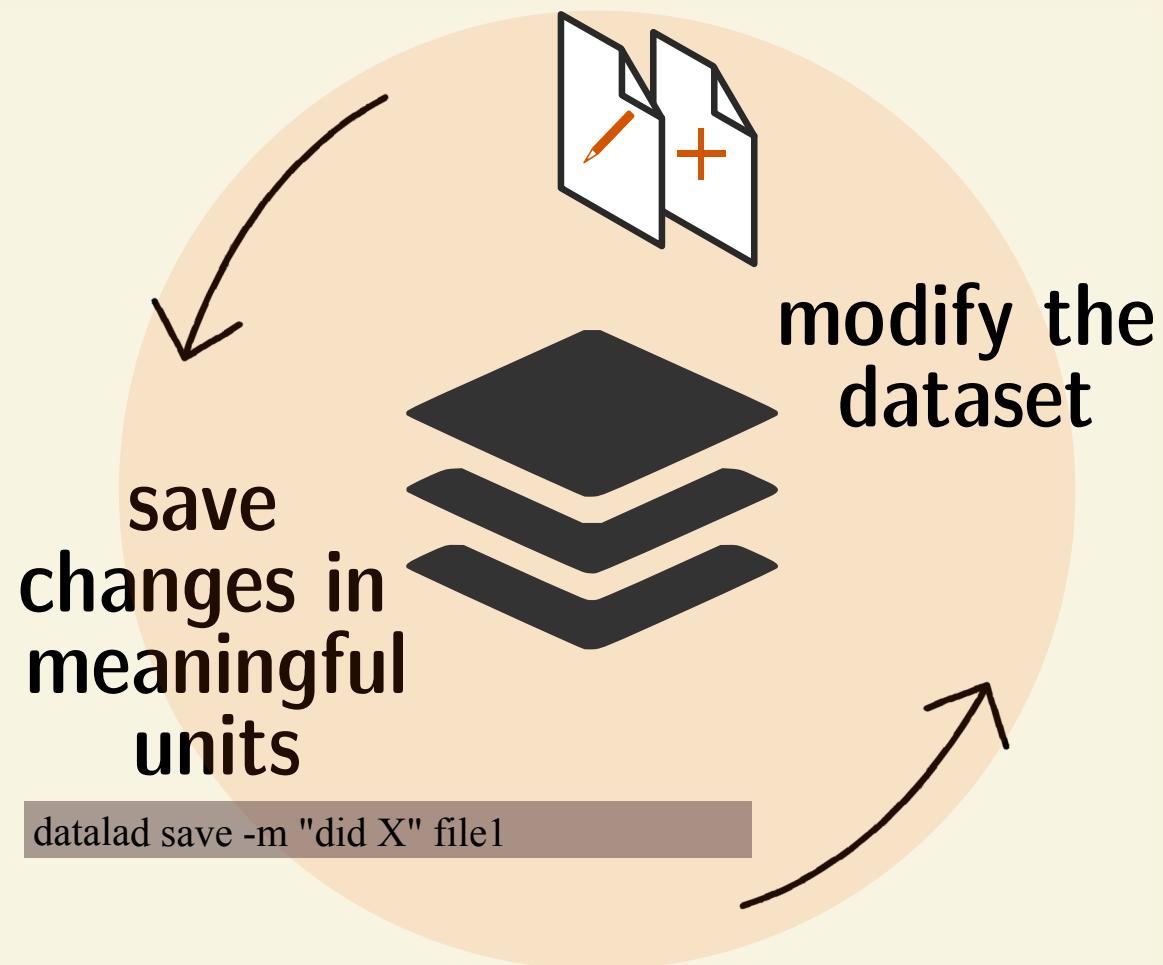
# LOCAL VERSION CONTROL

Procedurally, version control is easy with DataLad!



# LOCAL VERSION CONTROL

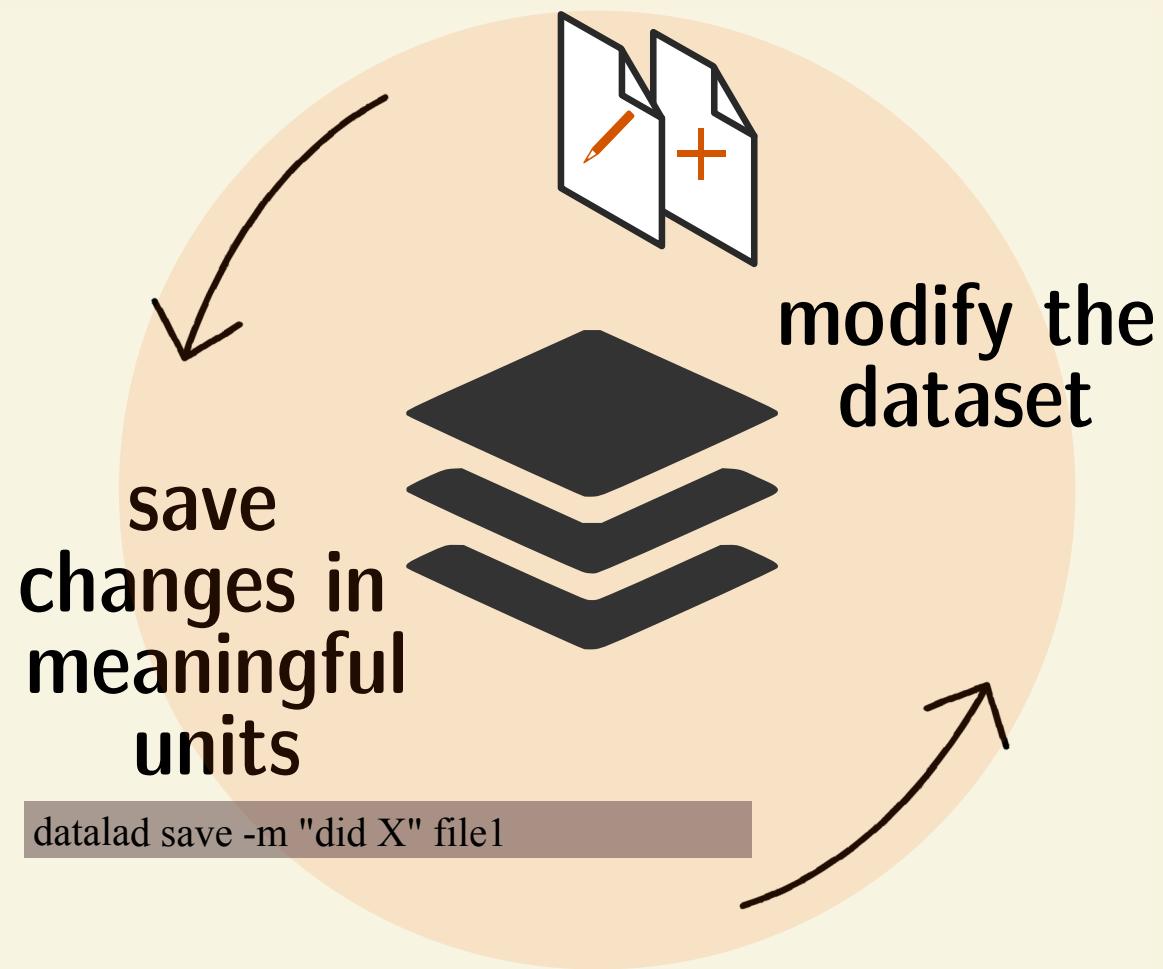
Procedurally, version control is easy with DataLad!



Advice:

# LOCAL VERSION CONTROL

Procedurally, version control is easy with DataLad!

