

"Motivation & Basics of version control" starts in

00 | 15 | 56 | 08

Days Hours Minutes Seconds

A digital timer is displayed horizontally, divided into four sections by vertical lines. The first section, labeled 'Days', shows '00'. The second section, labeled 'Hours', shows '15'. The third section, labeled 'Minutes', shows '56'. The fourth section, labeled 'Seconds', shows '08'. All text is in a large, bold, black font.

PARTICIPATION MODES

<http://etc.ch/XcWI>



PREREQUISITES: INSTALLATION AND CONFIGURATION

- Your installed version of DataLad should be 0.16.1

```
datalad --version  
0.16.1
```

- DataLad relies on Git to create a revision history with detailed information on what was changes, when, and how. Therefore, you should tell Git who you are and configure a Git identity (name and email)

```
$ git config --list  
user.name=Adina Wagner  
user.email=adina.wagner@t-online.de  
[...]
```

- Set a Git identity using

```
$ git config set --global user.name "Adina Wagner"  
$ git config set --global user.email "adina.wagner@t-online.de"
```

Find installation and configuration instructions at handbook.datalad.org

USING DATALAD

- DataLad can be used from the command line

```
datalad create mydataset
```

- ... or with its Python API

```
import datalad.api as dl
dl.create(path="mydataset")
```

- ... and other programming languages can use it via system call

```
# in R
> system("datalad create mydataset")
```

USING DATALAD

- Every DataLad command consists of a main command followed by a sub-command. The main and the sub-command can have options.

The screenshot shows the DataLad API documentation. It includes a header with the command structure: `datalad [--GLOBAL-OPTION <opt. flag spec.>] COMMAND [ARGUMENTS] [--OPTION <opt. flag spec>]`. Below this are two sections: **GLOBAL OPTIONS** and **COMMAND OPTIONS**. The **GLOBAL OPTIONS** section lists: `-c KEY=VALUE` (Set config variables), `-f/--output-format default|json|json_pp|tailored` (Specify the format for command result rendering), and `-l/--log-level critical|error|warning|info|debug` (Set logging verbosity level). The **COMMAND OPTIONS** section lists: `-d/--dataset` (Dataset location: path to root, or ^ for superdataset), `-D/--description` (A location description (e.g., "my backup server")), `-f/--force` (Force execution of a command (Dangerzone!)), `-m/--message` (A description about a change made to the dataset), `-r/--recursive` (Perform an operation recursively across subdatasets), and `-R/--recursion-limit <n>` (Limit recursion to n subdataset levels). A callout box on the right states: "Each datalad invocation can have two sets of options: general options are given first, command-specific ones go after the subcommand."

- Example (main command, subcommand, several subcommand options):

```
$ datalad save -m "Saving changes" --recursive
```

- Use `--help` to find out more about any (sub)command and its options, including detailed description and examples (`q` to close). Use `-h` to get a short overview of all options

```
$ datalad save -h
Usage: datalad save [-h] [-m MESSAGE] [-d DATASET] [-t ID] [-r] [-R LEVELS]
                  [-u] [-F MESSAGE_FILE] [--to-git] [-J NJOBS] [--amend]
                  [--version]
                  [PATH ...]

Use '--help' to get more comprehensive information.
```

DATALAD DATASETS

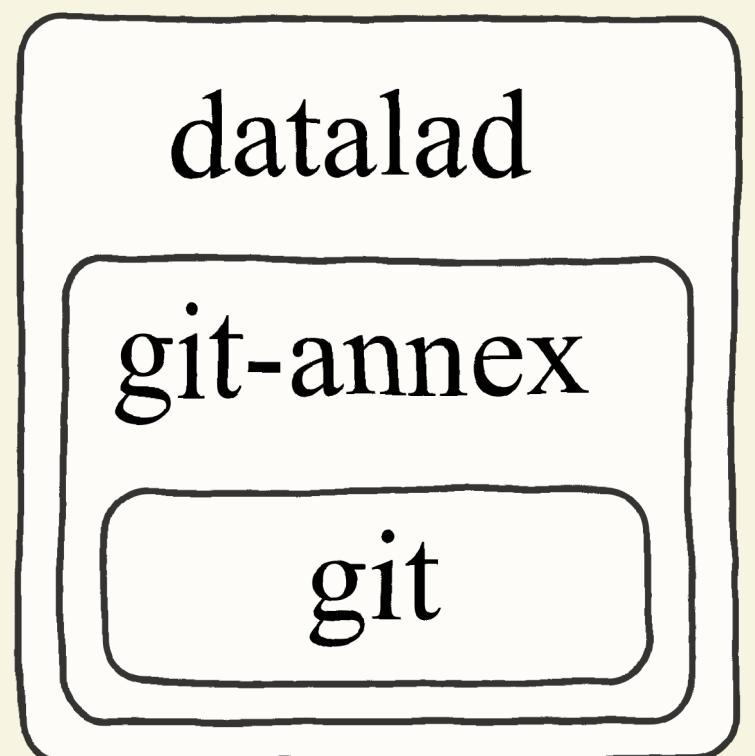
- DataLad's core data structure
 - Dataset = A directory managed by DataLad
 - Any directory of your computer can be managed by DataLad.
 - Datasets can be *created* (from scratch) or *installed*
 - Datasets can be nested: *linked subdirectories*
- Let's start by creating a dataset:

```
$ datalad create -c text2git my-dataset
```

Code:

psychoinformatics-de.github.io/rdm-course/01-content-tracking-with-datalad/index.html#getting-started-create-an-empty-dataset

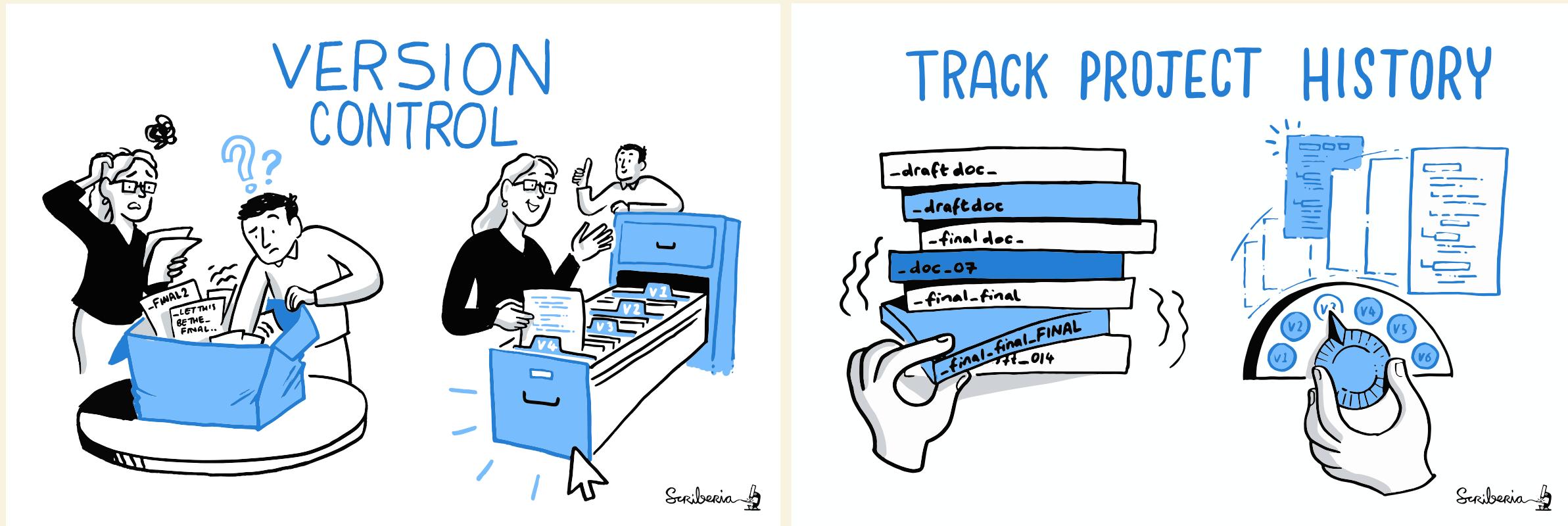
DATALAD DATASETS



A DataLad dataset is a joined Git + git-annex repository

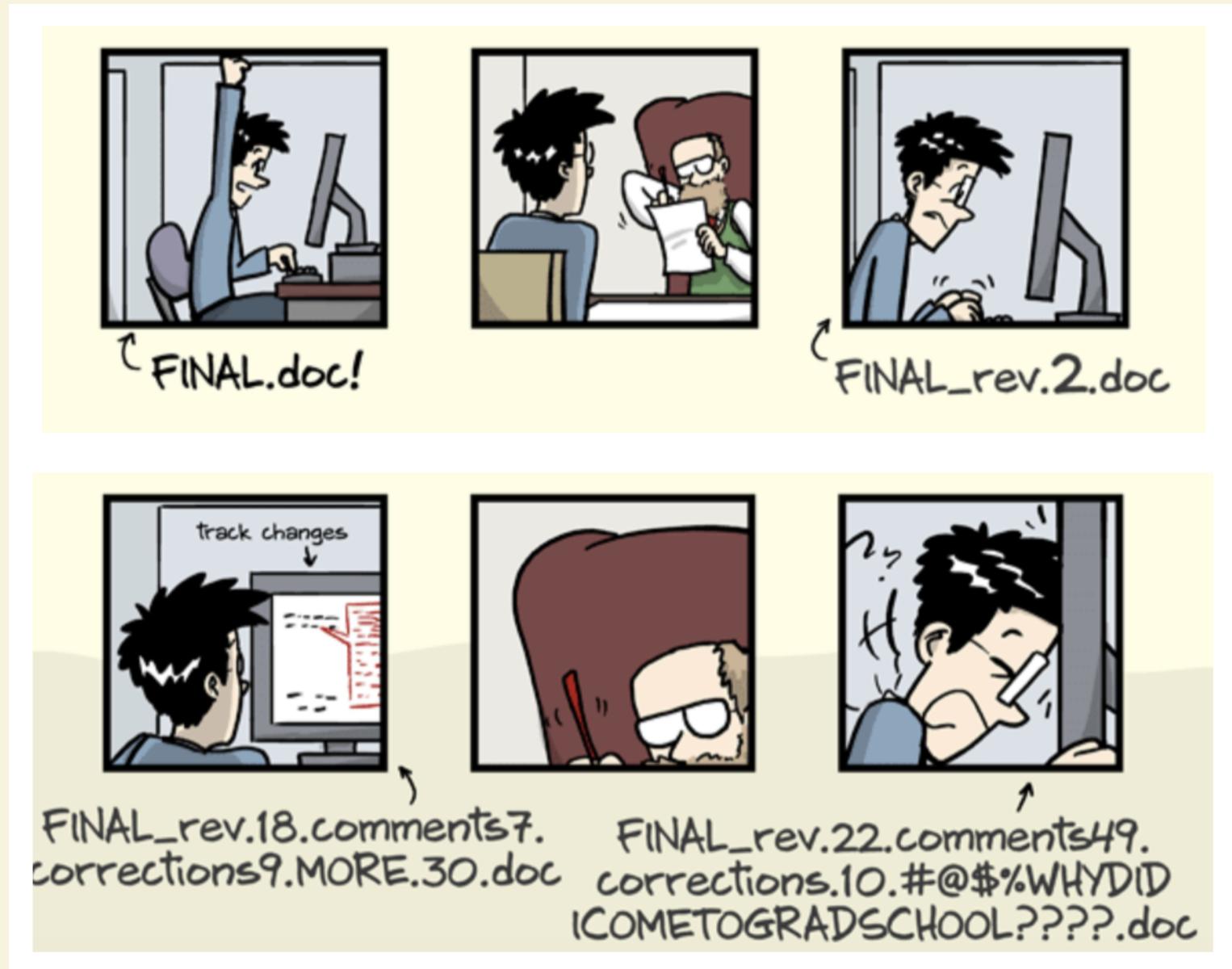
WHAT IS VERSION CONTROL?

Image credit: Illustration adapted from Scriberia and The Turing Way



- keep things organized
- keep track of changes
- revert changes or go back to previous states

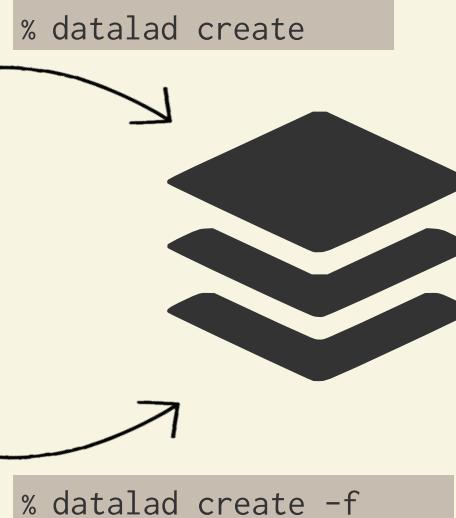
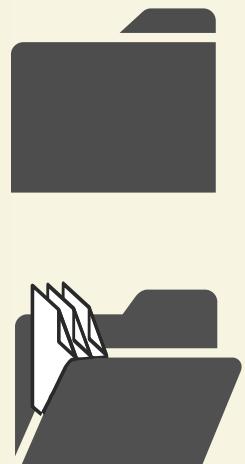
WHY VERSION CONTROL?



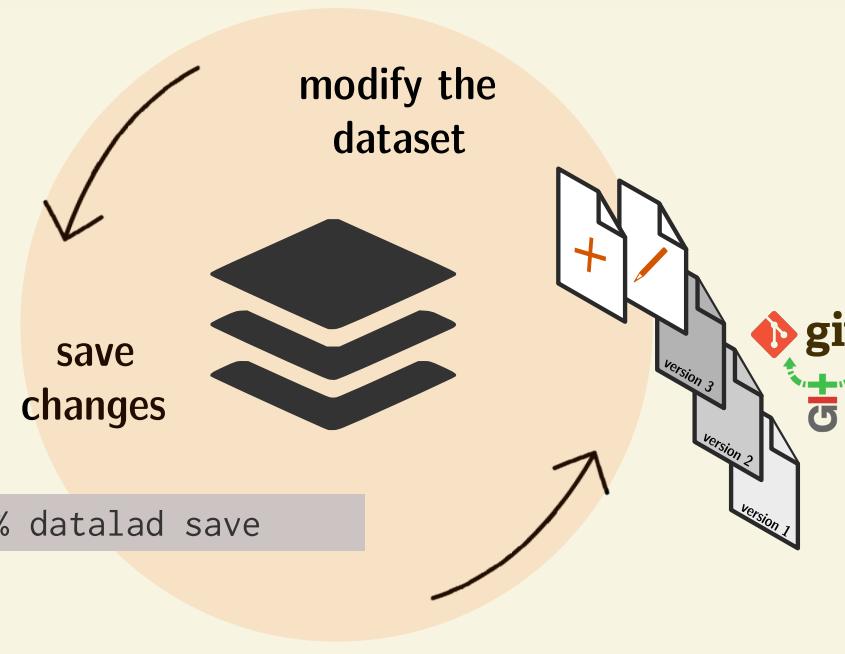
VERSION CONTROL

- DataLad knows two things: Datasets and files

create new, empty datasets to populate...