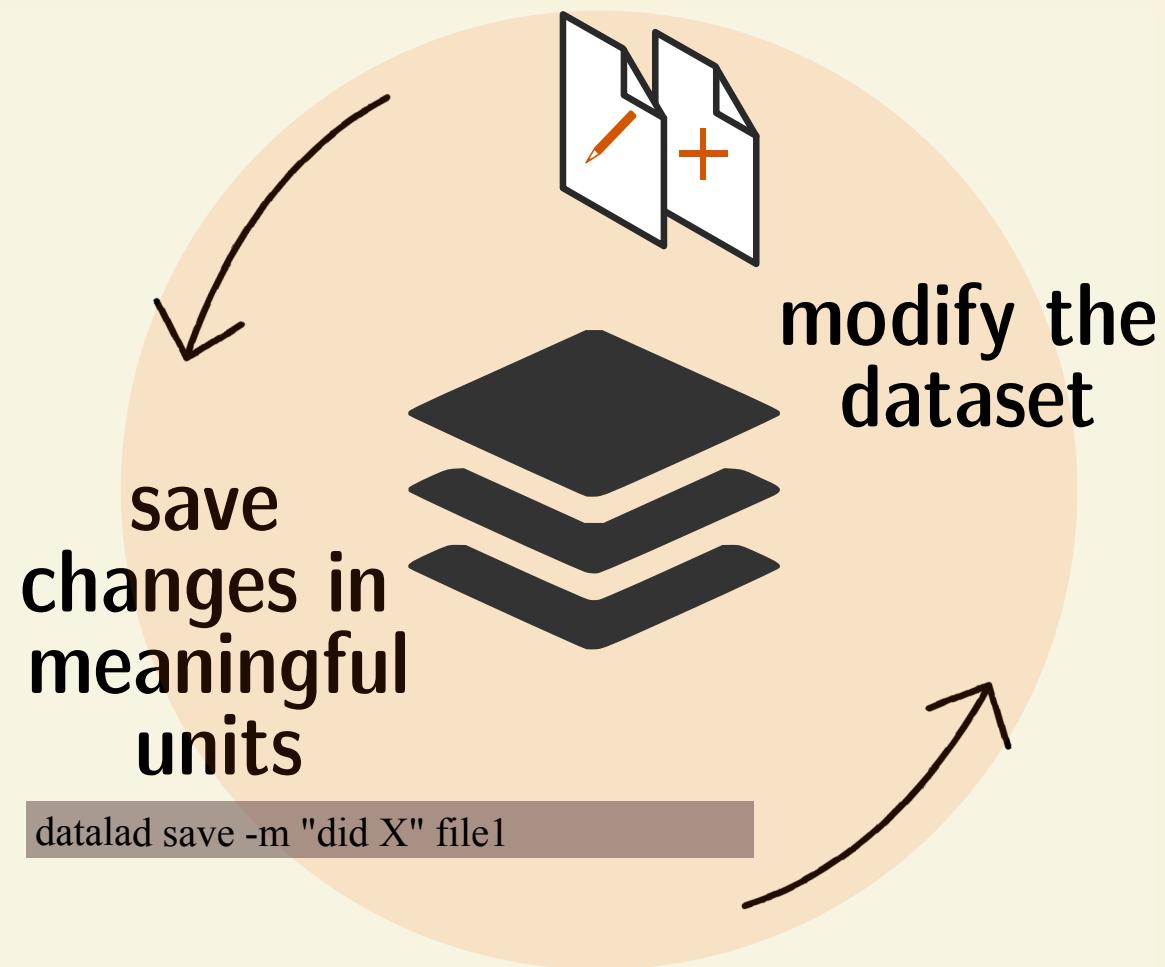


- Save *meaningful* units of change

Advice:

# LOCAL VERSION CONTROL

Procedurally, version control is easy with DataLad!



- Advice:
- Save *meaningful* units of change
  - Attach helpful commit messages

# **SUMMARY - LOCAL VERSION CONTROL**

## **SUMMARY - LOCAL VERSION CONTROL**

**datalad create** creates an empty dataset.

## **SUMMARY - LOCAL VERSION CONTROL**

**datalad create** creates an empty dataset.

Configurations (-c yoda, -c text2git) are useful.

## SUMMARY - LOCAL VERSION CONTROL

**datalad create** creates an empty dataset.

Configurations (-c yoda, -c text2git) are useful.

A dataset has a *history* to track files and their modifications.

## SUMMARY - LOCAL VERSION CONTROL

**datalad create** creates an empty dataset.

Configurations (-c yoda, -c text2git) are useful.

A dataset has a *history* to track files and their modifications.

Explore it with Git (**git log**) or external tools (e.g., **tig**).

## SUMMARY - LOCAL VERSION CONTROL

**datalad create** creates an empty dataset.

Configurations (-c yoda, -c text2git) are useful.

A dataset has a *history* to track files and their modifications.

Explore it with Git (`git log`) or external tools (e.g., `tig`).

**datalad save** records the dataset or file state to the history.

## SUMMARY - LOCAL VERSION CONTROL

**datalad create** creates an empty dataset.

Configurations (-c yoda, -c text2git) are useful.

A dataset has a *history* to track files and their modifications.

Explore it with Git (`git log`) or external tools (e.g., `tig`).

**datalad save** records the dataset or file state to the history.

Concise commit messages should summarize the change for future you and others.

## SUMMARY - LOCAL VERSION CONTROL

**datalad create** creates an empty dataset.

Configurations (-c yoda, -c text2git) are useful.

A dataset has a *history* to track files and their modifications.

Explore it with Git (`git log`) or external tools (e.g., `tig`).

**datalad save** records the dataset or file state to the history.

Concise commit messages should summarize the change for future you and others.

**datalad status** reports the current state of the dataset.

## SUMMARY - LOCAL VERSION CONTROL

**datalad create** creates an empty dataset.

Configurations (-c yoda, -c text2git) are useful.

A dataset has a *history* to track files and their modifications.

Explore it with Git (`git log`) or external tools (e.g., `tig`).

**datalad save** records the dataset or file state to the history.

Concise commit messages should summarize the change for future you and others.

**datalad status** reports the current state of the dataset.

A clean dataset status is good practice.

# FROM HERE

"FINAL".doc



FINAL.doc!



FINAL\_rev.2.doc



FINAL\_rev.6.COMMENTS.doc



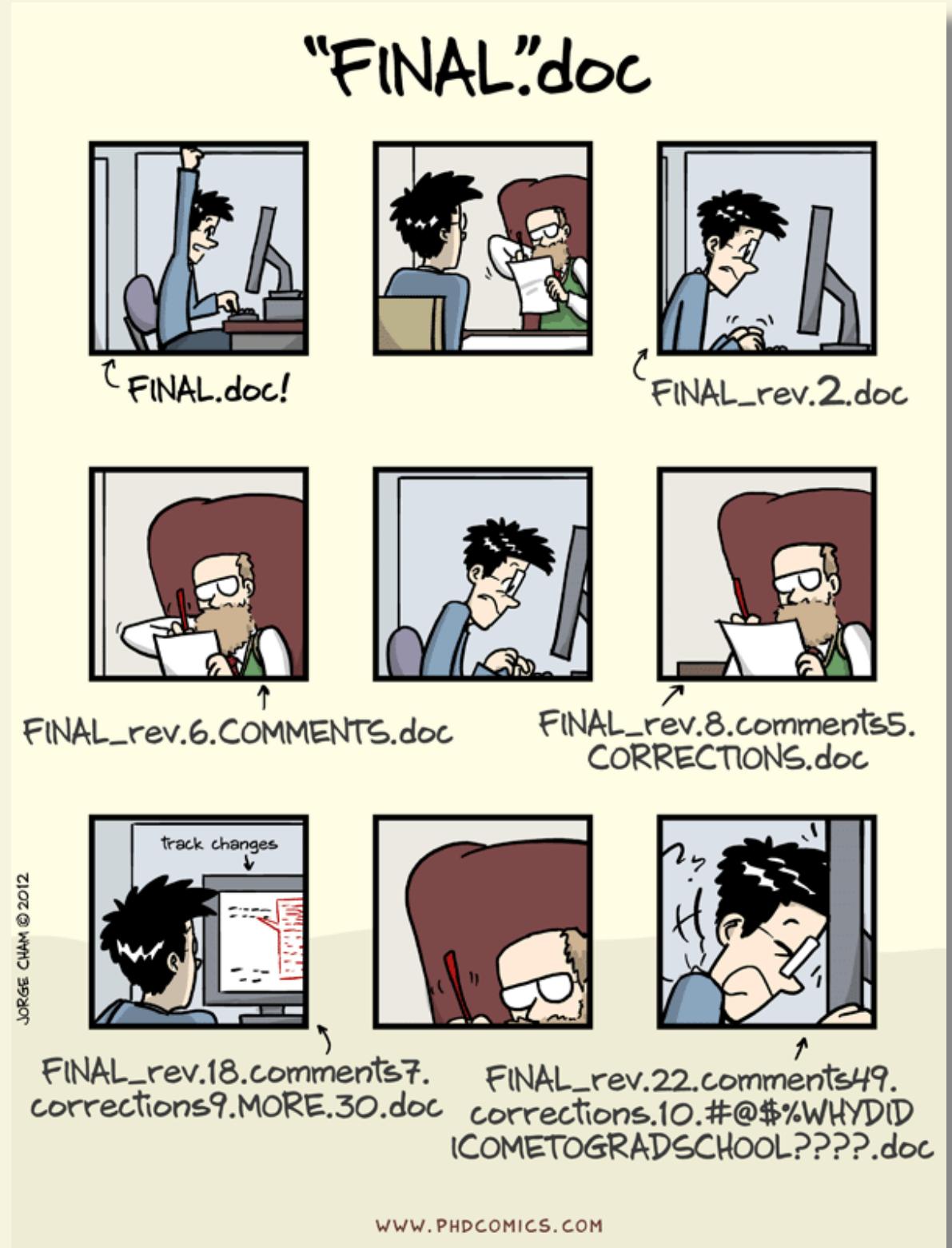
FINAL\_rev.8.comments5.  
CORRECTIONS.doc