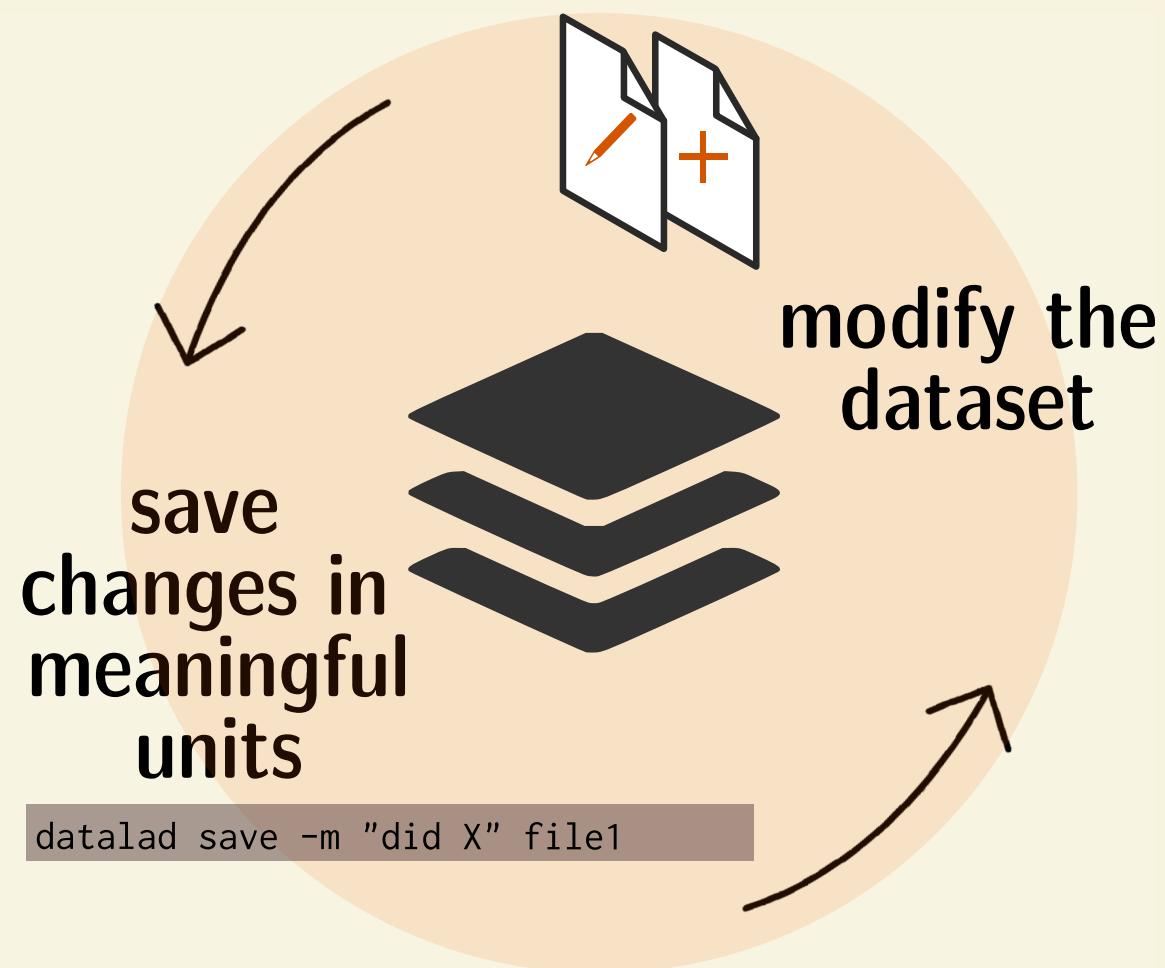
.... or transform existing directories into datasets



- Every file you put into a in a dataset can be easily version-controlled, regardless of size, with the same command: *datalad save*

LOCAL VERSION CONTROL

Procedurally, version control is easy with DataLad!



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

THIS MEANS: YOU CAN ALSO VERSION CONTROL DATA!

```
$ datalad save \  
-m "Adding raw data from neuroimaging study 1" \  
sub-*  
add(ok): sub-1/anat/T1w.json (file)  
add(ok): sub-1/anat/T1w.nii.gz (file)  
add(ok): sub-1/anat/T2w.json (file)  
add(ok): sub-1/anat/T2w.nii.gz (file)  
add(ok): sub-1/func/sub-1-run-1_bold.json (file)  
add(ok): sub-1/func/sub-1-run-1_bold.nii.gz (file)  
add(ok): sub-10/anat/T1w.json (file)  
add(ok): sub-10/anat/T1w.nii.gz (file)  
add(ok): sub-10/anat/T2w.json (file)  
add(ok): sub-10/anat/T2w.nii.gz (file)  
[110 similar messages have been suppressed]  
save(ok): . (dataset)  
action summary:  
add (ok: 120)  
save (ok: 1)
```

VERSION CONTROL

- Your dataset can be a complete research log, capturing everything that was done, when, by whom, and how

Date	Author	Change summary (commit message)
2020-06-05 10:58 +0200	Adina Wagner	M [master] {upstream/master} {upstream/HEAD} Merge pull request #12 from psychoinformatics-de
2020-06-05 08:24 +0200	Adina Wagner	o [finalround] {upstream/finalround} add results from computing with mean instead of median
2020-06-05 09:09 +0200	Michael Hanke	o Change wording, clarify comment
2020-06-05 07:26 +0200	Michael Hanke	M Merge remote-tracking branch 'gh-mine/finalround'
2020-05-28 16:39 +0200	Asim H Dar	o Added datalad.get() so S2SRMS() pulls data and can run standalone
2020-05-18 08:25 +0200	Adina Wagner	o {gh-asim/finalround} S2SRMS: example implementation of the S2SRMS method suggested by R2
2020-05-01 17:38 +0200	Adina Wagner	o Minor edits as suggested by reviewer 2
2020-05-29 09:04 +0200	Adina Wagner	M Merge pull request #13 from psychoinformatics-de/adswa-patch-1
2020-05-24 09:15 +0200	Adina Wagner	o {upstream/adswa-patch-1} Fix installation instructions
2020-05-24 09:53 +0200	Adina Wagner	M Merge pull request #14 from psychoinformatics-de/bf-data
2020-05-24 09:33 +0200	Adina Wagner	o [bf-data] One-time datalad import
2020-05-24 09:32 +0200	Adina Wagner	o install and get relevant subdataset data
2020-03-18 10:19 +0100	Michael Hanke	M Merge pull request #8 from psychoinformatics-de/adswa-patch-1
2019-12-19 10:22 +0100	Adina Wagner	o {gh-asim/adswa-patch-1} add sklearn to requirements
2020-03-18 10:13 +0100	Michael Hanke	o Tune new figure caption
2020-03-18 10:03 +0100	Michael Hanke	M Merge pull request #11 from ElectronicTeaCup/revision_2
2020-03-18 09:59 +0100	Adina Wagner	M [revision_2] {gh-asim/revision_2} Merge branch 'revision_2' of github.com:ElectronicTe
2020-03-18 09:58 +0100	Michael Hanke	o Last detections
2020-03-18 09:59 +0100	Adina Wagner	o name parameter in caption

- Interact with the history:
- reset your dataset (or subset of it) to a previous state,
- throw out changes or bring them back,
- find out what was done when, how, why, and by whom
- Identify precise versions: Use data in the most recent version, or the one from 2018, or...
- ...

PREVIEW: START TO RECORD PROVENANCE

- Have you ever saved a PDF to read later onto your computer, but forgot where you got it from?
- Digital Provenance = "*The tools and processes used to create a digital file, the responsible entity, and when and where the process events occurred*"
- The history of a dataset already contains provenance, but there is more to record - for example: Where does a file come from? `datalad download-url` is helpful

SUMMARY - LOCAL VERSION CONTROL

datalad create creates an empty dataset.

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

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 download-url obtains web content and records its origin.

It even takes care of saving the change.

datalad status reports the current state of the dataset.

A clean dataset status (no modifications, not untracked files) is good practice.

QUESTIONS!

Awkward silence can be bridged with awkward MC questions :)

TEASER: TIME-TRAVELLING

Comprehensive walk-through handbook.datalad.org/basics/101-137-history.html

- Mistakes are not forever anymore: Past changes can transparently be undone
- Become a time-bender: Travel back in time or rewrite history
- Prerequisite: Understand Git IDs and "refs"
 - Commit hash/Commit SHA: A 40-character string identifying each commit
 - Branch names, e.g., *main*
 - Tags, e.g., v.0.1
 - A pointer to the checked-out (current) commit on the current branch, *HEAD*

```
commit ac95f2fa94453e0e730f9183de7e37deb3f7b3df
Refs: [ux-exportfigshare], {gh-adswa/ux-exportfigshare}, 0.15.6-824-gac95f2fa9
Author: Adina Wagner <adina.wagner@t-online.de>
AuthorDate: Mon Mar 14 10:47:42 2022 +0100
Commit: Adina Wagner <adina.wagner@t-online.de>
CommitDate: Mon Mar 14 11:03:40 2022 +0100

UX: Yield impossible result instead of RunTimeError in dirty ds

This aims at standardizing result reporting/behavior over operations that
refuse to operate in dirty datasets. export-to-figshare does it similarly
to run now, and yields an impossible result. Fixes #6474

Before:
datalad export-to-figshare
[ERROR] RuntimeError(Paranoid authors of DataLad refuse to proceed in a dir

Now:
datalad export-to-figshare
export_to_figshare(impossible): /tmp/super (dataset) [clean dataset required
...
datalad/distributed/export_to_figshare.py | 12 ++++++-----
```

Code: psychoinformatics-de.github.io/rdm-course/01-content-tracking-with-datalad/index.html#breaking-things-and-repairing-them

SUMMARY: INTERACTING WITH GIT'S HISTORY (TEASER)

Interactions with Git's history require Git commands, but are immensely powerful

More in handbook.datalad.org/basics/101-137-history.html

git restore is a dangerous (!), but sometimes useful command:

It removes unsaved modifications to restore files to a past, saved state. What has been removed by it can not be brought back to life!

git revert [hash] transparently undoes a past commit

It will create a new entry in the revision history about this.

Commands that will be introduced later:

git checkout lets you time-travel.

Commands that are out of scope but useful to know:

git rebase changes and **git reset** rewinds history without creating a commit about it (see Handbook chapter for examples).

A life-saver that is not well-known: git reflog

A time-limited backlog of every past performed action, can undo every mistake except **git restore** and **git clean**.

QUESTIONS!

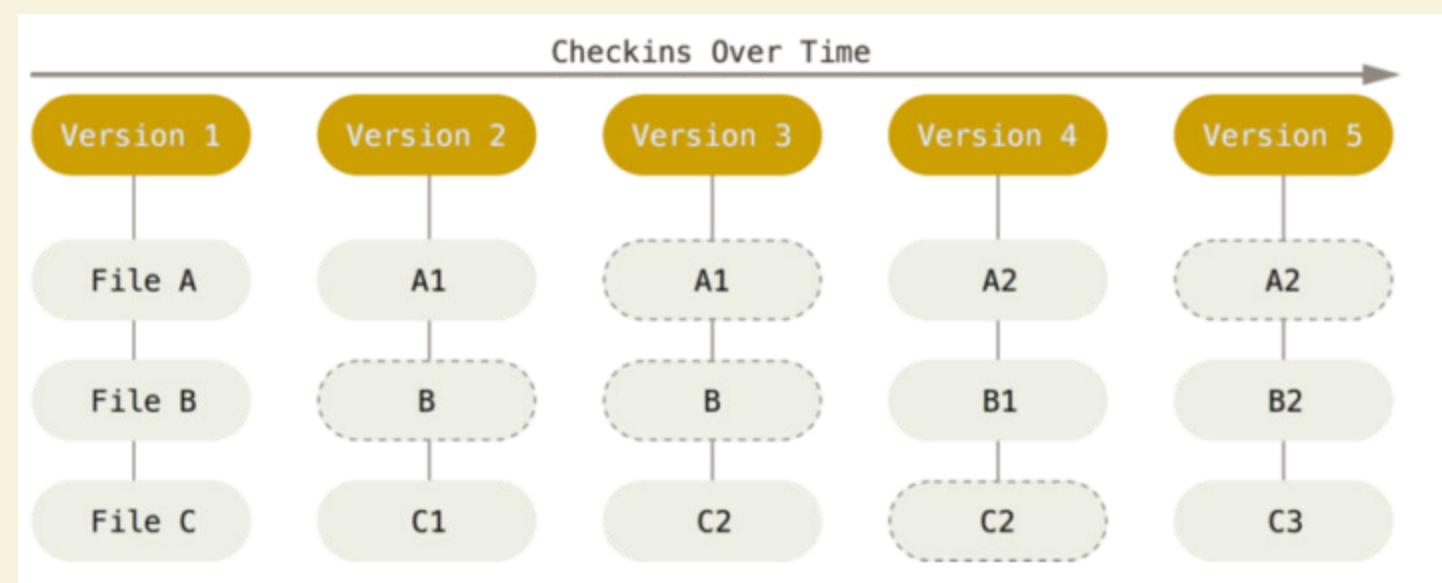
Awkward silence can be bridged with awkward MC questions :)

A LOOK UNDERNEATH THE HOOD

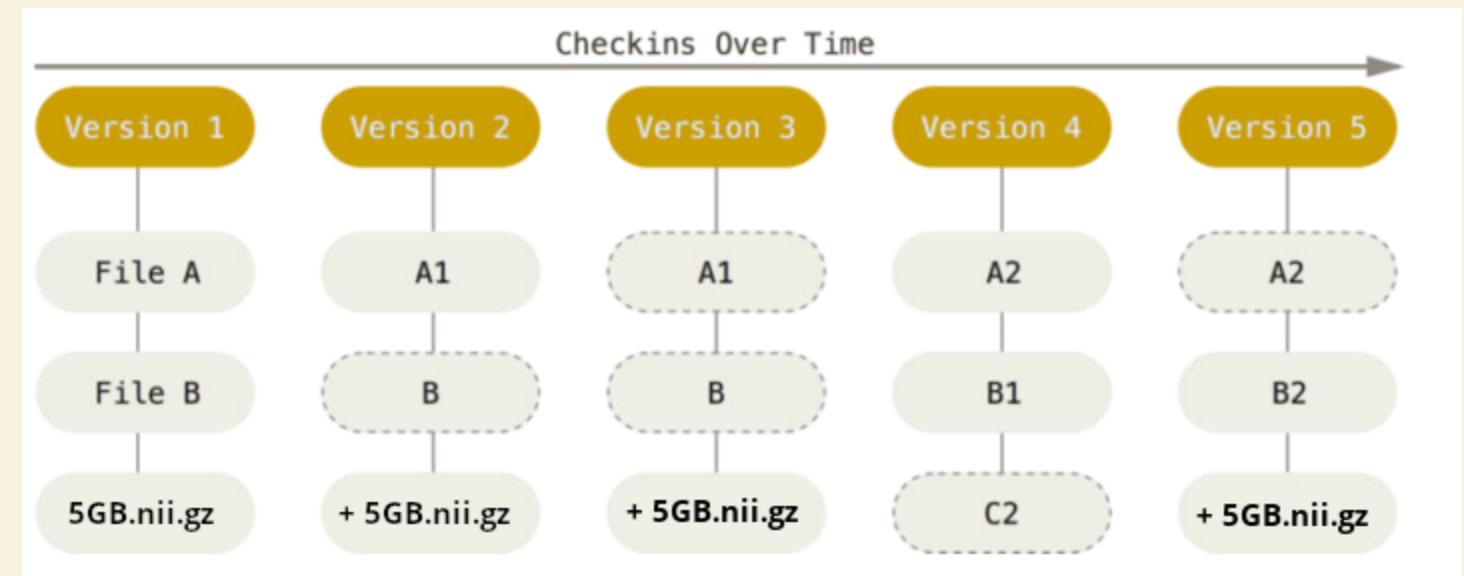
(IN-DEPTH EXPLANATIONS HOW AND WHY THINGS WORK, WITH PLENTY OF TEASERS TO ADDITIONAL FEATURES)

THERE ARE TWO VERSION CONTROL TOOLS AT WORK - WHY?

Git does not handle large files well.



THERE ARE TWO VERSION CONTROL TOOLS AT WORK - WHY?



Git does not handle large files well.

And repository hosting services refuse to handle large files:

```
adina@muninn in /tmp/myresearch on git:master
> git push gh-adswa master
Enumerating objects: 3, done.
Counting objects: 100% (3/3), done.
Delta compression using up to 8 threads
Compressing objects: 100% (2/2), done.
Writing objects: 100% (3/3), 497.87 KiB | 161.00 KiB/s, done.
Total 3 (delta 0), reused 0 (delta 0), pack-reused 0
remote: error: Trace: 64a78dd41ece8e5493fe33f97397a90ef9c91260ba32786970dbdcf5c4e0dd
remote: error: See http://git.io/iEPt8g for more information.
remote: error: File output.dat is 500.00 MB; this exceeds GitHub's file size limit of 100.00 MB
remote: error: GH001: Large files detected. You may want to try Git Large File Storage - https://git-
To github.com:adswa/myresearch.git
 ! [remote rejected] master -> master (pre-receive hook declined)
error: failed to push some refs to 'github.com:adswa/myresearch.git'
```



git-annex to the rescue! Let's take a look how it works

CONSUMING DATASETS

- A dataset can be created from scratch/existing directories:

```
$ datalad create mydataset  
[INFO] Creating a new annex repo at /home/adina/mydataset  
create(ok): /home/adina/mydataset (dataset)
```

- but datasets can also be installed from paths or from URLs:

```
$ datalad clone https://github.com/datalad-datasets/human-connectome-project-openaccess  
install(ok): /tmp/HCP (dataset)
```

Hint: Did you know that you can get the Human Connectome Project Open Access Data as a Dataset?