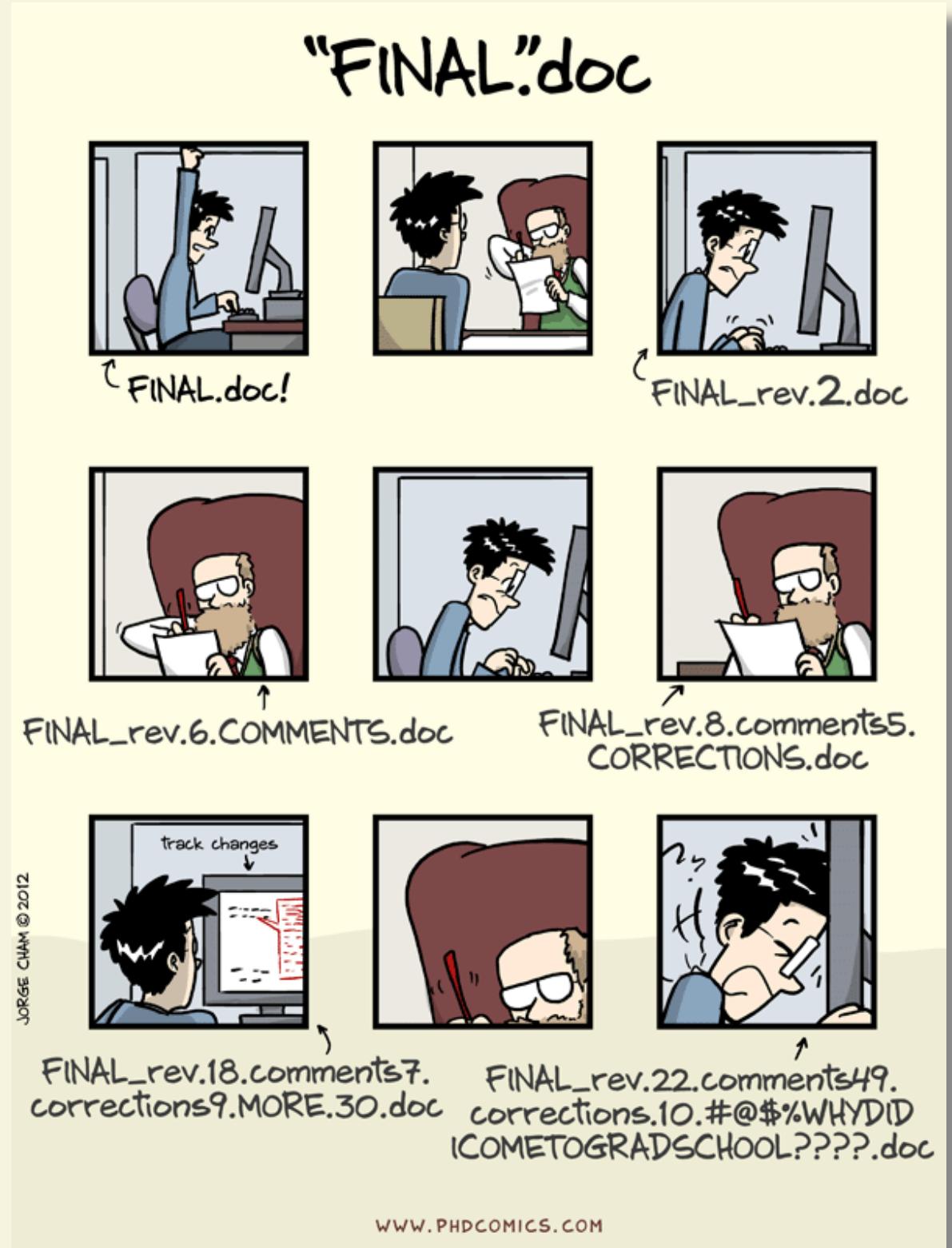
JORGE CHAM © 2012

FINAL\_rev.18.comments7.  
corrections9.MORE.30.doc      FINAL\_rev.22.comments49.  
corrections.10.#@\$%WHYDID  
ICOMETOGRAD SCHOOL????.doc

# FROM HERE TO THIS:



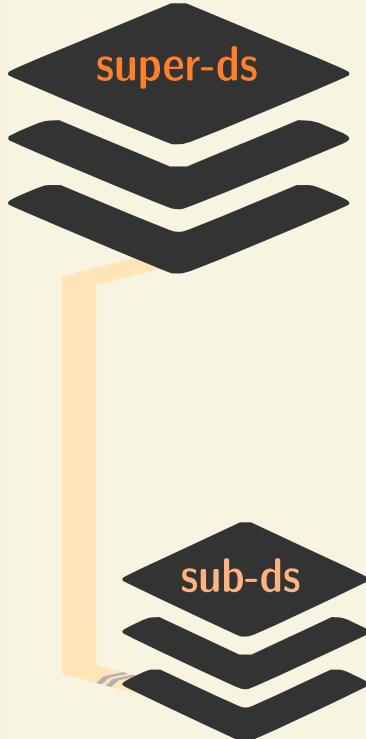
# FROM HERE TO THIS:



BUT: Version control is only one aspect of data management

# **CONSUMING DATASETS**

# CONSUMING DATASETS



DataLad-101/  
books/  
byte-of-python.pdf  
progit.pdf  
TLCL.pdf  
recordings/  
longnow/   
Long\_Now\_\_Conv [...] /  
...  
Long\_Now\_\_Seminars [...] /  
2003\_12\_13 [...]  
2003\_11\_15 [...]  
...  
notes.txt

! Dataset structure is fully flexible to be able to accommodate domain standards or personal preferences.

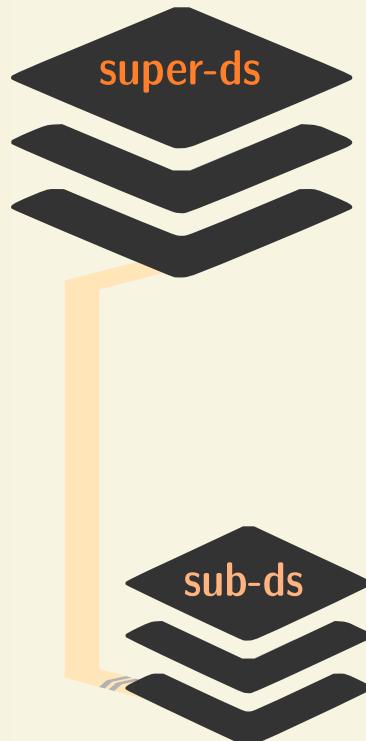
! A dataset can be populated with any type of files, and these files can be saved to the dataset.

! Published repositories can be installed as subdatasets. This nesting can be arbitrarily deep. Datasets can be installed from a path, URL., or data collection.

! DataLad can obtain required subdataset content on demand. Only content elements actually required for an analysis are present. Directory structure is expanded recursively as needed.

! Any content is referenced via the dataset that contains it. Dataset state provides unambiguous version specification for the subdataset.