CONSUMING DATASETS

- Here's how to get a dataset:

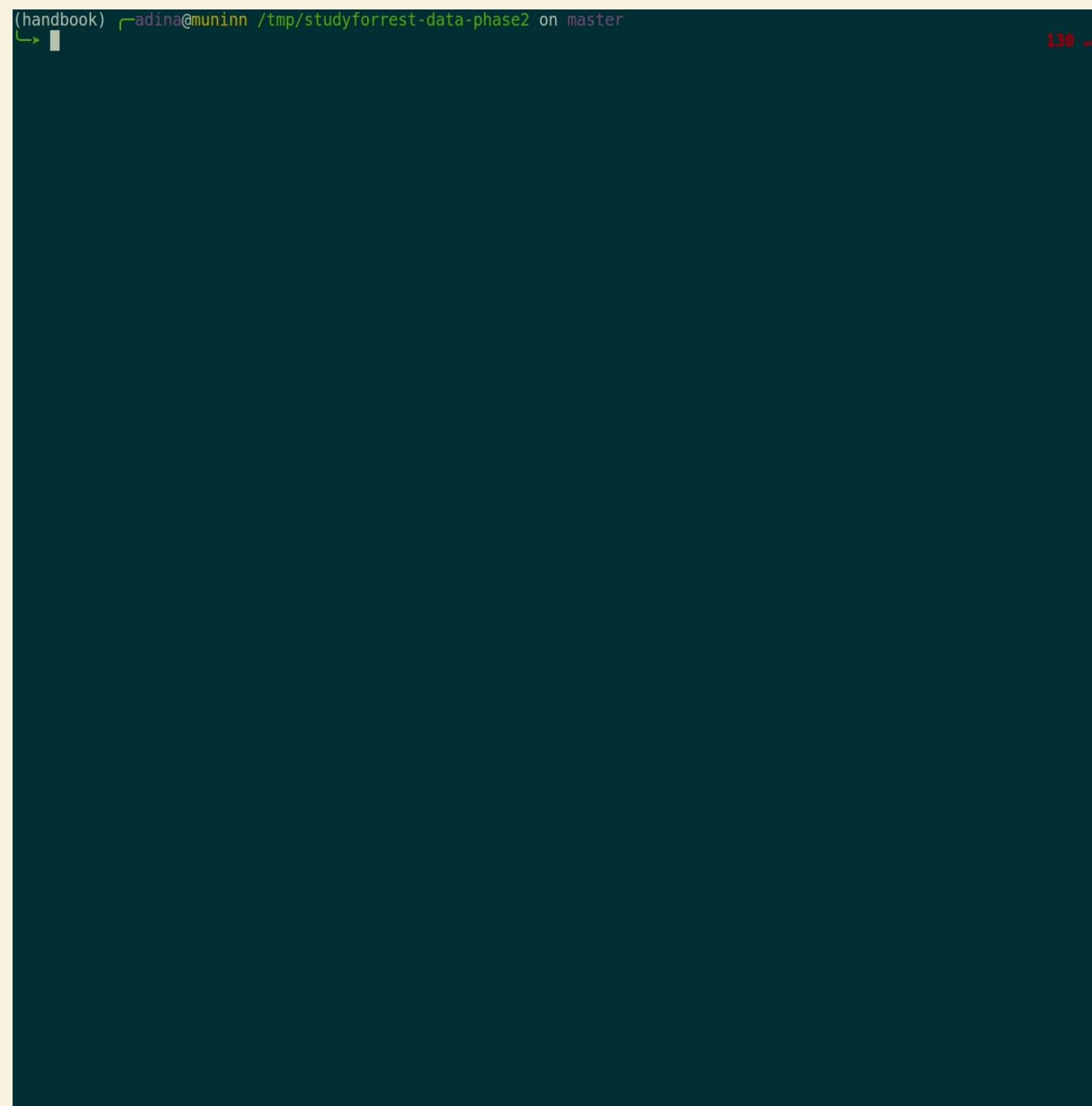
The screenshot shows a GitHub repository page for 'psychoinformatics-de / studyforrest-data-phase2'. The repository has 1 pull request, 4 issues, 1 action, 2 branches, and 1 tag. The 'Code' tab is selected. The main area displays a list of 20 commits from user 'mih'. The commits are as follows:

Commit	Message	Date
.datalad	[DATALAD] dataset aggregate metadata update	2 years ago
code	Fix type in physio log converter (fixes gh-11)	3 years ago
src	Recover lost segment from eyetracker (closes gh-3)	5 years ago
stimuli	Add BIDS-compatible stimuli/ directory (with symlinks)	4 years ago
sub-01	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-02	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-03	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-04	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-05	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-06	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-09	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-10	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-14	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-15	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-16	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-17	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-18	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-19	BF: Re-import respiratory trace after bug fix in converte...	3 years ago
sub-20	BF: Re-import respiratory trace after bug fix in converte...	3 years ago

The repository has 1 star, 6 forks, and 8 open issues. The 'About' section describes the data as Phase2 data for studyforrest.org, involving movie, eyetracking, retmapping, and visual localizers, and is BIDS-compliant. It links to studyforrest.org, a Readme file, and a View license. A 'First public release' was made on Mar 26, 2016. There are no packages published, 3 contributors (mih, dakot, adswa), and the repository uses Python as its primary language.

CONSUMING DATASETS

- Here's how a dataset looks after installation:

A dark-themed terminal window with a black background. At the top, there is a light-colored header bar containing the text "(handbook)" in green, "radina@muninn" in purple, "/tmp/studyforrest-data-phase2" in green, "on master" in white, and a red progress bar icon on the right. The main body of the terminal is completely blank, showing only the dark background.

PLENTY OF DATA, BUT LITTLE DISK-USAGE

- Cloned datasets are lean. "Meta data" (file names, availability) are present, but **no file content**:

```
$ datalad clone git@github.com:psychoinformatics-de/studyforrest-data-phase2.git
install(ok): /tmp/studyforrest-data-phase2 (dataset)
$ cd studyforrest-data-phase2 && du -sh
18M  .
```

- files' contents can be retrieved on demand:

```
$ datalad get sub-01/ses-movie/func/sub-01_ses-movie_task-movie_run-1_bold.nii.gz
get(ok): /tmp/studyforrest-data-phase2/sub-01/ses-movie/func/sub-01_ses-movie_task-movie_run-1_bold.nii.gz
```

- Have more access to your computer than you have disk-space:

```
# eNKI dataset (1.5TB, 34k files):
$ du -sh
1.5G  .
# HCP dataset (~200TB, >15 million files)
$ du -sh
48G  .
```

PLENTY OF DATA, BUT LITTLE DISK-USAGE

Drop file content that is not needed:

```
$ datalad drop sub-01/ses-movie/func/sub-01_ses-movie_task-movie_run-1_bold.nii.gz  
drop(ok): /tmp/studyforrest-data-phase2/sub-01/ses-movie/func/sub-01_ses-movie_task-movie_run-1_bold.nii.g
```

When files are dropped, only "meta data" stays behind, and they can be re-obtained on demand.

```
dl.get('input/sub-01')  
[really complex analysis]  
dl.drop('input/sub-01')
```

GIT VERSUS GIT-ANNEX

Data in datasets is either stored in Git or git-annex

By default, everything is *annexed*, i.e., stored in a dataset annex by git-annex

Git	git-annex
handles small files well (text, code)	handles all types and sizes of files well
file contents are in the Git history and will be shared upon git/datalad push	file contents are in the annex. Not necessarily shared
Shared with every dataset clone	Can be kept private on a per-file level when sharing the dataset
Useful: Small, non-binary, frequently modified, need-to-be-accessible (DUA, README) files	Useful: Large files, private files