# CONSUMING DATASETS



DataLad-101/  
books/

byte-of-python.pdf  
progit.pdf  
TLCL.pdf

recordings/

longnow/   
Long\_Now\_\_Conv [...]/  
...  
Long\_Now\_\_Seminars [...]/  
2003\_12\_13 [...]  
2003\_11\_15 [...]  
...

notes.txt

! Dataset structure is fully flexible to be able to accommodate domain standards or personal preferences.

! A dataset can be populated with any type of files, and these files can be saved to the dataset.

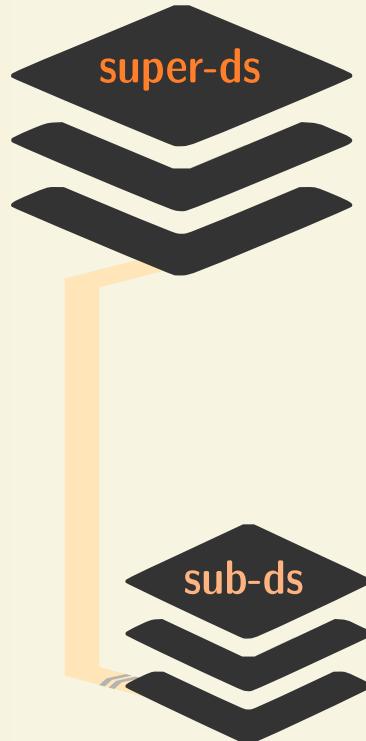
! Published repositories can be installed as subdatasets. This nesting can be arbitrarily deep. Datasets can be installed from a path, URL., or data collection.

! DataLad can obtain required subdataset content on demand. Only content elements actually required for an analysis are present. Directory structure is expanded recursively as needed.

! Any content is referenced via the dataset that contains it. Dataset state provides unambiguous version specification for the subdataset.

- Datasets are light-weight: Upon installation, only small files and meta data about file availability are retrieved.

# CONSUMING DATASETS



DataLad-101/  
books/

byte-of-python.pdf  
progit.pdf  
TLCL.pdf

recordings/

longnow/   
Long\_Now\_\_Conv [...]/  
...  
Long\_Now\_\_Seminars [...]/  
2003\_12\_13 [...]  
2003\_11\_15 [...]  
...

notes.txt

! Dataset structure is fully flexible to be able to accommodate domain standards or personal preferences.

! A dataset can be populated with any type of files, and these files can be saved to the dataset.

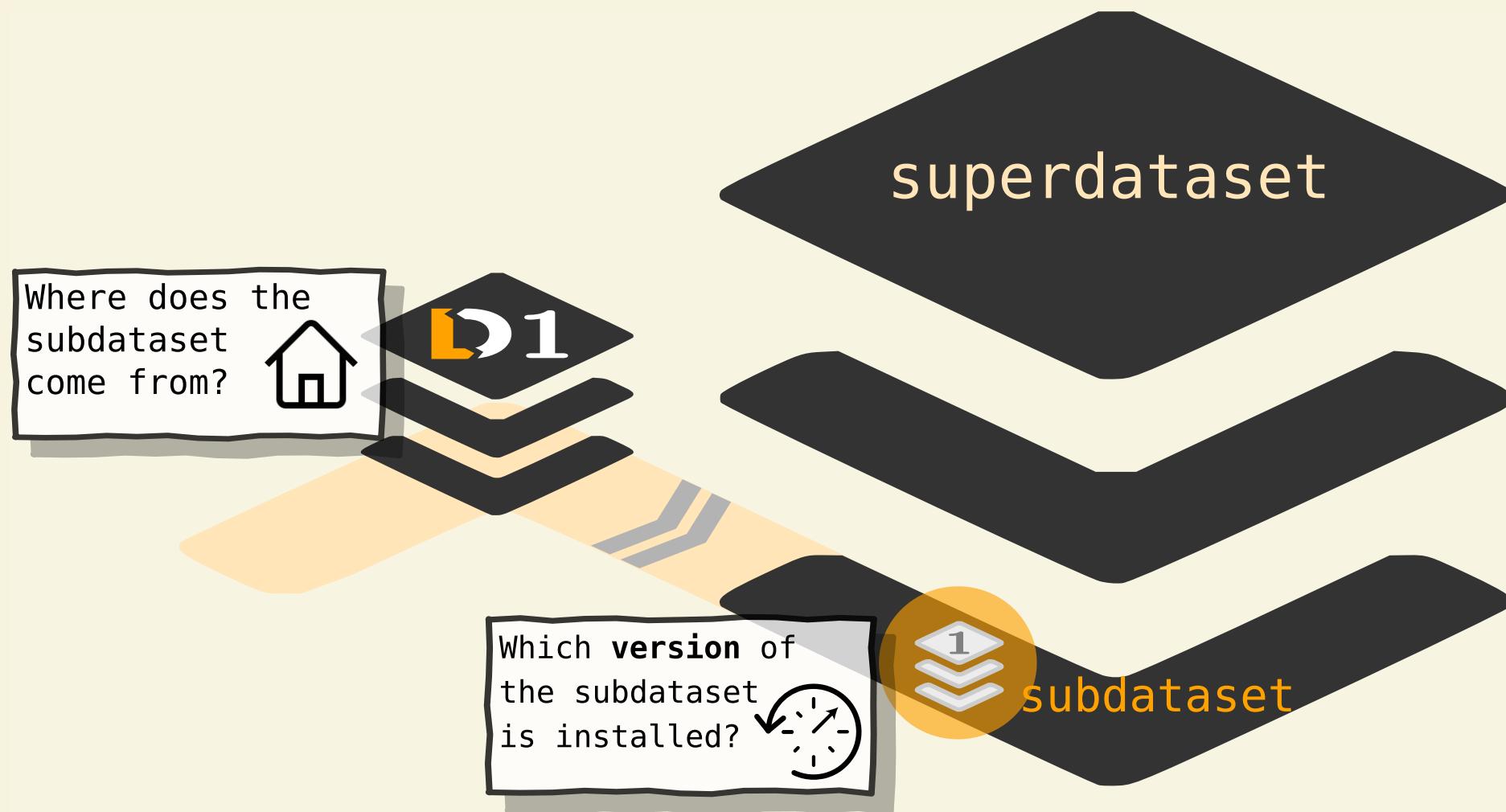
! Published repositories can be installed as subdatasets. This nesting can be arbitrarily deep. Datasets can be installed from a path, URL., or data collection.

! DataLad can obtain required subdataset content on demand. Only content elements actually required for an analysis are present. Directory structure is expanded recursively as needed.

! Any content is referenced via the dataset that contains it. Dataset state provides unambiguous version specification for the subdataset.

- Datasets are light-weight: Upon installation, only small files and meta data about file availability are retrieved.
- Content can be obtained on demand via `datalad get`.

# DATASET NESTING



# **SUMMARY - DATASET CONSUMPTION & NESTING**

## SUMMARY - DATASET CONSUMPTION & NESTING

`datalad install` installs a dataset.

## SUMMARY - DATASET CONSUMPTION & NESTING

**datalad install** installs a dataset.

It can be installed “on its own”: Specify the **--source/-s** of the dataset, and an optional **path** for it to be installed to.

## SUMMARY - DATASET CONSUMPTION & NESTING

**datalad install** installs a dataset.

It can be installed “on its own”: Specify the **--source/-s** of the dataset, and an optional **path** for it to be installed to.

Datasets can be installed as subdatasets within an existing dataset.

## SUMMARY - DATASET CONSUMPTION & NESTING

**datalad install** installs a dataset.

It can be installed “on its own”: Specify the **--source/-s** of the dataset, and an optional **path** for it to be installed to.

Datasets can be installed as subdatasets within an existing dataset.

The **--dataset/-d** option needs a path to the root of the superdataset.

## **SUMMARY - DATASET CONSUMPTION & NESTING**

**datalad install** installs a dataset.

It can be installed “on its own”: Specify the **--source/-s** of the dataset, and an optional **path** for it to be installed to.

**Datasets can be installed as subdatasets within an existing dataset.**

The **--dataset/-d** option needs a path to the root of the superdataset.

**Only small files and metadata about file availability are present locally after an install.**

## SUMMARY - DATASET CONSUMPTION & NESTING

**datalad install** installs a dataset.

It can be installed “on its own”: Specify the **--source/-s** of the dataset, and an optional **path** for it to be installed to.

Datasets can be installed as subdatasets within an existing dataset.

The **--dataset/-d** option needs a path to the root of the superdataset.

**Only small files and metadata about file availability are present locally after an install.**

To retrieve actual file content of larger files, **datalad get** downloads large file content on demand.

## **SUMMARY - DATASET CONSUMPTION & NESTING**

**datalad install** installs a dataset.

It can be installed “on its own”: Specify the **--source/-s** of the dataset, and an optional **path** for it to be installed to.

**Datasets can be installed as subdatasets within an existing dataset.**

The **--dataset/-d** option needs a path to the root of the superdataset.

**Only small files and metadata about file availability are present locally after an install.**

To retrieve actual file content of larger files, **datalad get** downloads large file content on demand.

**datalad status** can report on total and retrieved repository size

## **SUMMARY - DATASET CONSUMPTION & NESTING**

**datalad install** installs a dataset.

It can be installed “on its own”: Specify the **--source/-s** of the dataset, and an optional **path** for it to be installed to.

**Datasets can be installed as subdatasets within an existing dataset.**

The **--dataset/-d** option needs a path to the root of the superdataset.

**Only small files and metadata about file availability are present locally after an install.**

To retrieve actual file content of larger files, **datalad get** downloads large file content on demand.

**datalad status** can report on total and retrieved repository size using **--annex** and **--annex all** options.

## SUMMARY - DATASET CONSUMPTION & NESTING

**datalad install** installs a dataset.

It can be installed “on its own”: Specify the **--source/-s** of the dataset, and an optional **path** for it to be installed to.

Datasets can be installed as subdatasets within an existing dataset.

The **--dataset/-d** option needs a path to the root of the superdataset.

**Only small files and metadata about file availability are present locally after an install.**

To retrieve actual file content of larger files, **datalad get** downloads large file content on demand.

**datalad status** can report on total and retrieved repository size using **--annex** and **--annex all** options.

Datasets preserve their history.

## SUMMARY - DATASET CONSUMPTION & NESTING

**datalad install** installs a dataset.

It can be installed “on its own”: Specify the **--source/-s** of the dataset, and an optional **path** for it to be installed to.

Datasets can be installed as subdatasets within an existing dataset.

The **--dataset/-d** option needs a path to the root of the superdataset.

**Only small files and metadata about file availability are present locally after an install.**

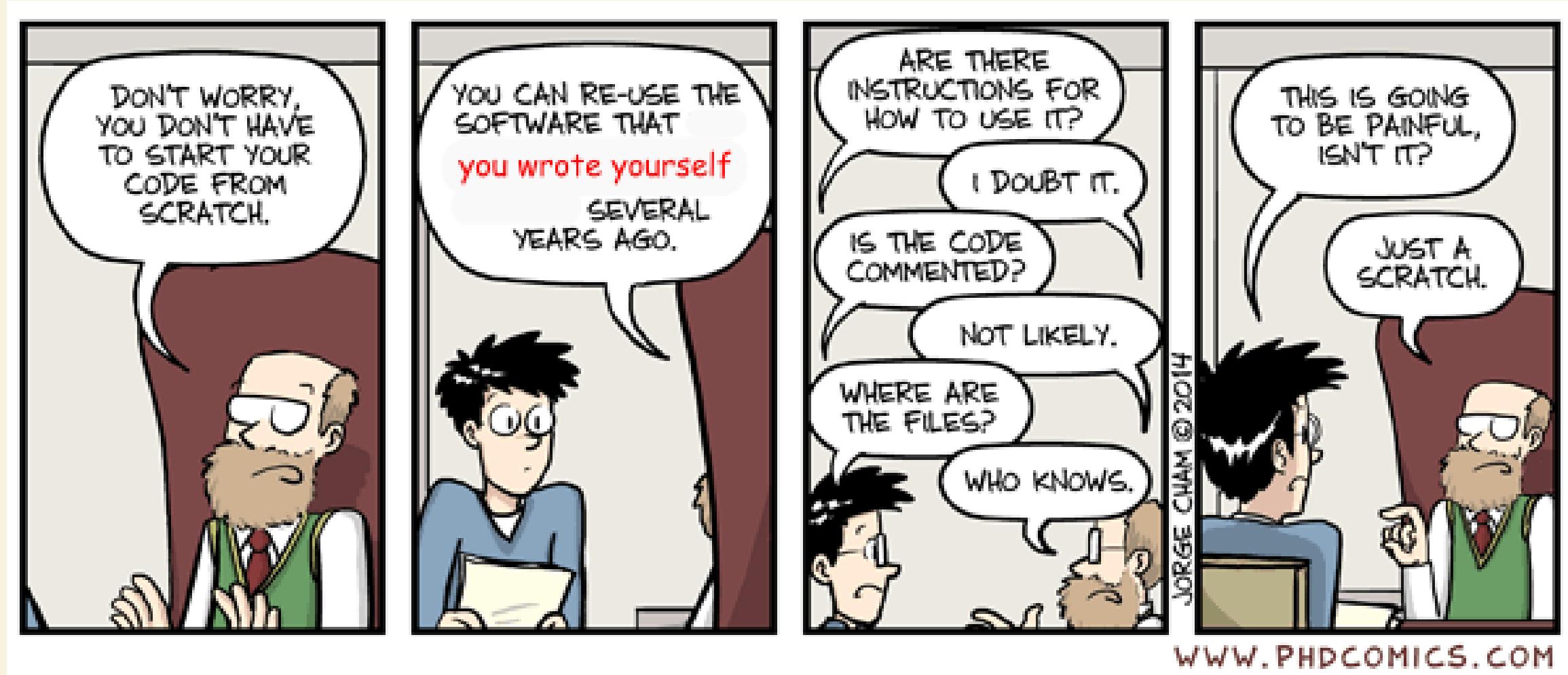
To retrieve actual file content of larger files, **datalad get** downloads large file content on demand.

**datalad status** can report on total and retrieved repository size using **--annex** and **--annex all** options.

Datasets preserve their history.

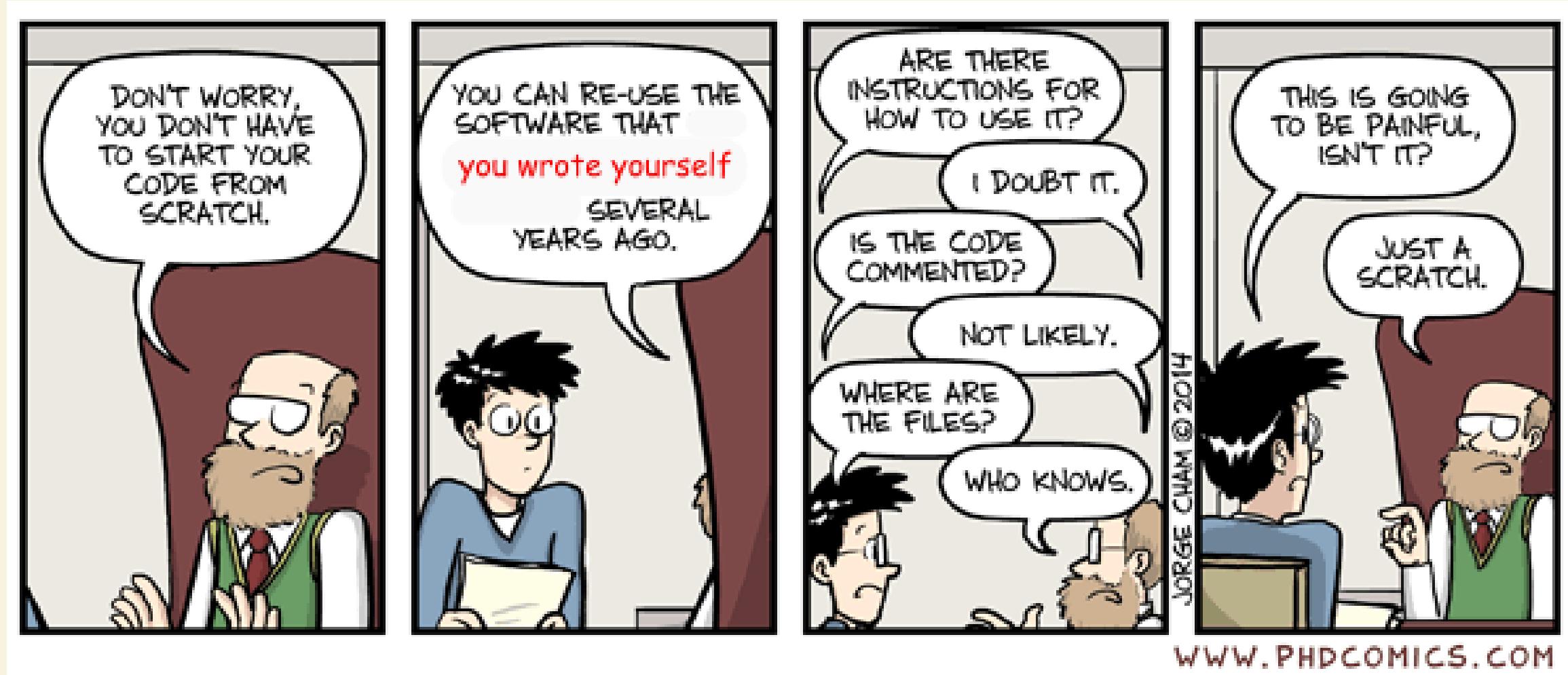
The superdataset records only the *version state* of the subdataset.