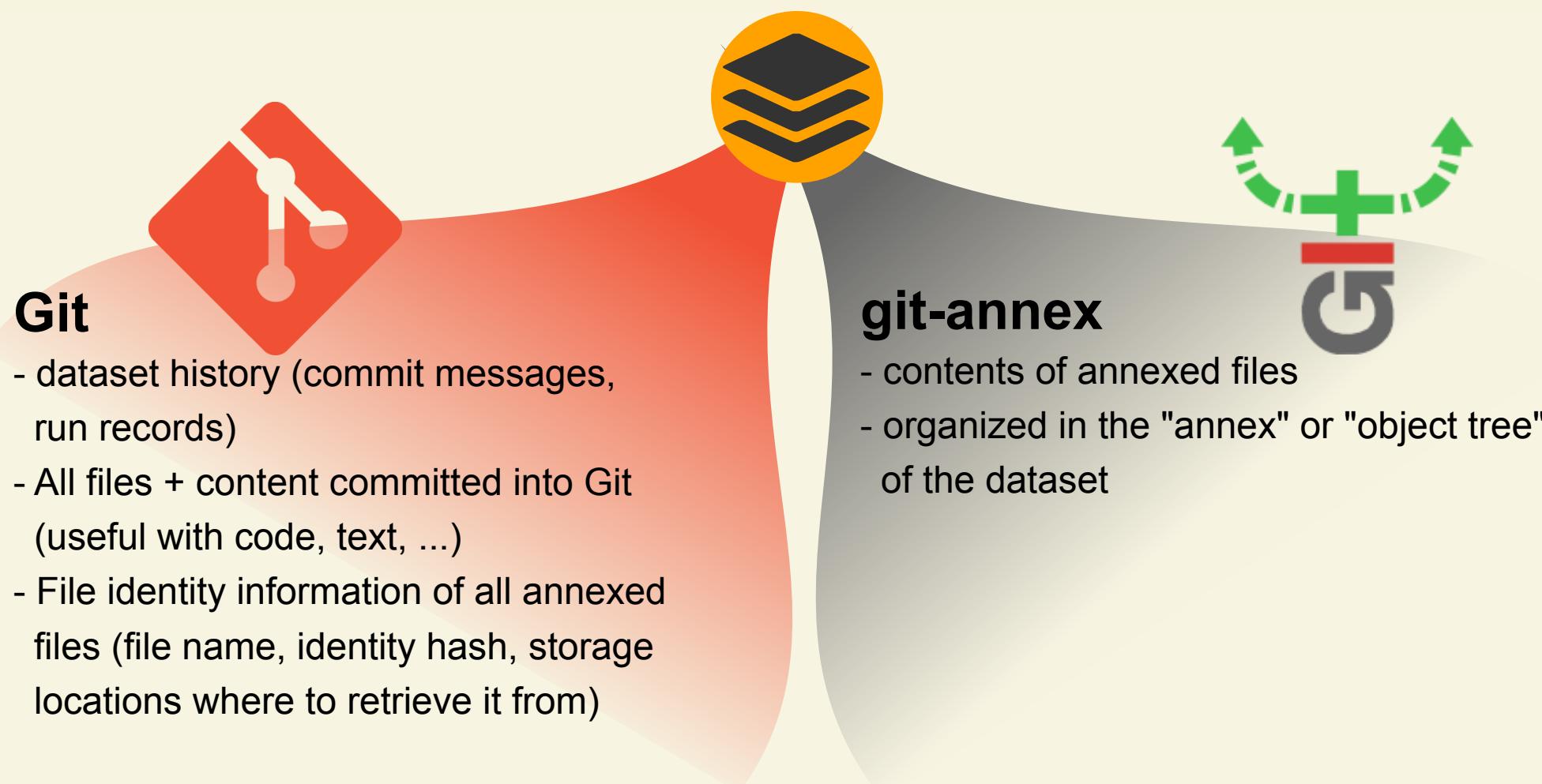
GIT VERSUS GIT-ANNEX

Useful background information for demo later. Read this [handbook chapter](#) for details

Git and Git-annex handle files differently: annexed files are stored in an annex. File content is hashed & only content-identity is committed to Git.

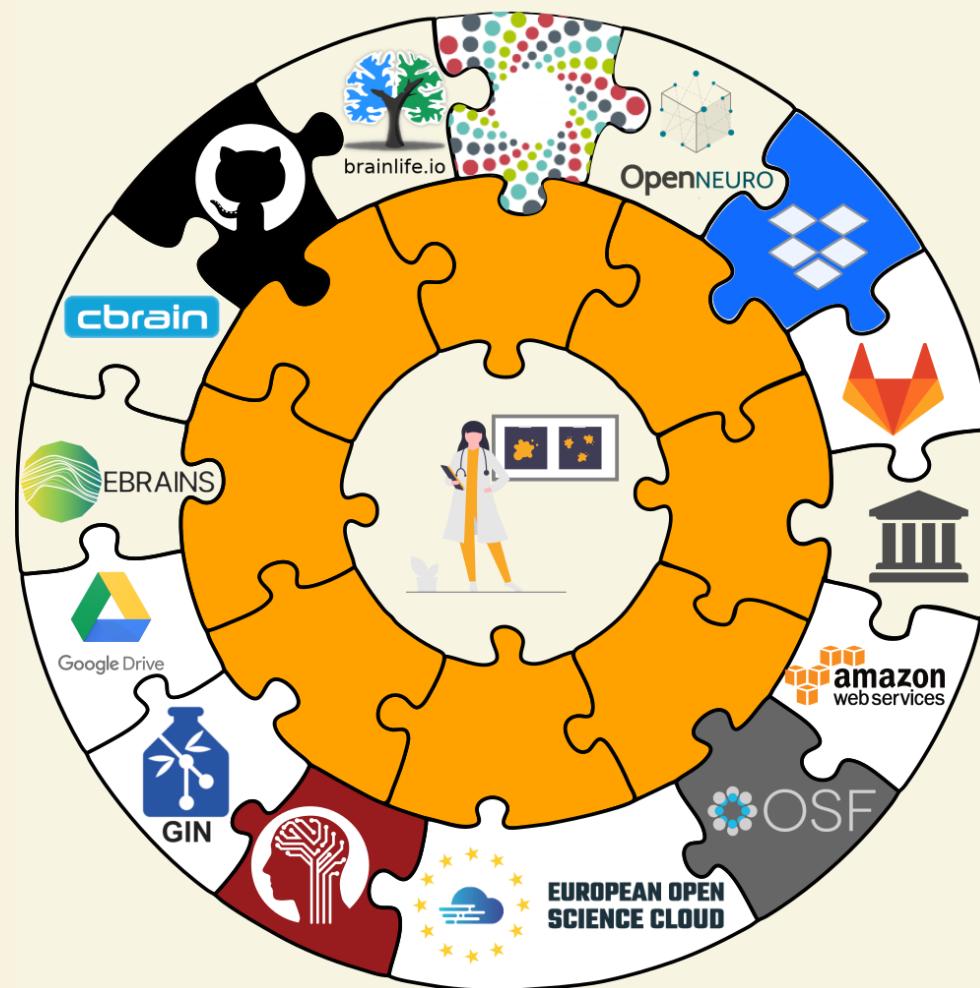
- Files stored in Git are modifiable, files stored in Git-annex are content-locked
- Annexed contents are not available right after cloning, only content identity and availability information (as they are stored in Git). Everything that is annexed needs to be retrieved with `datalad get` from wherever it is stored.

GIT VERSUS GIT-ANNEX



GIT VERSUS GIT-ANNEX

When sharing datasets with someone without access to the same computational infrastructure, annexed data is not necessarily stored together with the rest of the dataset (more in the **session on publishing**).



Transport logistics exist to interface with all major storage providers. If the one you use isn't supported, let us know!

GIT VERSUS GIT-ANNEX

Users can decide which files are annexed:

- Pre-made run-procedures, provided by DataLad (e.g., `text2git`, `yoda`) or created and shared by users ([Tutorial](#))
- Self-made configurations in `.gitattributes` (e.g., based on file type, file/path name, size, ...; [rules and examples](#))
- Per-command basis (e.g., via `datalad save --to-git`)

TEXT2G/TEXT VERSUS BINARY FILES

DISK-SPACE AWARE WORKFLOWS

- Clone the input data:

```
$ datalad clone git@github.com:datalad-datasets/machinelearning-books.git  
install(ok): /tmp/machinelearning-books (dataset)  
$ cd machinelearning-books && du -sh  
348K .
```

```
$ ls  
A.Shashua-Introduction_to_Machine_Learning.pdf  
B.Efron_T.Hastie-Computer_Age_Statistical_Inference.pdf  
C.E.Rasmussen_C.K.I.Williams-Gaussian_Processes_for_Machine_Learning.pdf  
D.Barber-Bayesian_Reasoning_and_Machine_Learning.pdf  
[...]
```

- retrieve annexed file's contents on demand:

```
$ datalad get A.Shashua-Introduction_to_Machine_Learning.pdf  
get(ok): /tmp/machinelearning-books/A.Shashua-Introduction_to_Machine_Learning.pdf (file) [from web...]
```

- Drop annexed file's contents when done:

```
$ datalad drop A.Shashua-Introduction_to_Machine_Learning.pdf  
drop(ok): /tmp/machinelearning-books/A.Shashua-Introduction_to_Machine_Learning.pdf (file) [checking https://]
```

DISTRIBUTED AVAILABILITY

- git-annex conceptualizes file availability information as a decentral network. A file can exist in multiple different locations. *git annex whereis* tells you which are known:

```
$ git annex whereis inputs/images/chinstrap_02.jpg
whereis inputs/images/chinstrap_02.jpg (1 copy)
  00000000-0000-0000-0000-000000000001 -- web
    c1bfc615-8c2b-4921-ab33-2918c0cbfc18 -- adina@muninn:/tmp/my-dataset [here]

  web: https://unsplash.com/photos/8PxCm4HsPX8/download?force=true
ok
```

- If a file has no other known storage locations, *drop* will warn
 - Here is a file with a registered remote location (the web)

```
$ datalad drop inputs/images/chinstrap_02.jpg
drop(ok): /home/my-dataset/inputs/images/chinstrap_02.jpg (file)
$ datalad get inputs/images/chinstrap_02.jpg
get(ok): inputs/images/chinstrap_02.jpg (file)
```

- Here is a file without a registered remote location (the web)

```
$ datalad drop inputs/images/chinstrap_01.jpg
drop(error): inputs/images/chinstrap_01.jpg (file)
  [unsafe; Could only verify the existence of 0 out of 1 necessary copy;
  (Use --reckless availability to override this check, or adjust numcopies.)]
```

- Delineation and advantages of decentral versus central RDM: In defense of decentralized research data management

DATA PROTECTION

Why are annexed contents write-protected? (part I)

- Where the filesystem allows it, annexed files are symlinks:

```
$ ls -l inputs/images/chinstrap_01.jpg
lrwxrwxrwx 1 adina adina 132 Apr  5 20:53 inputs/images/chinstrap_01.jpg -> ../../.git/annex/objects/1z/
xP/MD5E-s725496--2e043a5654cec96aadad554fda2a8b26.jpg/MD5E-s725496--2e043a5654cec96aadad554fda2a8b26.jpg
```

(PS: especially useful in datasets with many identical files)

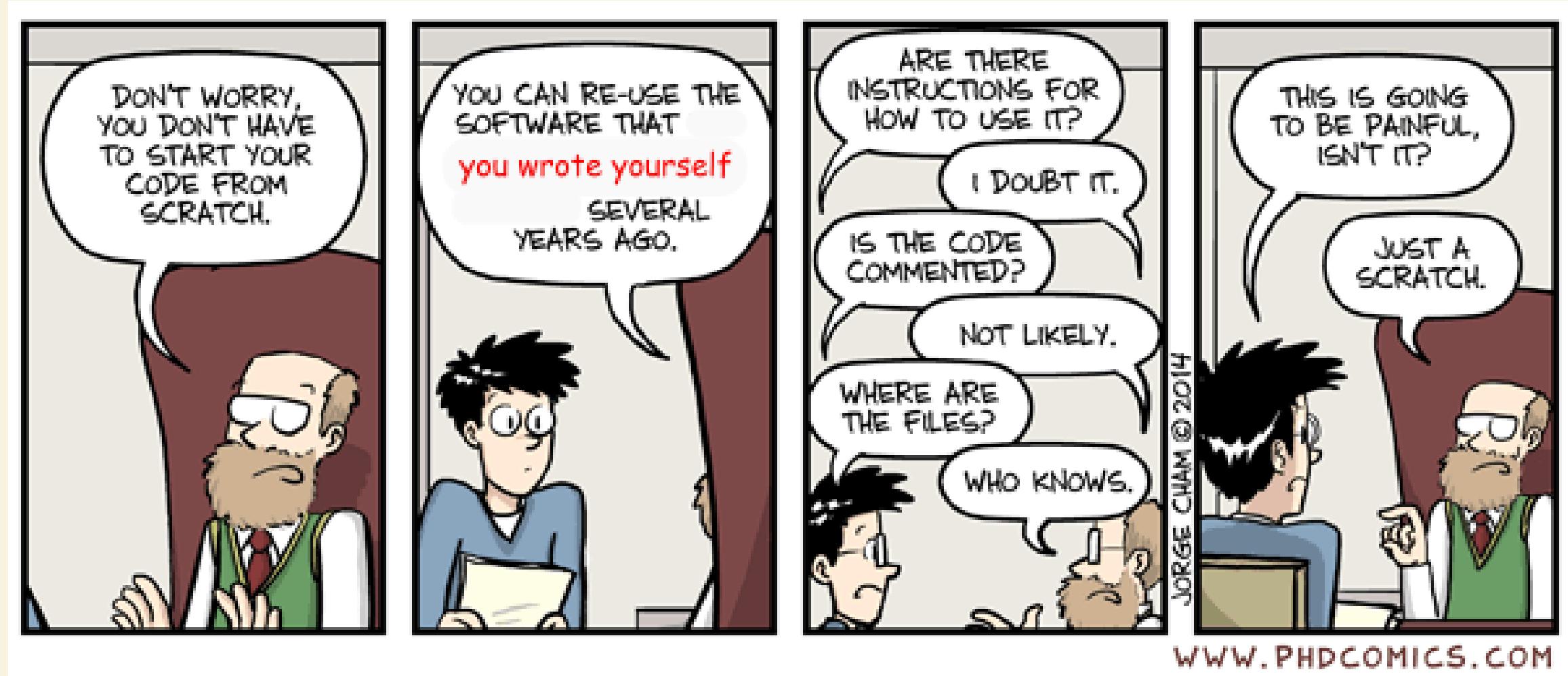
- The symlink reveals git-annex internal data organization based on identity hash:

```
$ md5sum inputs/images/chinstrap_01.jpg
2e043a5654cec96aadad554fda2a8b26  inputs/images/chinstrap_01.jpg
```

- git-annex write-protects files to keep this symlink functional - Changing file contents without git-annex knowing would make the hash change and the symlink point to nothing
- To (temporarily) remove the write-protection one can *unlock* the file

DETOUR & TEASER: REPRODUCIBLE DATA ANALYSIS

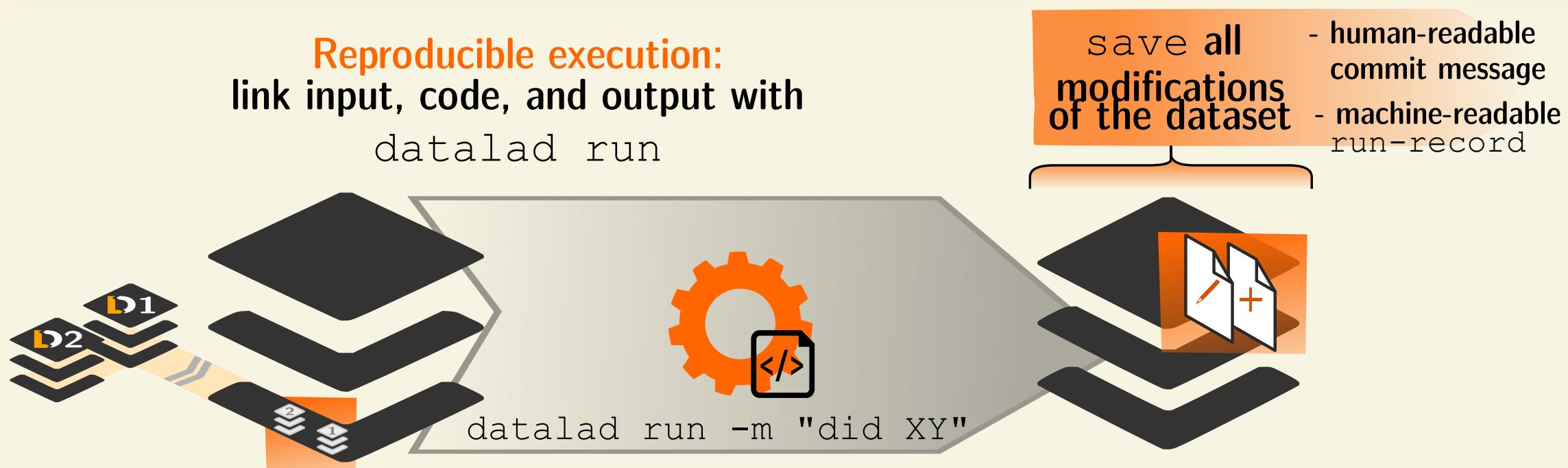
Your past self is the worst collaborator:



Code: psychoinformatics-de.github.io/rdm-course/01-content-tracking-with-datalad/index.html#data-processing

REPRODUCIBLE EXECUTION & PROVENANCE CAPTURE

datalad run wraps a command execution and records its impact on a dataset.