# REPRODUCIBLE DATA ANALYSIS



# REPRODUCIBLE DATA ANALYSIS

Image credit: Full comic at <http://phdcomics.com/comics.php?f=1979>



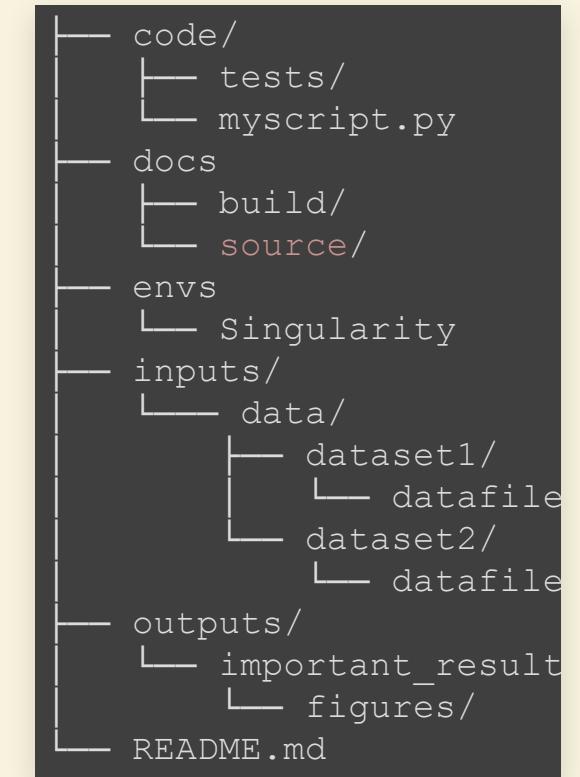
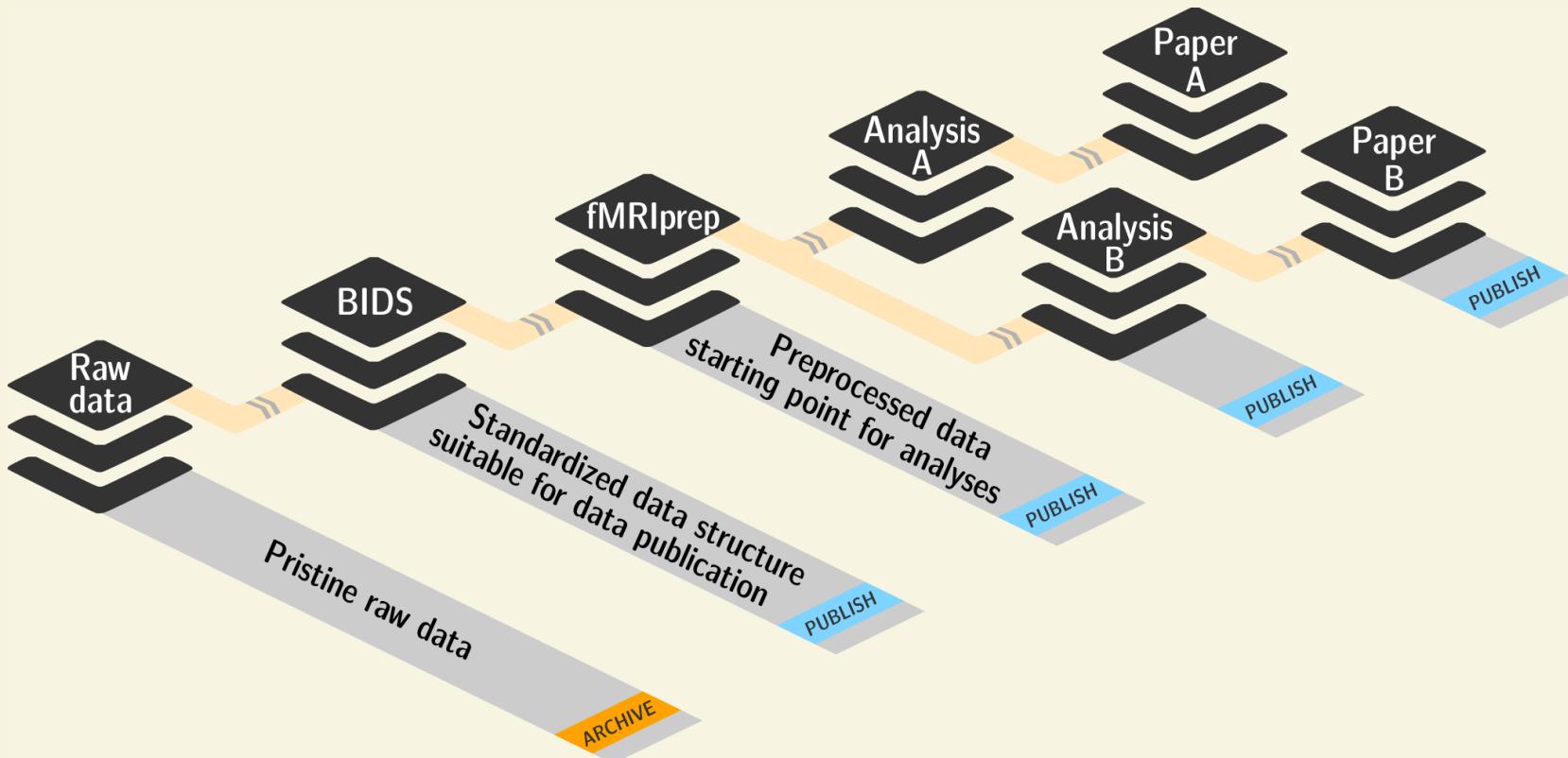
Code to follow along:

[http://handbook.datalad.org/en/latest/code\\_from\\_chapters/10\\_yoda\\_code.html](http://handbook.datalad.org/en/latest/code_from_chapters/10_yoda_code.html)

# BASIC ORGANIZATIONAL PRINCIPLES FOR DATASETS

Keep everything clean and modular

- An analysis is a superdataset, its components are subdatasets, and its structure modular

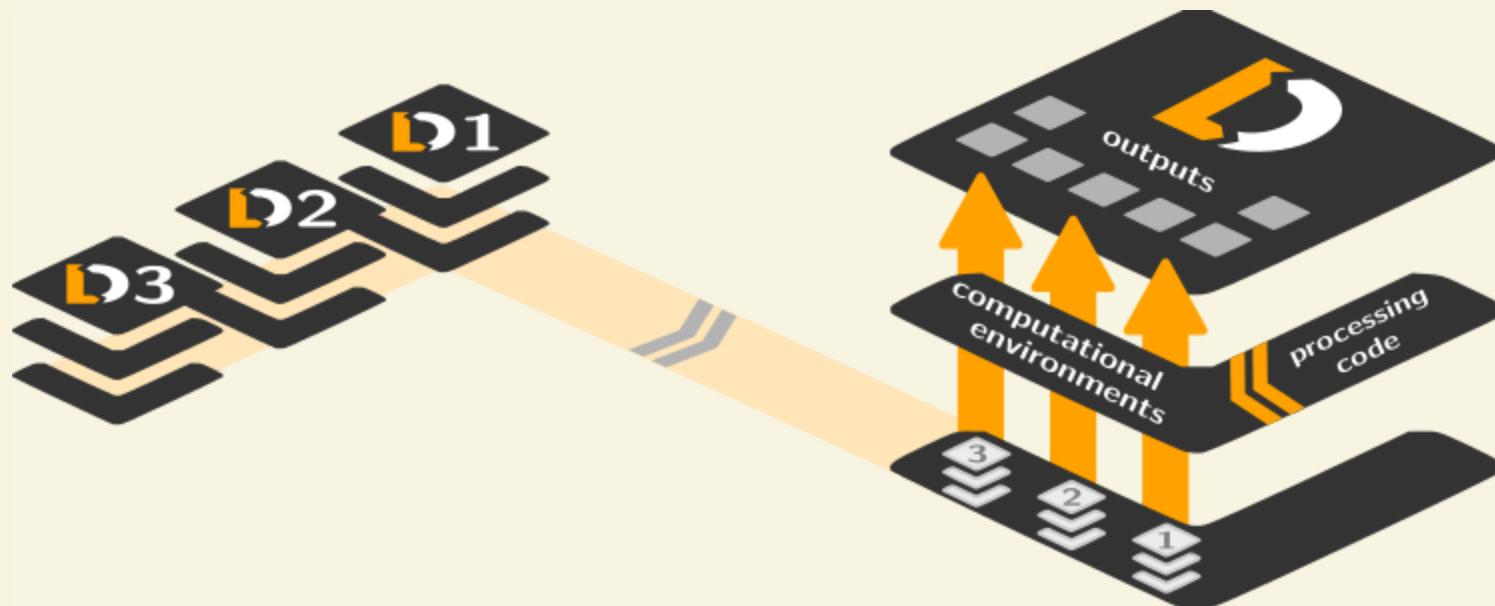


- do not touch/modify raw data: save any results/computations *outside* of input datasets
- Keep a superdataset self-contained: Scripts reference subdatasets or files with *relative paths*

# BASIC ORGANIZATIONAL PRINCIPLES FOR DATASETS

Record where you got it from, where it is now, and what you do to it

- Link datasets (as subdatasets), record data origin
- Collect and store provenance of all contents of a dataset that you create

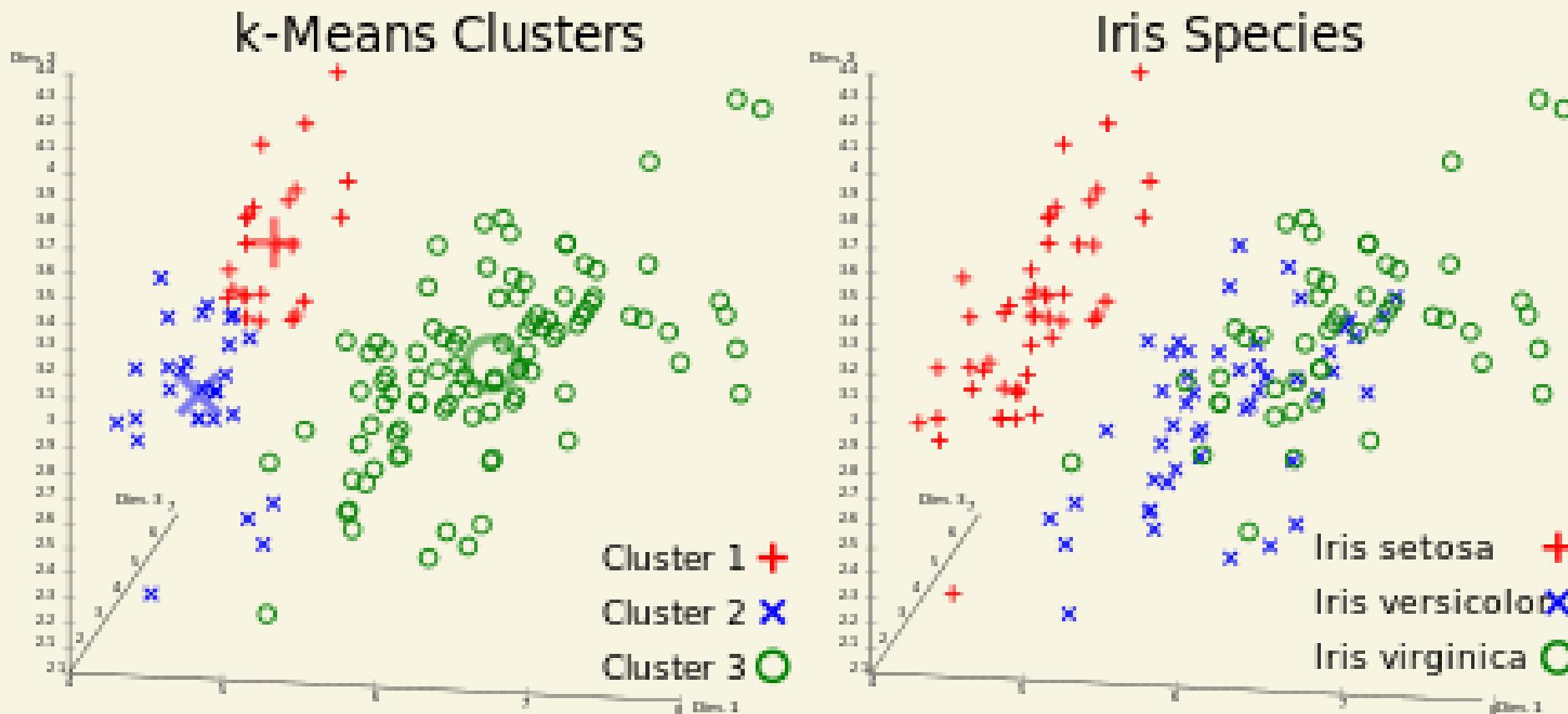
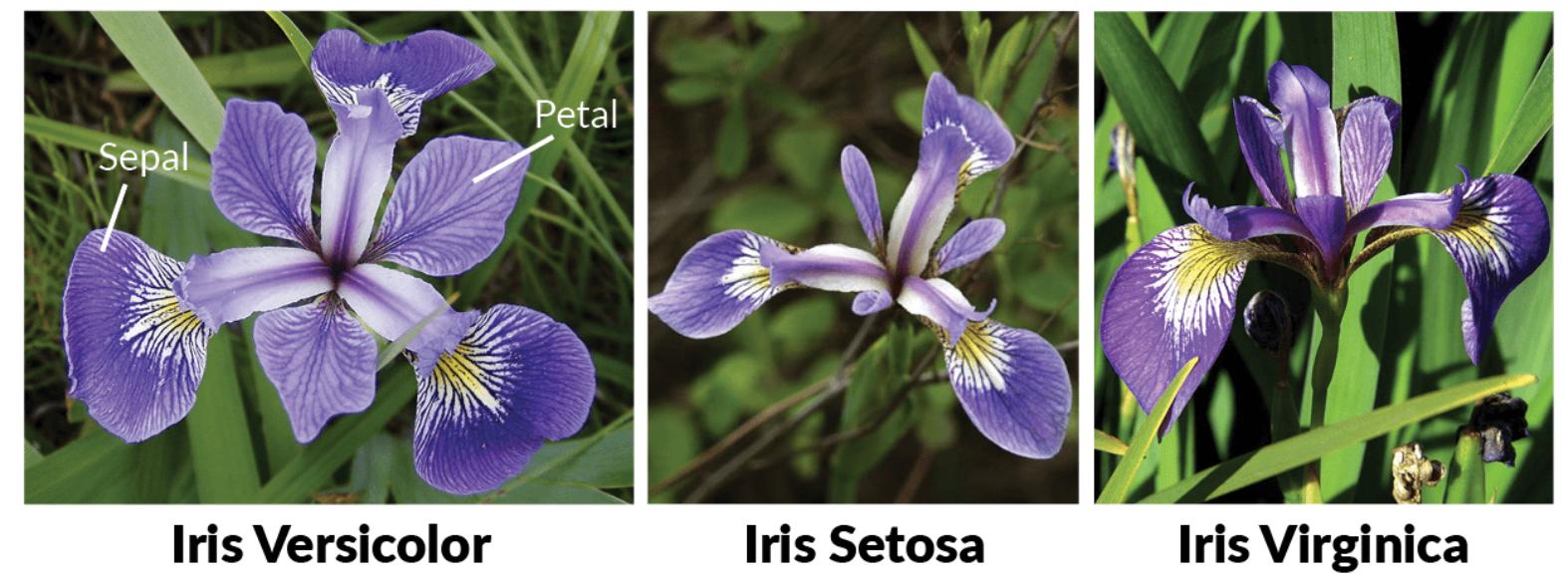


Document everything:

- Which script produced which output? From which data? In which software environment? ...

Find out more about organizational principles in the [YODA principles!](#)

# A CLASSIFICATION ANALYSIS ON THE IRIS FLOWER DATASET

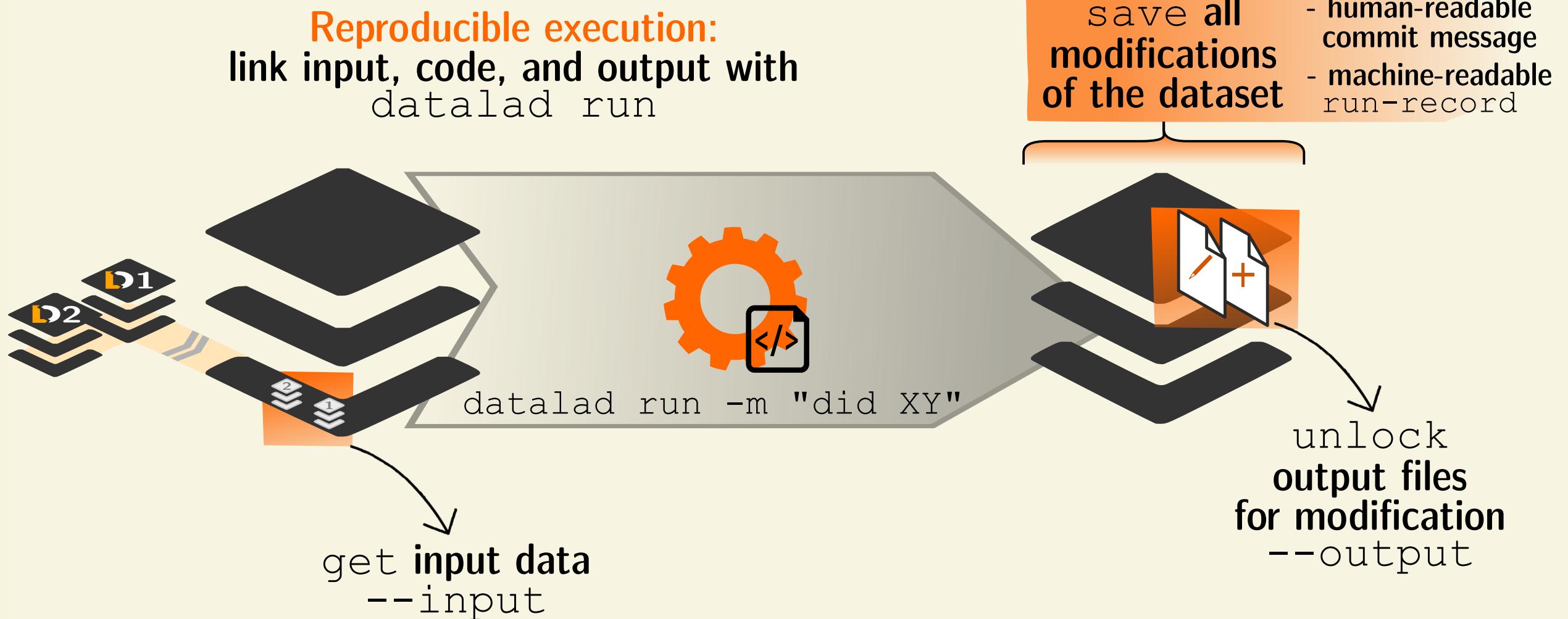


# **REPRODUCIBLE EXECUTION & PROVENANCE CAPTURE**

datalad run

# REPRODUCIBLE EXECUTION & PROVENANCE CAPTURE

datalad run



# **HOW TO GET STARTED WITH BIDS AND DATALAD**

# **HOW TO GET STARTED WITH BIDS AND DATALAD**

**Check out BIDS compliant datasets - with DataLad!**

# HOW TO GET STARTED WITH BIDS AND DATALAD

Check out BIDS compliant datasets - with DataLad!

```
$ datalad install ///openneuro/ds000001
[INFO    ] Cloning http://datasets.datalad.org/openneuro/ds000001 [1 other candidates] into '/tmp
[INFO    ] access to 1 dataset sibling s3-PRIVATE not auto-enabled, enable with:
|       datalad siblings -d "/tmp/ds000001" enable -s s3-PRIVATE
install(ok): /tmp/ds000001 (dataset)

$ cd ds000001
$ ls sub-01/*
sub-01/anat:
sub-01_inplaneT2.nii.gz  sub-01_T1w.nii.gz

sub-01/func:
sub-01_task-balloonanalogrisktask_run-01_bold.nii.gz
sub-01_task-balloonanalogrisktask_run-01_events.tsv
sub-01_task-balloonanalogrisktask_run-02_bold.nii.gz
sub-01_task-balloonanalogrisktask_run-02_events.tsv
sub-01_task-balloonanalogrisktask_run-03_bold.nii.gz
sub-01_task-balloonanalogrisktask_run-03_events.tsv
```

# HOW TO GET STARTED WITH BIDS AND DATALAD

Check out BIDS compliant datasets - with DataLad!

```
$ datalad install ///openneuro/ds000001
[INFO    ] Cloning http://datasets.datalad.org/openneuro/ds000001 [1 other candidates] into '/tmp
[INFO    ] access to 1 dataset sibling s3-PRIVATE not auto-enabled, enable with:
|       datalad siblings -d "/tmp/ds000001" enable -s s3-PRIVATE
install(ok): /tmp/ds000001 (dataset)

$ cd ds000001
$ ls sub-01/*
sub-01/anat:
sub-01_inplaneT2.nii.gz  sub-01_T1w.nii.gz

sub-01/func:
sub-01_task-balloonanalogrisktask_run-01_bold.nii.gz
sub-01_task-balloonanalogrisktask_run-01_events.tsv
sub-01_task-balloonanalogrisktask_run-02_bold.nii.gz
sub-01_task-balloonanalogrisktask_run-02_events.tsv
sub-01_task-balloonanalogrisktask_run-03_bold.nii.gz
sub-01_task-balloonanalogrisktask_run-03_events.tsv
```

Read the DataLad handbook

# HOW TO GET STARTED WITH BIDS AND DATALAD

Check out BIDS compliant datasets - with DataLad!

```
$ datalad install ///openneuro/ds000001
[INFO    ] Cloning http://datasets.datalad.org/openneuro/ds000001 [1 other candidates] into '/tmp
[INFO    ] access to 1 dataset sibling s3-PRIVATE not auto-enabled, enable with:
|           datalad siblings -d "/tmp/ds000001" enable -s s3-PRIVATE
install(ok): /tmp/ds000001 (dataset)

$ cd ds000001
$ ls sub-01/*
sub-01/anat:
sub-01_inplaneT2.nii.gz  sub-01_T1w.nii.gz

sub-01/func:
sub-01_task-balloonanalogrisktask_run-01_bold.nii.gz
sub-01_task-balloonanalogrisktask_run-01_events.tsv
sub-01_task-balloonanalogrisktask_run-02_bold.nii.gz
sub-01_task-balloonanalogrisktask_run-02_events.tsv
sub-01_task-balloonanalogrisktask_run-03_bold.nii.gz
sub-01_task-balloonanalogrisktask_run-03_events.tsv
```

Read the DataLad handbook

An interactive, hands-on crash-course (free and open source)

# ACKNOWLEDGEMENTS

## The BIDS standard

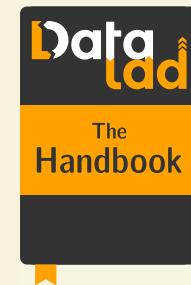
- Chris Gorgolewski
- Russ Poldrack
- 100(0?)+ additional contributors



- Michael Hanke
- Yaroslav Halchenko
- Joey Hess (git-annex)
- Benjamin Poldrack
- Kyle Meyer
- 22+ additional contributors

## The DataLad Handbook

- Laura Waite
- Michael Hanke
- 11+ additional contributors



**THANK YOU!**

**QUESTIONS?**