REPRODUCIBLE EXECUTION & PROVENANCE CAPTURE

datalad run wraps a command execution and records its impact on a dataset.

```
commit 9fbc0c18133aa07b215d81b808b0a83bf01b1984 (HEAD -> main)
Author: Adina Wagner [adina.wagner@t-online.de]
Date:   Mon Apr 18 12:31:47 2022 +0200

[DATALAD RUNCMD] Convert the second image to greyscale

==== Do not change lines below ====
{
  "chain": [],
  "cmd": "python code/greyscale.py inputs/images/chinstrap_02.jpg outputs/im>
  "dsid": "418420aa-7ab7-4832-a8f0-21107ff8cc74",
  "exit": 0,
  "extra_inputs": [],
  "inputs": [],
  "outputs": [],
  "pwd": "."
}
^^^ Do not change lines above ^^^

diff --git a/outputs/images_greyscale/chinstrap_02_grey.jpg b/outputs/images_gr>
new file mode 120000
index 000000..5febcb72
--- /dev/null
+++ b/outputs/images_greyscale/chinstrap_02_grey.jpg
@@ -0,0 +1 @@
+../../.git/annex/objects/19/mp/MD5E-s758168--8e840502b762b2e7a286fb5770f1ea69.>
\ No newline at end of file
```

The resulting commit's hash (or any other identifier) can be used to automatically re-execute a computation (more on this tomorrow)

DATA PROTECTION

Why are annexed contents write-protected? (part 2)

- When you try to modify an annexed file without unlocking you will see "Permission denied" errors.

```
Traceback (most recent call last):
  File "/home/bob/Documents/rdm-warmup/example-dataset/code/greyscale.py", line 20, in module
    grey.save(args.output_file)
  File "/home/bob/Documents/rdm-temporary/venv/lib/python3.9/site-packages/PIL/Image.py", line 2232, in save
    fp = builtins.open(filename, "w+b")
PermissionError: [Errno 13] Permission denied: 'outputs/images_greyscale/chinstrap_02_grey.jpg'
```

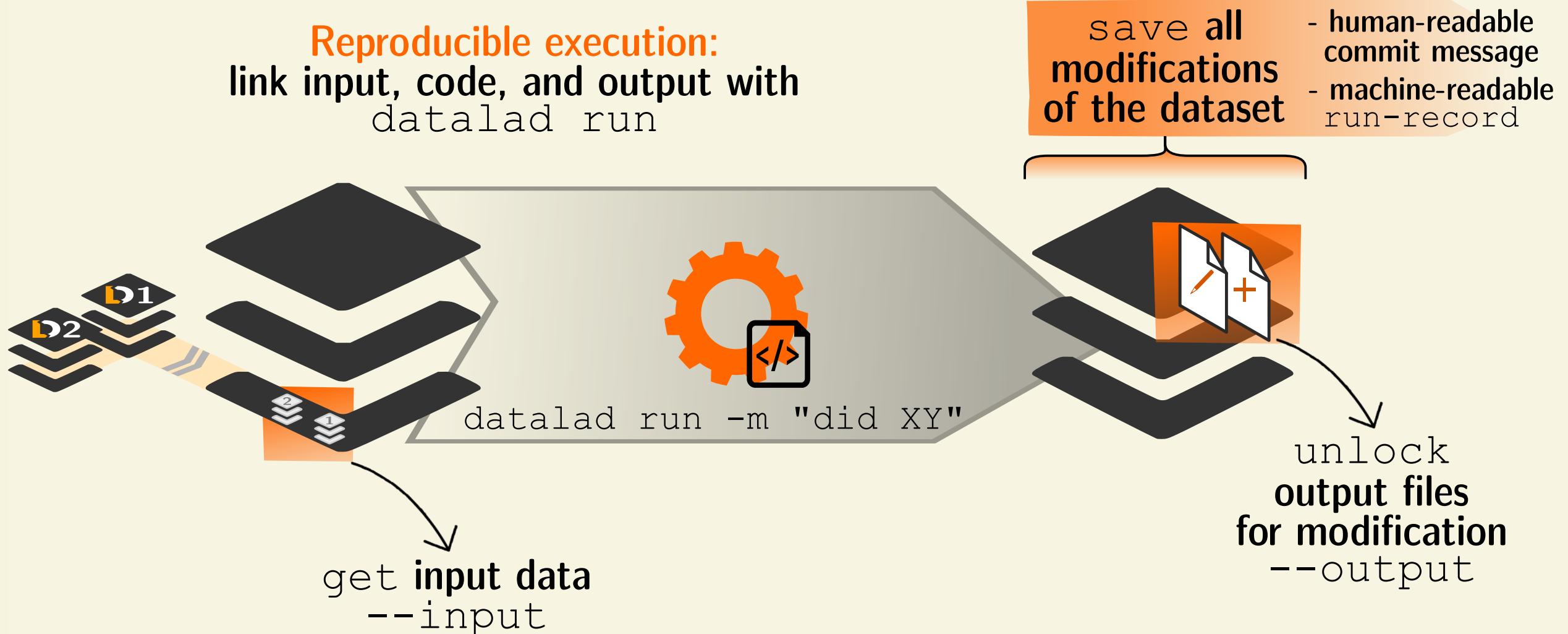
- Use *datalad unlock* to make the file modifiable. Underneath the hood (given the file system initially supported symlinks), this removes the symlink:

```
$ datalad unlock outputs/images_greyscale/chinstrap_02_grey.jpg
$ ls outputs/images_greyscale/chinstrap_02_grey.jpg
-rw-r--r-- 1 adina adina 758168 Apr 18 12:31 outputs/images_greyscale/chinstrap_02_grey.jpg
```

- *datalad save* locks the file again. Locking and unlocking ensures that git-annex always finds the right version of a file.

REPRODUCIBLE EXECUTION & PROVENANCE CAPTURE

datalad run wraps a command execution and records its impact on a dataset.
In addition, it can take care of data retrieval and unlocking



DATALAD RERUN

- `datalad rerun` is helpful to spare others and yourself the short- or long-term memory task, or the forensic skills to figure out how you performed an analysis
- But it is also a digital and machine-readable provenance record
- Important: The better the run command is specified, the better the provenance record
- Note: `run` and `rerun` only create an entry in the history if the command execution leads to a change.
- Task: Use `datalad rerun` to rerun the script execution. Find out if the output changed

SUMMARY - UNDERNEATH THE HOOD

Files are either kept in Git or in git-annex.

datalad save is used for both, but configurations (e.g., *text2git*), dataset rules (e.g., in a *.gitattributes* file, or flags change the default behavior of annexing everything

Annexed files behave differently from files kept in Git:

They can be retrieved and dropped from local or remote locations, they are write-protected, their content is unknown to Git (and thus easy to keep private).

***datalad clone* installs datasets from URLs or local or remote paths**

Annexed files contents can be retrieved or dropped on demand, file contents of files stored in Git are available right away.

***datalad unlock* makes annexed files modifiable, *datalad save* locks them again.**

(It is generally easier to get accidentally saved files out of the annex than out of Git - see handbook.datalad.org/basics/101-136-filesystem.html for examples)

***datalad run* records the impact of any command execution in a dataset.**

Data/directories specified as `--input` are retrieved prior to command execution, data/directories specified as `--output` unlocked.

***datalad rerun* can automatically re-execute run-records later.**

They can be identified with any commit-ish (hash, tag, range, ...)

QUESTIONS!

Awkward silence can be bridged with awkward MC questions :)

BEFORE WE CONTINUE...

Let your energy level define how we progress: