

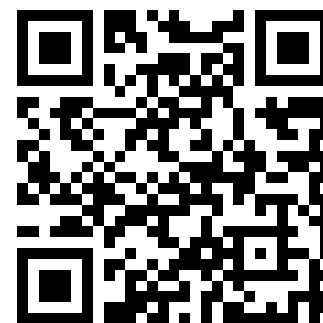
# MANAGE YOUR DATA WITH DATALAD:

## FROM LOCAL VERSION CONTROL TO DATA PUBLICATION

Adina Wagner

 [mas.to/@adswa](https://mas.to/@adswa)

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



PDF: DOI [10.5281/zenodo.8431164](https://doi.org/10.5281/zenodo.8431164) (Scan the QR code)  
Rendering: [files.inm7.de/adina/talks/html/cuttinggardens.html](https://files.inm7.de/adina/talks/html/cuttinggardens.html)

# LOGISTICS

- Collaborative, public notes, networking, & anonymous questions at [etherpad.wikimedia.org/p/datalad@cuttinggardens](https://etherpad.wikimedia.org/p/datalad@cuttinggardens)
- We are using a JupyterHub at [datalad-hub.inm7.de](https://datalad-hub.inm7.de). Draw a username! You can log in with a password of your choice (at least 8 characters).
- Format:
  - Mostly hands-on: Watch me live-code, and try out the software yourself in the browser. Conceptual wrap-up at the end.
  - Ask questions any time

# FURTHER RESOURCES AND STAY IN TOUCH

If you have questions after the workshop...

**Reach out to to the DataLad team via**

- **Matrix** (free, decentralized communication app, no app needed). We run a weekly Zoom office hour (Tuesday, 4pm Berlin time) from this room as well.
- **The development repository on GitHub**

**Reach out to the (Neuro-) user community with**

- A question on **neurostars.org** with a `datalad` tag

**Find more user tutorials or workshop recordings**

- On **DataLad's YouTube channel**
- In the **DataLad Handbook**
- In the **DataLad RDM course**
- In the **Official API documentation**
- In an overview of most tutorials, talks, videos at **github.com/datalad/tutorials**

# ACKNOWLEDGEMENTS

## DataLad software & ecosystem

- Psychoinformatics Lab,  
Research center Jülich
- Center for Open  
Neuroscience,  
Dartmouth College
- Joey Hess (git-annex)
- *>100 additional contributors*

## Funders



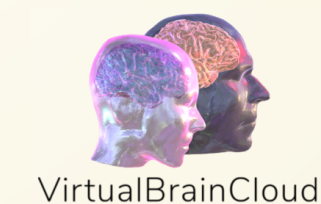
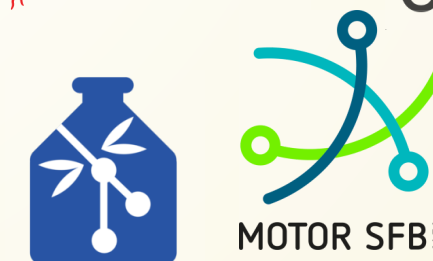
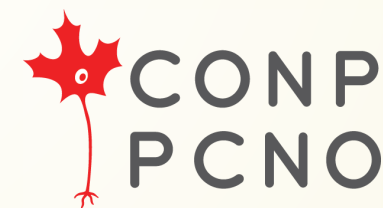
NSF 1429999



BMBF 01GQ1411



## Collaborators





# DATALAD USECASES

psychoinformatics-de / paper-remodnav

Public

Edit Pins 11 Unwatch 2 Fork 3 Star

Code Issues Pull requests Actions Projects Wiki Security Insights

master 6 branches 2 tags Go to file Add file Code

adswa Merge pull request #19 from psychoinformatics-de/reproducibility-im... 4edccb8 on Sep 14, 2021 451 commits

|                    |  |               |
|--------------------|--|---------------|
| code               | Generate figures with deterministic metadata for comp. reproducibility   | 10 months ago |
| data               | Point to latest label dataset  | 3 years ago   |
| img                | Include original SVGs into the repo to allow immediate manuscript ren... | 10 months ago |
| remodnav @ d289118 | Update remodnav with latest test dataset                                 | 3 years ago   |
| .gitignore         | Prevent permanent rebuilds of the figures                                | 3 years ago   |
| .gitmodules        | [DATALAD] modified subdataset properties                                 | 16 months ago |
| COPYING            | Declare CC-BY license  | 3 years ago   |
| Makefile           | Specify Python package versions in the state of final submission June... | 10 months ago |
| README.md          | DOC: improve readme, differentiate between recompile and recompute       | 10 months ago |

About

Code, data and manuscript for <https://doi.org/10.1101/619254>

Readme CC-BY-4.0 license 3 stars 11 watching 2 forks

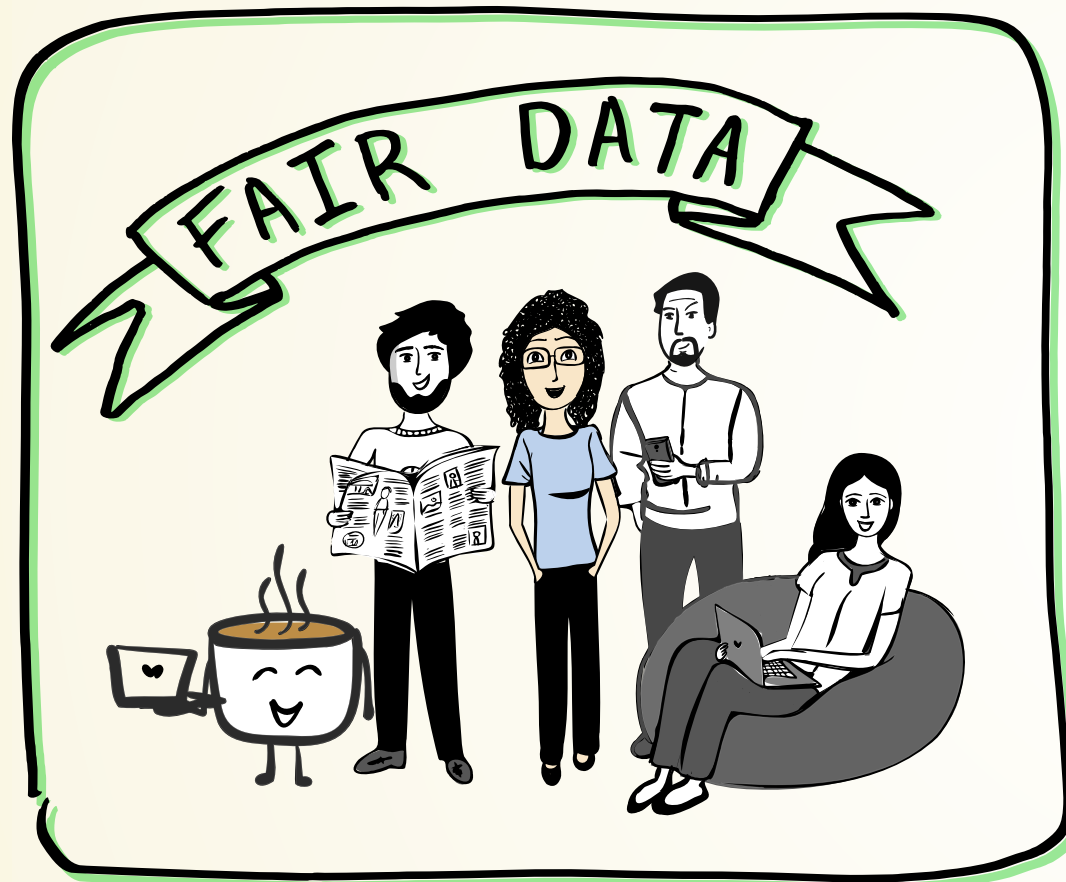
Releases

2 tags

Create a new release

# A COMMON USECASE

## ● THE TEAM



- Alice is a PhD student in a research team.
- She works on a fairly typical research project: Data collection & processing.
- First sample → final result = complex process

**HOW DOES ALICE GO ABOUT HER DAILY JOB?**

# A COMMON USECASE

- In her project, Alice likes to have an automated record of:
  - when a given file was last changed
  - where it came from
  - what input files were used to generate a given output
  - why some things were done.
- Even if she doesn't share her work, this is essential for her future self
- Her project is exploratory: Frequent changes to her analysis scripts
- She enjoys the comfort of being able to return to a previously recorded state

**THIS IS: \*LOCAL VERSION CONTROL\***

# A COMMON USECASE

- Alice's work is not confined to a single computer:
  - Laptop / desktop / remote server / dedicated back-up
  - Alice wants to automatically & efficiently synchronize
- Parts of the data are collected or analyzed by colleagues. This requires:
  - distributed synchronization with centralized storage
  - preservation of origin & authorship of changes
  - effective combination of simultaneous contributions

**THIS IS: \*DISTRIBUTED VERSION CONTROL\***

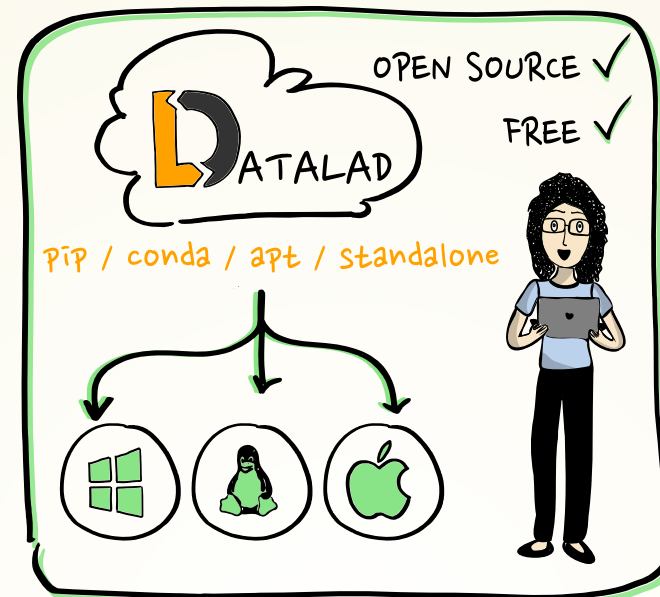
# A COMMON USECASE

- Alice applies local version control for her own work, and reproducibly records it
- She also applies distributed version control when working with colleagues and collaborators
- She often needs to work on a subset of data at any given time:
  - all files are kept on a server
  - a few files are rotated into and out of her laptop
- Alice wants to publish the data at project's end:
  - raw data / outputs / both
  - completely or selectively

**THIS IS: \*DATA MANAGEMENT (WITH DATALAD 😊)\***

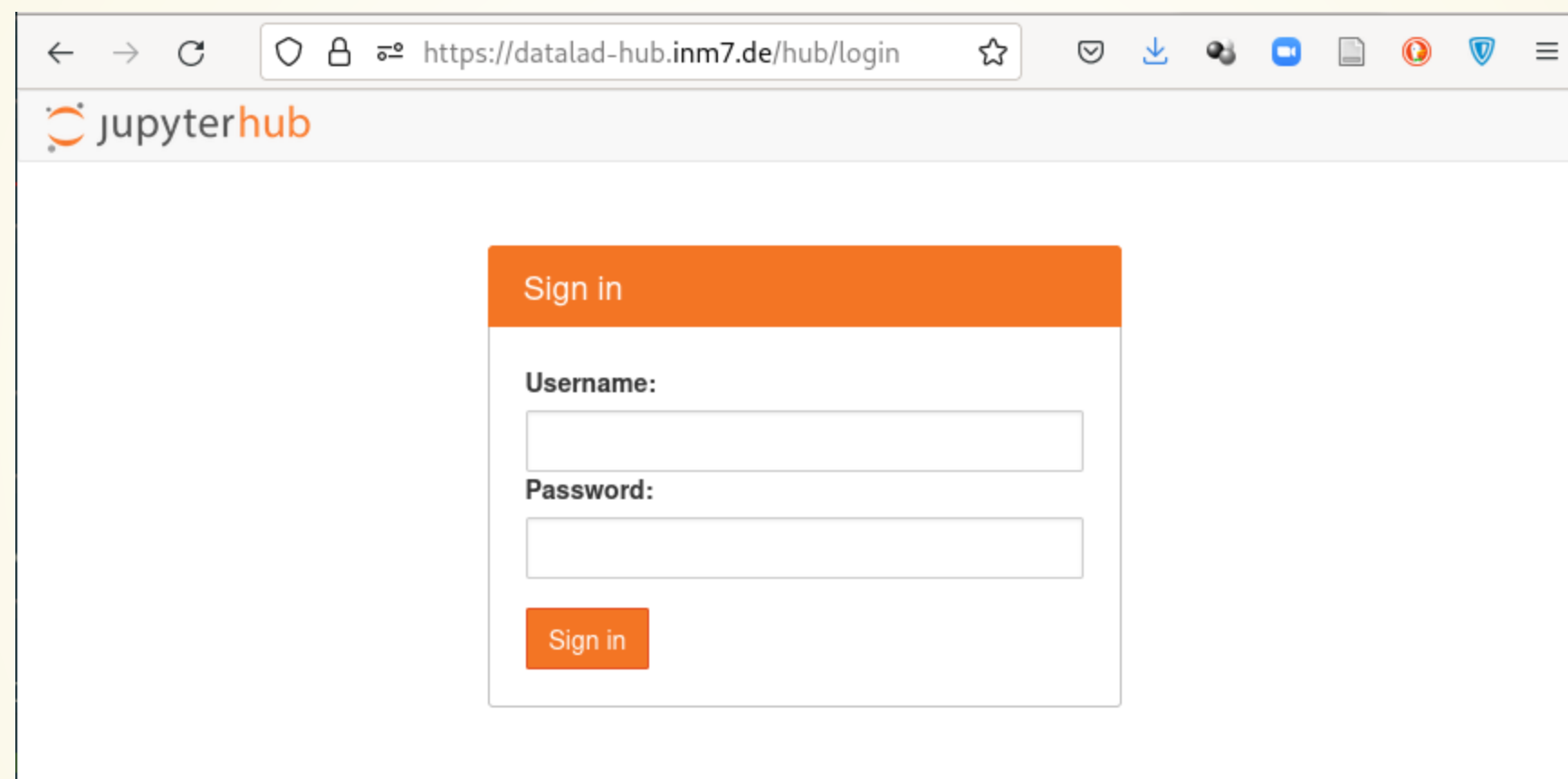


# DATALAD



- Domain-agnostic **command-line tool** (+ graphical user interface), built on top of **Git** & **Git-annex**
- Major features:
  - Version-controlling arbitrarily large content**  
Version control data & software alongside to code!
  - Transport mechanisms for sharing & obtaining data**  
Consume & collaborate on data (analyses) like software
  - (Computationally) reproducible data analysis**  
Track and share provenance of all digital objects
  - (... and much more)

# LET'S TRY IT OUT



The screenshot shows a web browser window with the address bar displaying `https://datalad-hub.inm7.de/hub/login`. The page features the JupyterHub logo at the top left. In the center, there is a login form with an orange header bar labeled "Sign in". Below this, the form contains two input fields: "Username:" and "Password:". At the bottom of the form is an orange button labeled "Sign in".

`datalad-hub.inm7.de`

**username:**

what you draw from the jaw

**password:**

Set at first login, at least 8 characters

**Important!** The Hub is a shared resource. Don't fill it up :)

# GIT IDENTITY SETUP

Check Git identity:

```
git config --get user.name  
git config --get user.email
```

copy

Configure Git identity:

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

copy

Configure DataLad to use latest features:

```
git config --global --add datalad.extensions.load next
```

copy

## USING DATALAD IN A TERMINAL

Check the installed version:

```
datalad --version
```

[copy](#)

For help on using DataLad from the command line:

```
datalad --help
```

[copy](#)

The help may be displayed in a pager - exit it by pressing "q"

For extensive info about the installed package, its dependencies, and extensions, use `datalad wtf`. Let's find out what kind of system we're on:

```
datalad wtf -S system
```

[copy](#)

# USING DATALAD VIA ITS PYTHON API

Open a Python environment:

```
ipython
```

[copy](#)

Import and start using:

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

[copy](#)

Exit the Python environment:

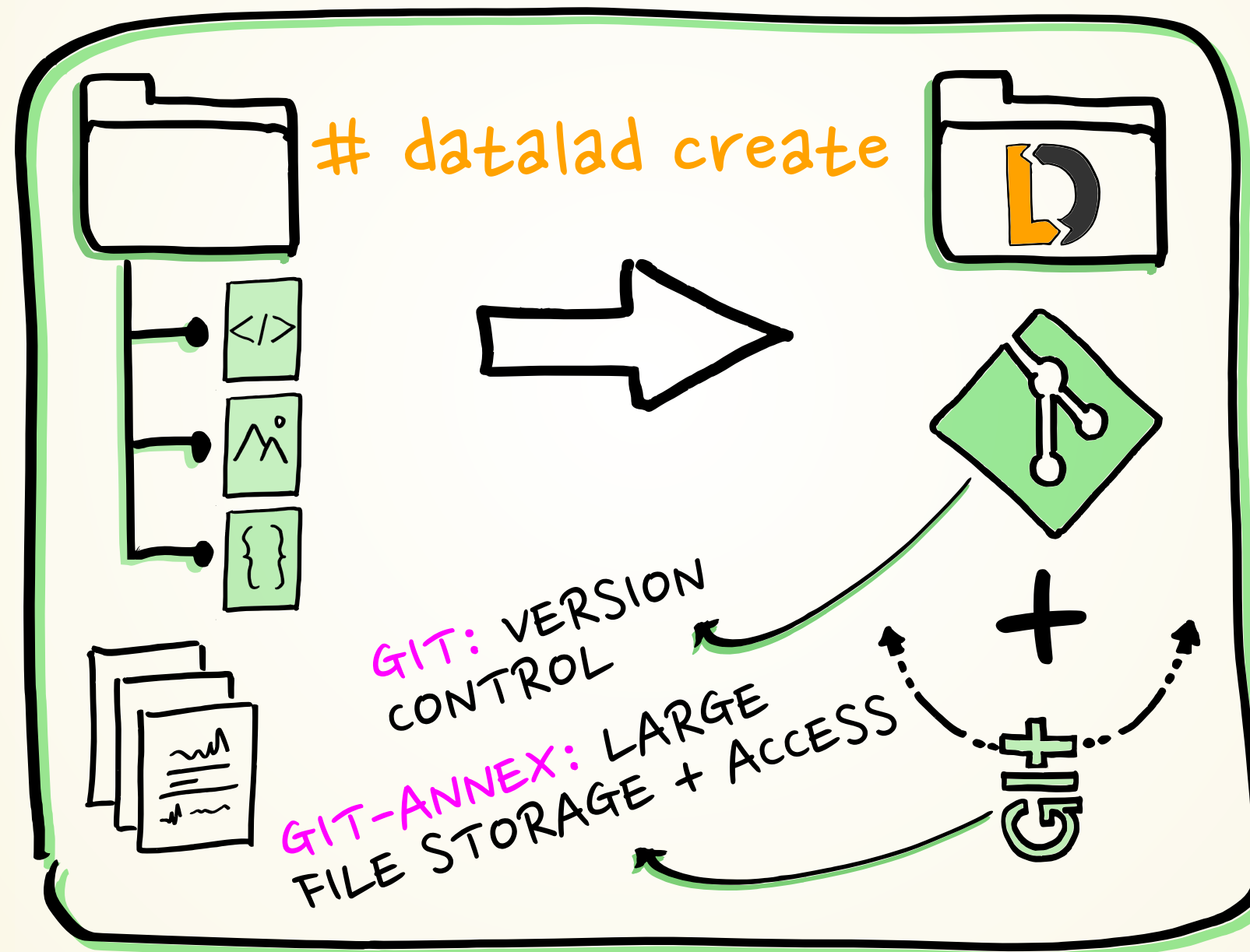
```
exit
```

[copy](#)



# DATALAD DATASETS...

## ● THE DATALAD DATASET



## ...DATALAD DATASETS

Create a dataset (here, with the yoda configuration, which adds a helpful structure and configuration for data analyses):



```
datalad create -c yoda my-analysis
```

[copy](#)

Let's have a look inside. Navigate using cd (change directory):

```
cd my-analysis
```

[copy](#)

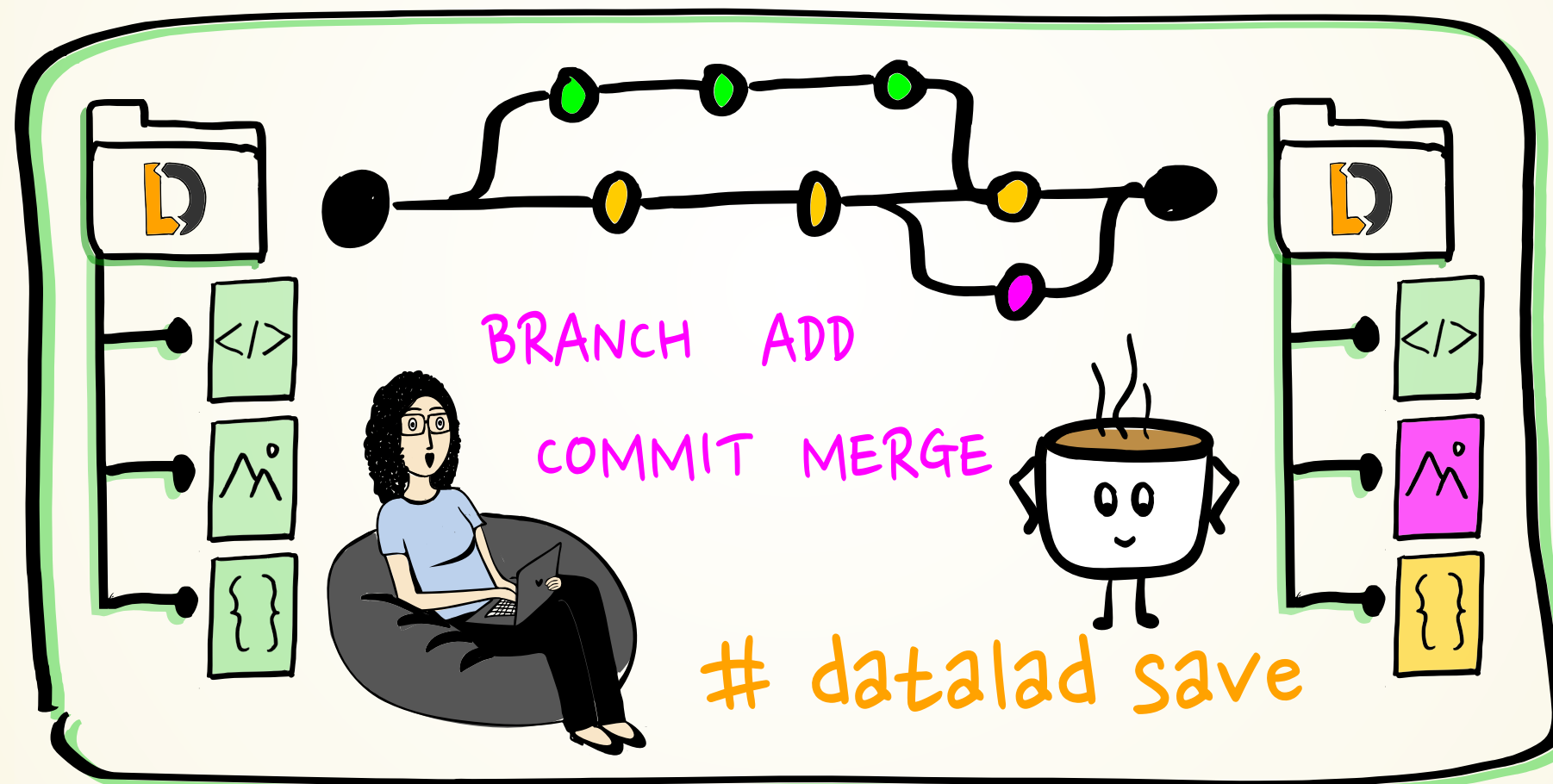
List the directory content, including hidden files, with ls:

```
ls -la .
```

[copy](#)

# VERSION CONTROL...

## ● VERSION CONTROL WITH GIT



## ...VERSION CONTROL

The yoda-configuration added a README placeholder in the dataset. Let's add Markdown text (a project title) to it:

```
echo "# My example DataLad dataset" >| README.md
```

[copy](#)

Now we can check the status of the dataset:

```
datalad status
```

[copy](#)

We can save the state with save

```
datalad save -m "Add project title into the README"
```

[copy](#)

Further modifications:

```
echo "Contains a small data analysis for my project" >> README.md
```

[copy](#)

You can also checkout what has changed:

```
git diff
```

[copy](#)

Save again:

```
datalad save -m "Add information on the dataset contents to the README"
```

[copy](#)

## ...VERSION CONTROL

Now, let's check the dataset history:

```
git log
```

[copy](#)

We can also make the history prettier:

```
tig
```

[copy](#)

(navigate with arrow keys and enter, press "q" to go back and exit the program)

Convenience functions make downloads easier. Let's add code for a data analysis from an external source:

```
datalad download-url -m "Add an analysis script" \  
-O code/mne_time_frequency_tutorial.py \  
https://raw.githubusercontent.com/datalad-handbook/resources/master/mne_time_frequency_tutorial.py
```

[copy](#)

Check out the file's history:

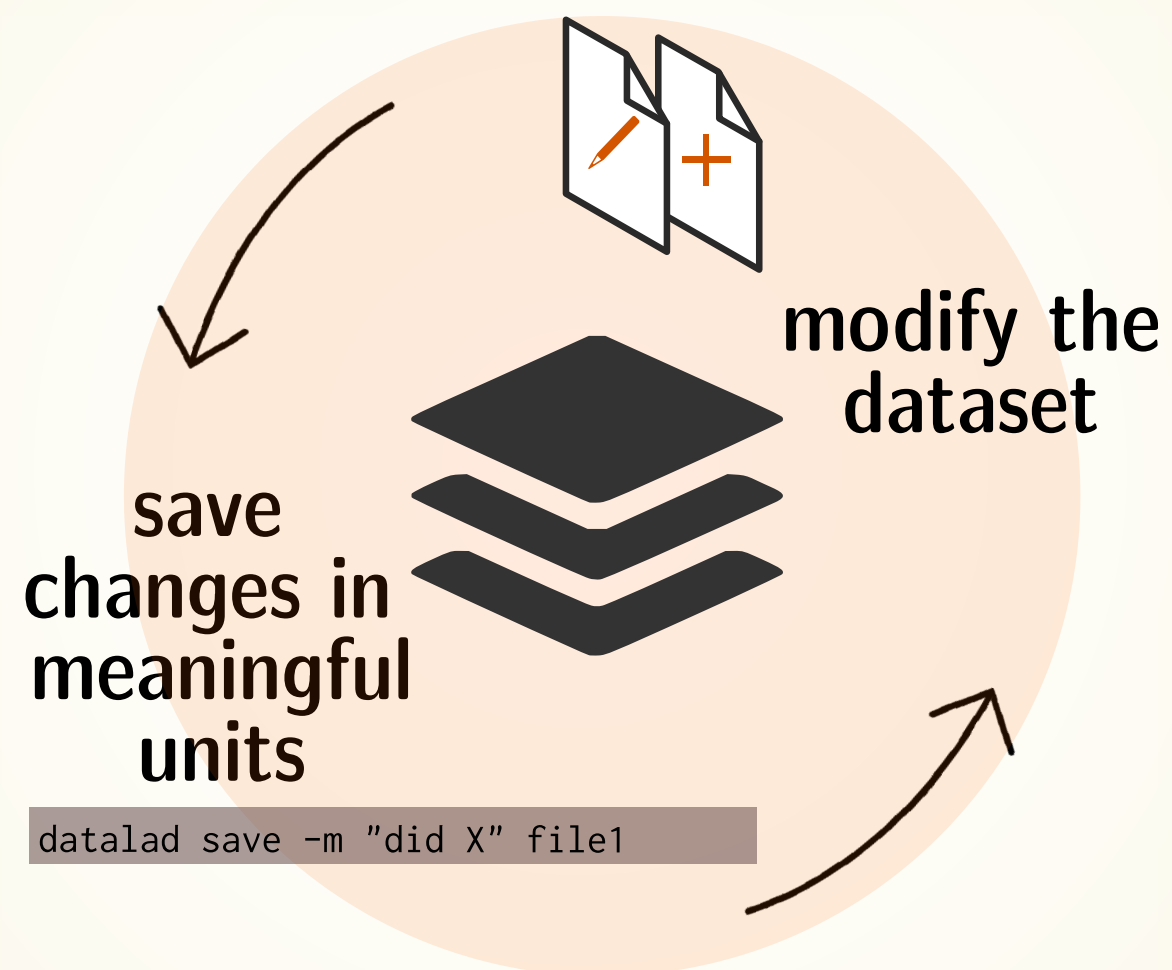
```
git log code/mne_time_frequency_tutorial.py
```

[copy](#)



# LOCAL VERSION CONTROL

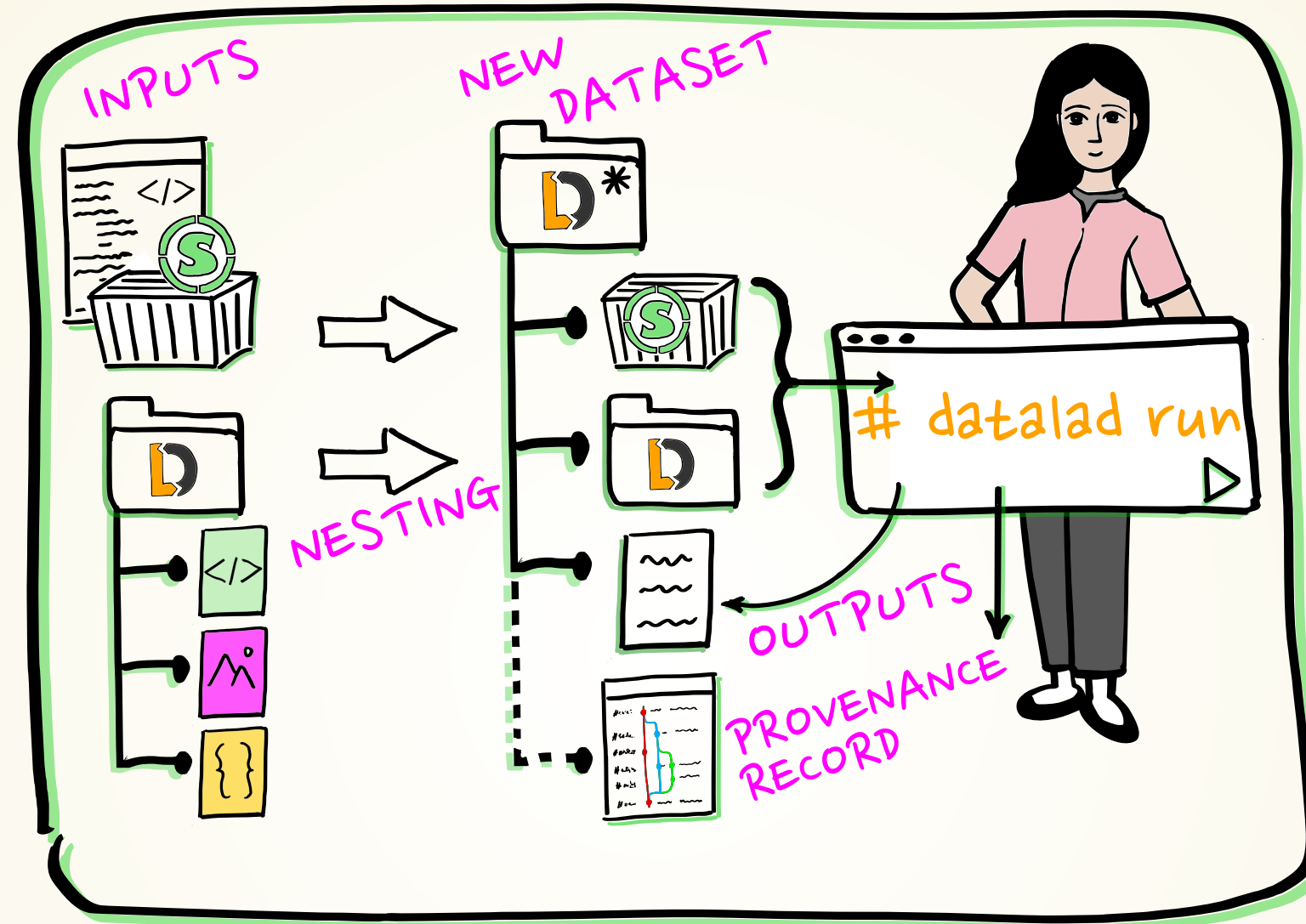
Procedurally, version control is easy with DataLad!



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

# COMPUTATIONALLY REPRODUCIBLE EXECUTION I...

## ● EXACT DEPENDENCIES+PROVENANCE



- which script/pipeline version
- was run on which version of the data
- to produce which version of the results?

## ... COMPUTATIONALLY REPRODUCIBLE EXECUTION I

A variety of processes can modify files. A simple example: Code formatting

```
black code/mne_time_frequency_tutorial.py
```

[copy](#)

Version control makes changes transparent:

```
git diff
```

[copy](#)

But its useful to keep track beyond that. Let's discard the latest changes...

```
git restore code/mne_time_frequency_tutorial.py
```

[copy](#)

... and record precisely what we did

```
datalad run -m "Reformat code with black" \  
"black code/mne_time_frequency_tutorial.py"
```

[copy](#)

let's take a look:

```
git show
```

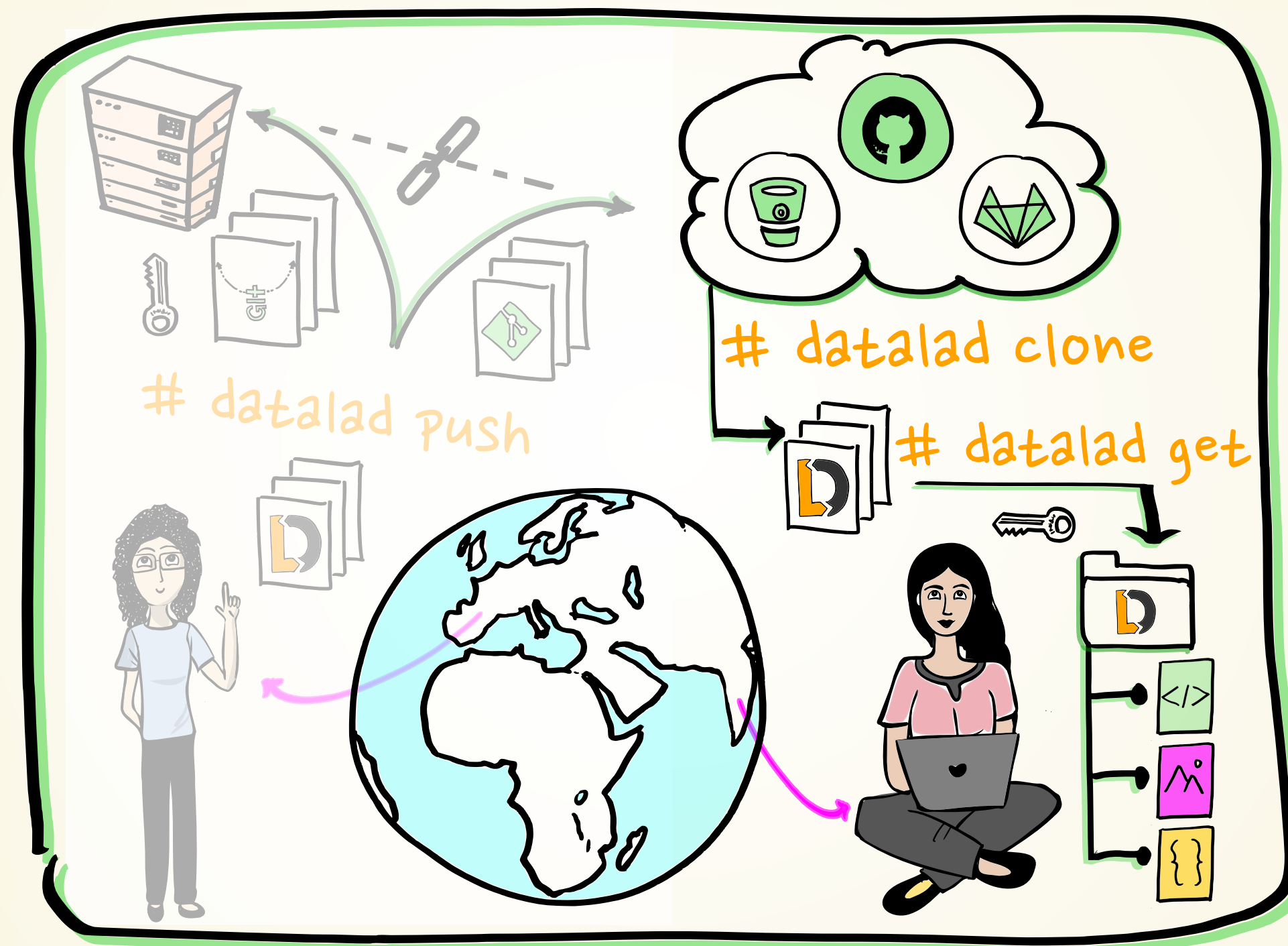
[copy](#)

... and repeat!

```
datalad rerun
```

[copy](#)

# DATA CONSUMPTION & TRANSPORT...



# ...DATA CONSUMPTION & TRANSPORT...

You can install a dataset from remote URL (or local path) using `clone`. Either as a stand-alone entity:

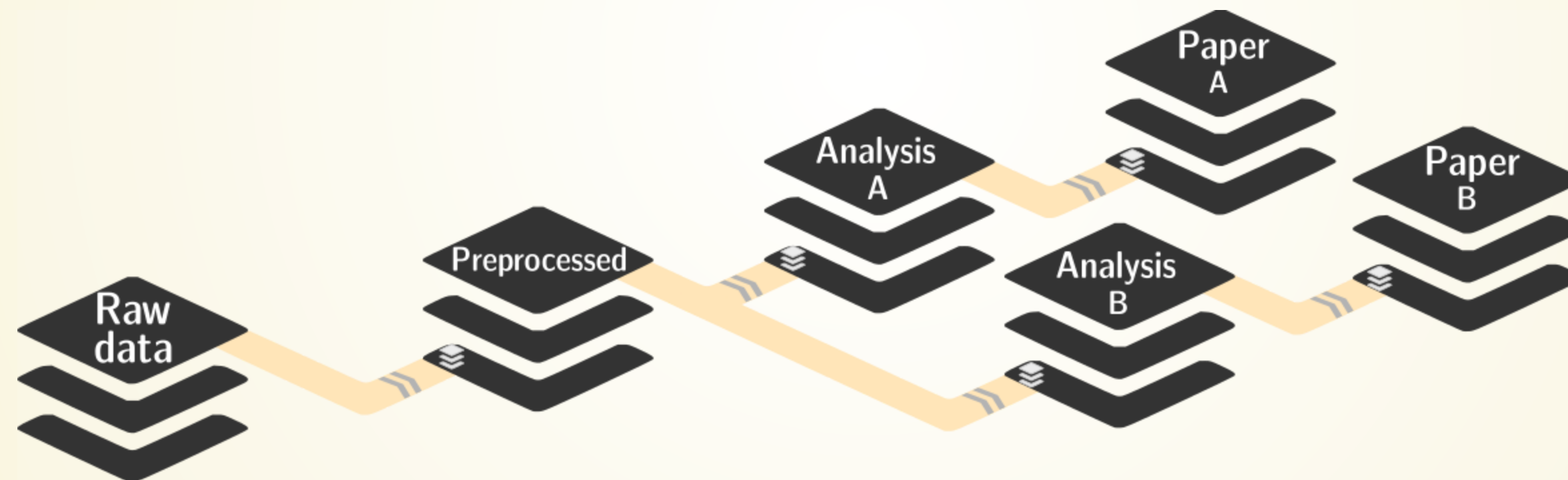
```
# just an example:  
datalad clone \  
https://github.com/psychoinformatics-de/studyforrest-data-phase2.git
```

copy

Or as linked dataset, nested in another dataset in a superdataset-subdataset hierarchy:

```
# just an example:  
datalad clone -d . \  
https://github.com/psychoinformatics-de/studyforrest-data-phase2.git
```

copy



- Helps with scaling (see e.g. the [Human Connectome Project dataset](#) )
- Version control tools struggle with >100k files
- Modular units improves intuitive structure and reuse potential
- Versioned linkage of inputs for reproducibility



## ...DATASET NESTING

Let's make a nest!

Clone **dataset from OpenNeuro** with MEG example data into a specific location ("input/") in the existing dataset, making it a *subdataset*:

```
datalad clone --dataset . \
  https://github.com/OpenNeuroDatasets/ds003104.git \
  input/
```

copy

Let's see what changed in the dataset, using the `subdatasets` command:

```
datalad subdatasets
```

copy

... and also `git show`:

```
git show
```

copy

We can now view the cloned dataset's file tree:

```
cd input  
ls
```

[copy](#)

...and also its history

```
tig
```

[copy](#)

Let's check the dataset size (with the du disk-usage command):

```
du -sh
```

[copy](#)

Let's check the *actual* dataset size:

```
datalad status --annex
```

[copy](#)

## ...DATA CONSUMPTION & TRANSPORT

We can retrieve actual file content with get:

```
datalad get sub-01
```

[copy](#)

If we don't need a file locally anymore, we can drop its content:

```
datalad drop sub-01
```

[copy](#)

No need to store all files locally, or archive results with Giga/Terra-Bytes of source data:

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

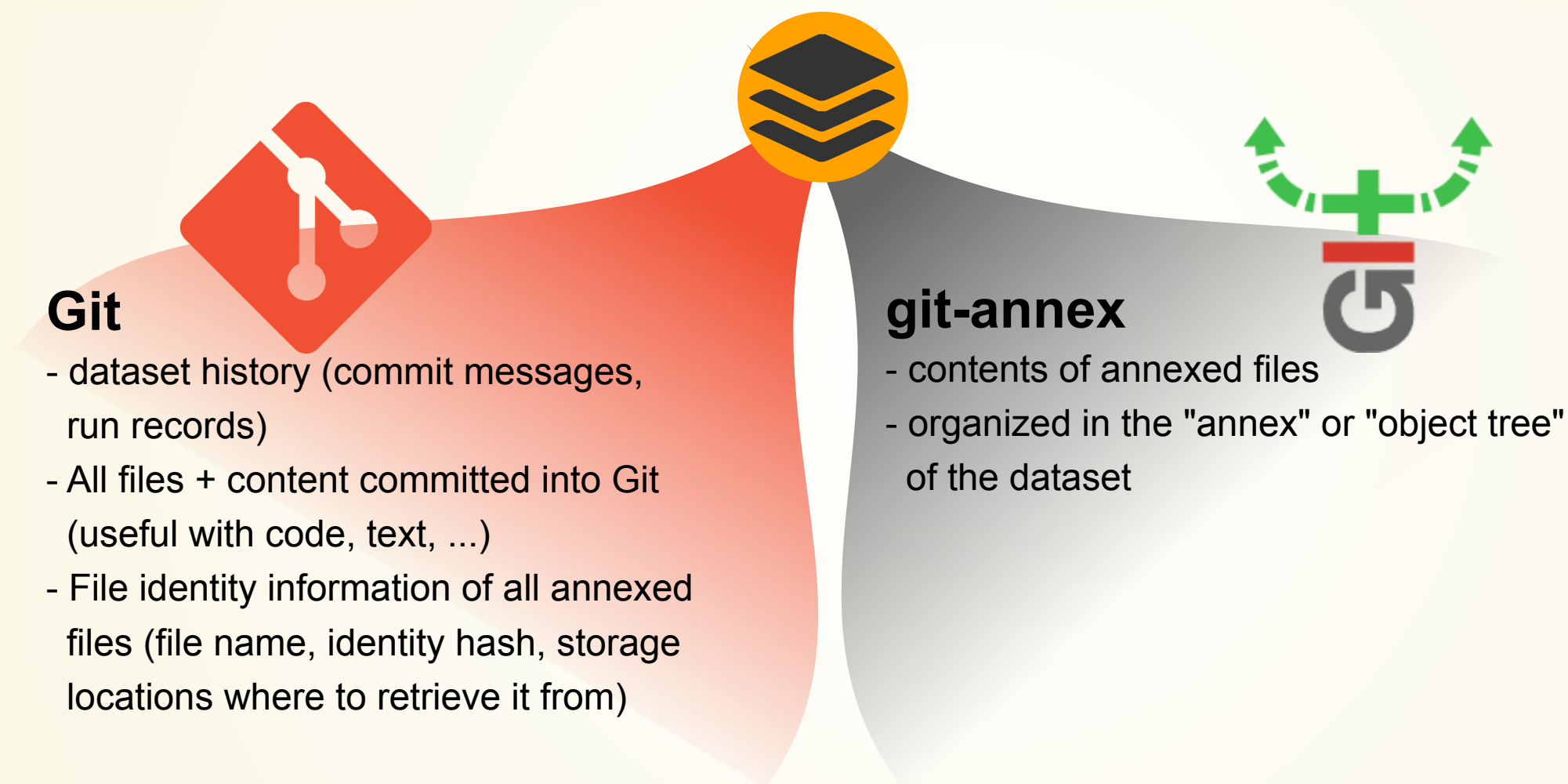
[copy](#)

If data is published anywhere, your data analysis can carry an actionable link to it, with barely any space requirements.

# 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



- With annexed data, only content identity (hash) and location information is put into Git, rather than file content. The annex, and transport to and from it is managed with **git-annex**

# GIT VERSUS GIT-ANNEX

Configurations (e.g., YODA), custom **rules**, or command parametrization determines if a file is annexed

Storing files in Git or git-annex has distinct advantages:

| Git   | git-annex   |
|---|---|
| handles <b>small</b> files well (text, code)  | handles <b>all</b> types and sizes of files well                        |
| file contents are in the Git history and will be <b>shared</b> upon git/datalad push      | file contents are in the annex. Not necessarily shared                  |
| Shared with every dataset clone   | <b>Can be kept private</b> 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                                      |

YODA configures the contents of the code / directory and the dataset descriptions (e.g., README files) to be in Git. There are many other configurations, and you can also **write your own**.



## ...COMPUTATIONALLY REPRODUCIBLE EXECUTION...

Try to execute the downloaded analysis script. Does it work?

```
cd ..  
python code/mne_time_frequency_tutorial.py
```

copy

- Software can be difficult or impossible to install (e.g. conflicts with existing software, or on HPC) for you or your collaborators
- Different software versions/operating systems can produce different results:  
[Glatard et al., doi.org/10.3389/fninf.2015.00012](https://doi.org/10.3389/fninf.2015.00012)
- **Software containers** encapsulate a software environment and isolate it from a surrounding operating system. Two common solutions: Docker, Singularity



## ...COMPUTATIONALLY REPRODUCIBLE EXECUTION...

- The `datalad run` can run any command in a way that links the command or script to the results it produces and the data it was computed from
- The `datalad rerun` can take this recorded provenance and recompute the command
- The `datalad containers-run` (from the extension "datalad-container") can capture software provenance in the form of software containers in addition to the provenance that `datalad run` captures

## ...COMPUTATIONALLY REPRODUCIBLE EXECUTION

With the `data-lad-container` extension, we can add software containers to datasets and work with them. Let's add a software container with Python software to run the script

```
data-lad containers-add python-env \  
--url https://files.inm7.de/adina/resources/mne \  
--call-fmt "singularity exec {img} {cmd}"
```

[copy](#)

inspect the list of registered containers:

```
data-lad containers-list
```

[copy](#)

Now, let's try out the `containers-run` command:

```
data-lad containers-run -m "run classification analysis in python environment" \  
--container-name python-env \  
--input "input/sub-01/meg/sub-01_task-somato_meg.fif" \  
--output "figures/*" \  
"python3 code/mne_time_frequency_tutorial.py {inputs}"
```

[copy](#)

What changed after the `containers-run` command has completed?

We can use `data-lad diff` (based on `git diff`):

```
data-lad diff -f HEAD~1
```

[copy](#)

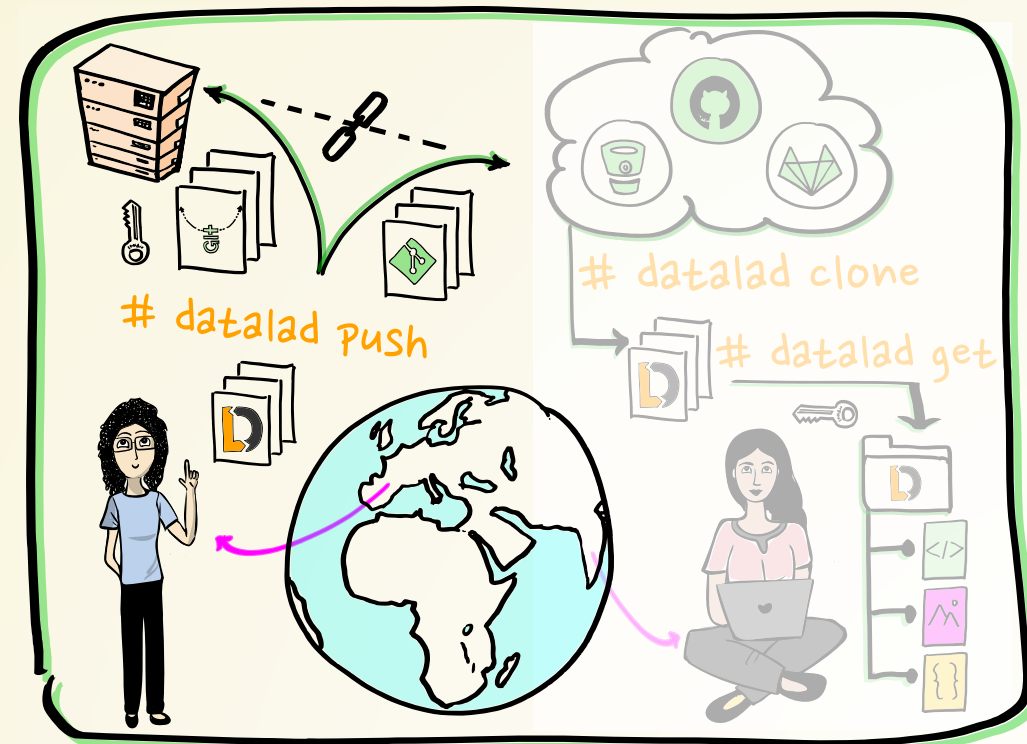
We see that some files were added to the dataset!

And we have a complete provenance record as part of the git history:

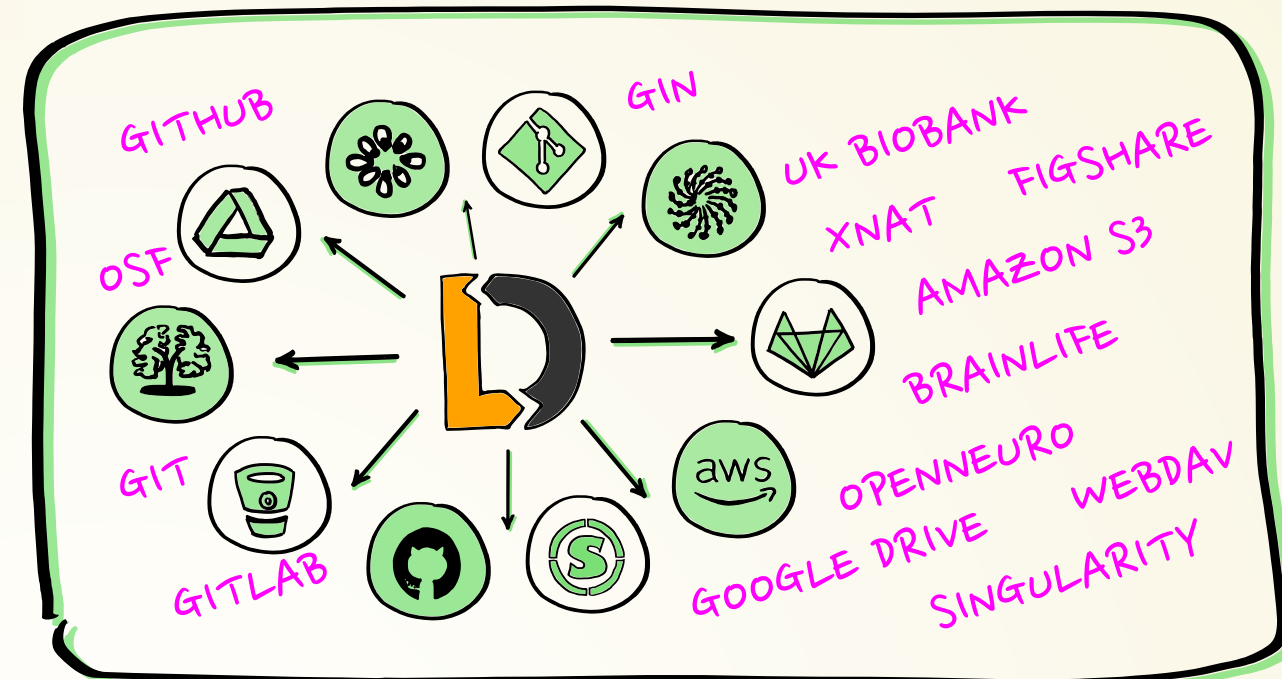
```
git log -n 1
```

[copy](#)

# PUBLISHING DATASETS...



## ● INTEGRATE AND EXTEND



We will use GIN: [gin.g-node.org](http://gin.g-node.org)

(but dozens of other places are supported: See chapter on third-party publishing

:

# PUBLISHING DATASETS...

- Create a GIN user account and log in: [gin.g-node.org/user/sign\\_up](https://gin.g-node.org/user/sign_up)
- Create an SSH key

```
ssh-keygen -t ed25519 -C "your-email"  
eval "$(ssh-agent -s)"  
ssh-add ~/.ssh/id_ed25519
```

- upload the SSH key to GIN

```
cat ~/.ssh/id_ed25519.pub
```

The screenshot shows the GitHub 'Manage SSH Keys' interface. The top navigation bar includes links for Dashboard, Issues, Pull Requests, Explore, Help, and News. The left sidebar shows the 'Settings' menu with options for Profile, Avatar, Password, Email Addresses, SSH Keys (which is highlighted), and Security. The main content area is titled 'Manage SSH Keys' and includes an 'Add Key' button. A warning message states: 'This is a list of SSH keys associated with your account. As these keys allow anyone using them to gain access to your repositories, it is highly important that you make sure you recognize them.' Below this, a single SSH key is listed with a key icon, the name 'private laptop', its SHA256 fingerprint 'SHA256:alG7pJosAS37t06hIU8grM72+xfk0apCJtT1V3ZjSkk', and a 'Delete' button. The key was added on Jan 06, 2020, and last used on Nov 17, 2020. At the bottom, a note says: 'Don't know how? Check out GitHub's guide to [create your own SSH keys](#) or solve [common problems](#) you might encounter using SSH.'

- Publish your dataset!

## ...PUBLISHING DATASETS

DataLad has convenience functions to create `sibling`-repositories on various infrastructure and third party services (GitHub, GitLab, OSF, WebDAV-based services, DataVerse, ...) , to which data can then be published with `push`.

```
datalad create-sibling-gin example-analysis --access-protocol ssh
```

[copy](#)

You can verify the dataset's siblings with the `siblings` command:

```
datalad siblings
```

[copy](#)

And we can push our complete dataset (Git repository and annex) to GIN:

```
datalad push --to gin
```

[copy](#)

**In case of fire**



1. git commit

(datalad save)



2. git push

(datalad push)



3. leave building



## USING PUBLISHED DATA...

Let's see how the analysis feels like to others:

```
cd ../  
datalad clone \  
  https://gin.g-node.org/adswa/example-analysis \  
  myclone
```

[copy](#)

```
cd myclone
```

[copy](#)

Get results:

```
datalad get figures/inter_trial_coherence.png
```

[copy](#)

```
datalad drop figures/inter_trial_coherence.png
```

[copy](#)

Or recompute results:

```
datalad rerun
```

[copy](#)



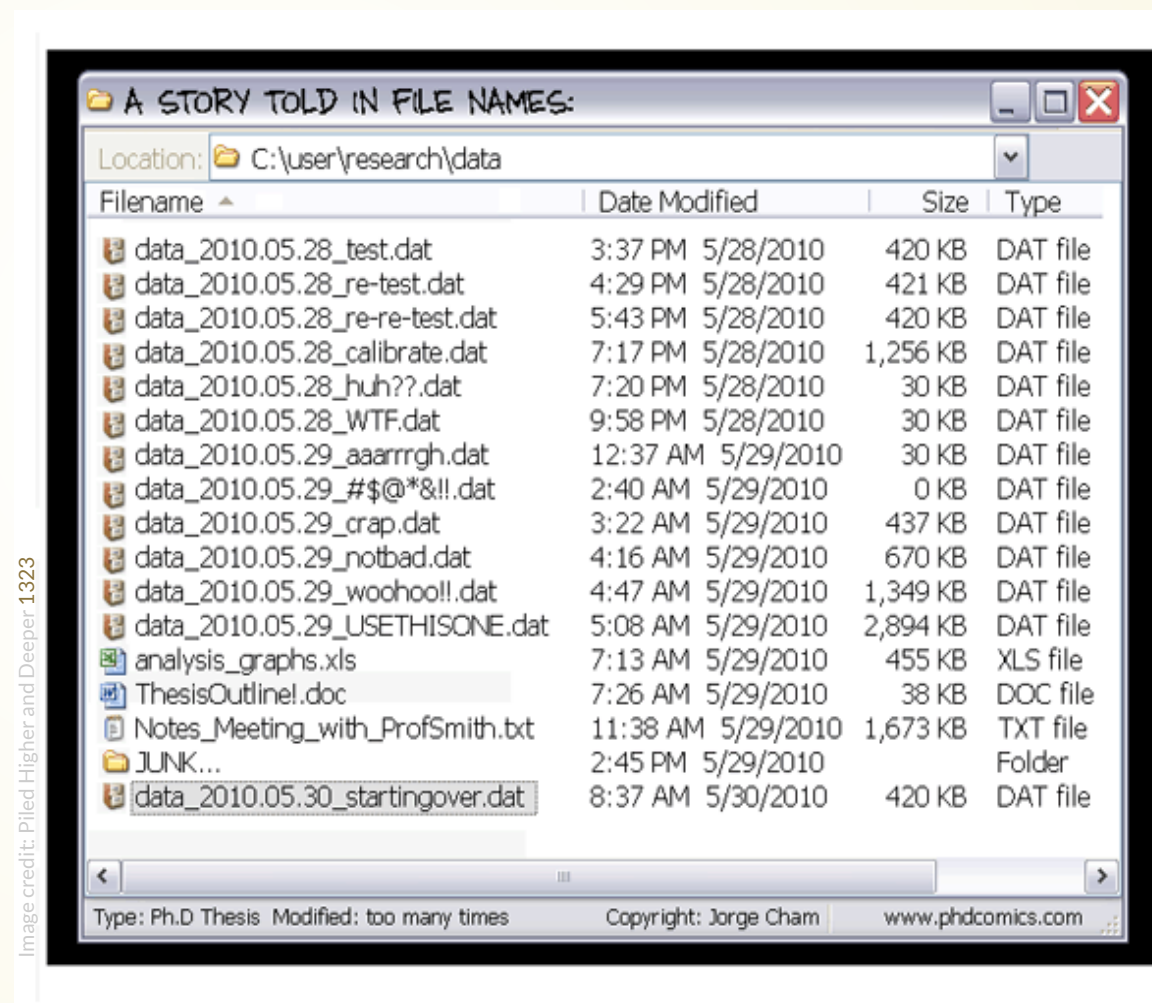
**HOW DOES THIS RELATE TO REPRODUCIBILITY?**

# EXHAUSTIVE TRACKING

The building blocks of a scientific result are rarely static

## Data changes

(errors are fixed, data is extended, naming standards change, an analysis requires only a subset of your data...)



# EXHAUSTIVE TRACKING

"Shit, which version of which script produced these outputs from which version of what data... and which software version?"

Image credit: CC-BY Scriberia and The Turing Way



# EXHAUSTIVE TRACKING

Once you track changes to data with version control tools, you can find out *why* it changed, *what* has changed, *when* it changed, and *which version* of your data was used at which point in time.

```
2020-03-13 10:46 +0100 Adina Wagner    o [DATALAD RUNCMD] add non-defaced
2020-03-13 10:29 +0100 Adina Wagner    o [DATALAD RUNCMD] reconvert DICOM
2018-05-11 09:23 +0200 Michael Hanke    o [master] {origin/HEAD} {origin/m
2018-05-11 09:19 +0200 Michael Hanke    o Enable DataLad metadata extracto
2018-05-11 09:17 +0200 Michael Hanke    o [DATALAD] new dataset
2018-05-11 09:17 +0200 Michael Hanke    o [DATALAD] Set default backend fo
2018-01-19 14:19 +0100 Michael Hanke    o <v1.5> Update changelog for 1.5
2018-01-19 14:09 +0100 Michael Hanke    o BF: Re-import respiratory trace
2018-01-14 18:59 +0100 Michael Hanke    o Fix type in physio log converter
2017-01-10 10:10 +0100 Michael Hanke    o ENH: Report per-stimulus events
2016-12-10 20:18 +0100 Michael Hanke    o Add BIDS-compatible stimuli/ dir
2016-11-15 07:04 +0100 Michael Hanke    o Minor tweaks to gaze overlay scr
2016-10-30 11:03 +0100 Michael Hanke    o Add "TaskName" meta data field f
2016-09-21 08:33 +0200 Michael Hanke    o Add task-*_physio.json files
2016-09-21 08:23 +0200 Michael Hanke    o BF: Fix task label in file names
2016-08-04 13:14 +0200 Michael Hanke    o Update changelog
2016-08-03 22:22 +0200 Michael Hanke    o Add cut position information to
2016-05-27 17:35 +0200 Michael Hanke    o {origin/_} Mention openfmri as d
2016-04-04 09:31 +0200 Michael Hanke    o Update publication links
2016-03-31 11:26 +0200 Michael Hanke    o Disable invalid test
[main] 6da25fb6fee2c698d35f52066698b6f94850f4d2 - commit 10 of 79    27%

commit 6da25fb6fee2c698d35f52066698b6f94850f4d2
Refs: v1.0-19-g6da25fb6
Author: Michael Hanke <michael.hanke@gmail.com>
AuthorDate: Fri Jan 19 14:09:53 2018 +0100
Commit: Michael Hanke <michael.hanke@gmail.com>
CommitDate: Fri Jan 19 14:11:23 2018 +0100

    BF: Re-import respiratory trace after bug fix in converter (fixes gh-
---
...er_task-movielocalizer_run-1_recording-cardresp_physio.tsv.gz | 2 +-
..._task-objectcategories_run-1_recording-cardresp_physio.tsv.gz | 2 +-
..._task-objectcategories_run-2_recording-cardresp_physio.tsv.gz | 2 +-
..._task-objectcategories_run-3_recording-cardresp_physio.tsv.gz | 2 +-
..._task-objectcategories_run-4_recording-cardresp_physio.tsv.gz | 2 +-
...calizer_task-retmapccw_run-1_recording-cardresp_physio.tsv.gz | 2 +-
...calizer_task-retmapclw_run-1_recording-cardresp_physio.tsv.gz | 2 +-
...calizer_task-retmapcon_run-1_recording-cardresp_physio.tsv.gz | 2 +-
...calizer_task-retmapexp_run-1_recording-cardresp_physio.tsv.gz | 2 +-
...2_ses-movie_task-movie_run-1_recording-cardresp_physio.tsv.gz | 2 +-
...2_ses-movie_task-movie_run-2_recording-cardresp_physio.tsv.gz | 2 +-
[diff] 6da25fb6fee2c698d35f52066698b6f94850f4d2 - line 1 of 2391    0%
```



# DIGITAL PROVENANCE

= "The tools and processes used to create a digital file, the responsible entity, and when and where the process events occurred"

- Have you ever saved a PDF to read later onto your computer, but forgot where you got it from? Or did you ever find a figure in your project, but forgot which analysis step produced it?

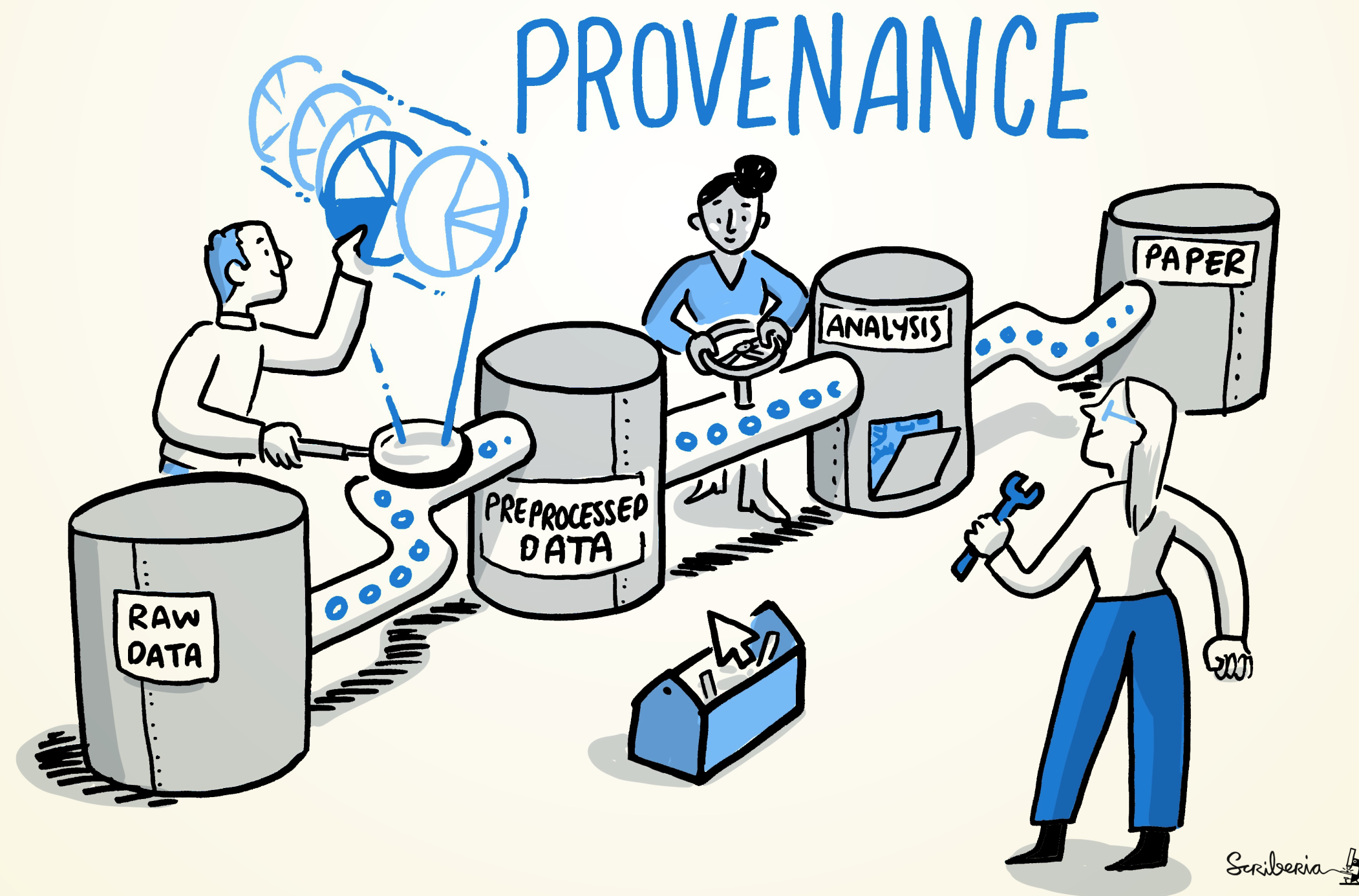
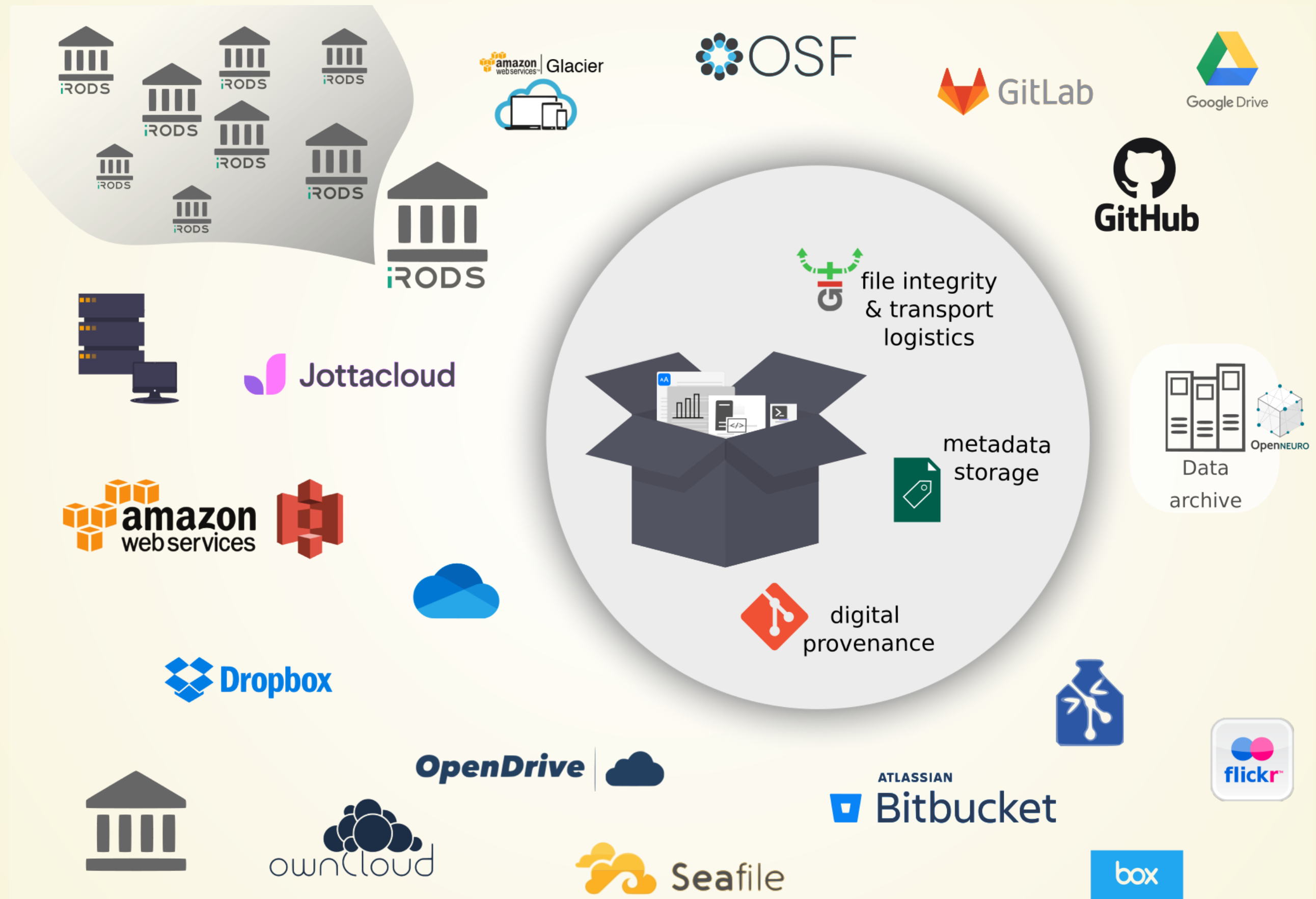


Image credit: Scriberia and The Turing Way (CC-BY)

# DATA TRANSPORT: SECURITY AND RELIABILITY - FOR DATA

Decentral version control for data integrates with a variety of services to let you store data in different places - creating a resilient network for data



"In defense of decentralized Research Data Management", [doi.org/10.1515/nf-2020-0037](https://doi.org/10.1515/nf-2020-0037)



# ULTIMATE GOAL: REUSABILITY

Teamscience on more than code:

Q: Adjustment to tiny labeling mistake (ref Zemblys, 2018)? #2

Closed

adswa opened this issue on Mar 8, 2019 · 1 comment



adswa commented on Mar 8, 2019

Contributor

😊 ...

Zemblys et al., 2018, report a minor labeling mistake in one image file [here](#) (or a pay-wall free version [here](#)) in Appendix 2:

We found an obvious labeling mistake in the one of the validation trials (file UH29\_img\_Europe\_labelled\_MN. We fixed this error by reassigning 75 samples, [3197,3272) (zero-based index), from the saccade to the fixation class.

I checked the data file in question and it appears to still contain the erroneous saccade labels. Just to reconfirm: this labeling error has not been fixed in the data file in this repository, correct?

If you wish, I can PR a fixed file, the issue at hand is intended to just reconfirm my assumption.

Thanks in advance!



adswa added a commit to adswa/remodnav that referenced this issue on Mar 8, 2019

 ENH/FIX: use file with fixed labels. ... 1b2b162

```
2019-03-08 12:38 +0100 Richard Andersson M [master] {origin/master} {origin/HEAD} Merge pull request #3 from AdinaWagner/datafix
2019-03-08 11:35 +0100 Adina Wagner | o ENH/FIX: relabel erroneous saccades to fixations, closes #2.
2018-12-05 15:27 +0100 Richard Andersson o- Uploaded a folder consting only of the data used in the original article
2017-08-22 19:40 +0200 Richard Andersson o Code added
2016-12-14 10:35 +0100 richardandersson o Added currently shared data and stimuli.
2016-12-14 10:29 +0100 richardandersson o Deleted -- to be reuploaded
2016-12-14 10:08 +0100 Richard Andersson o Add files via upload
2016-12-14 10:07 +0100 Richard Andersson o Add files via upload
2016-12-14 10:04 +0100 Richard Andersson I Initial commit
```

Fixed data

Original data

# THE YODA PRINCIPLES

# DATALAD DATASETS FOR DATA ANALYSIS

- A DataLad dataset can have any structure, and use as many or few features of a dataset as required.
- However, for **data analyses** it is beneficial to make use of DataLad features and structure datasets according to the **YODA principles**:



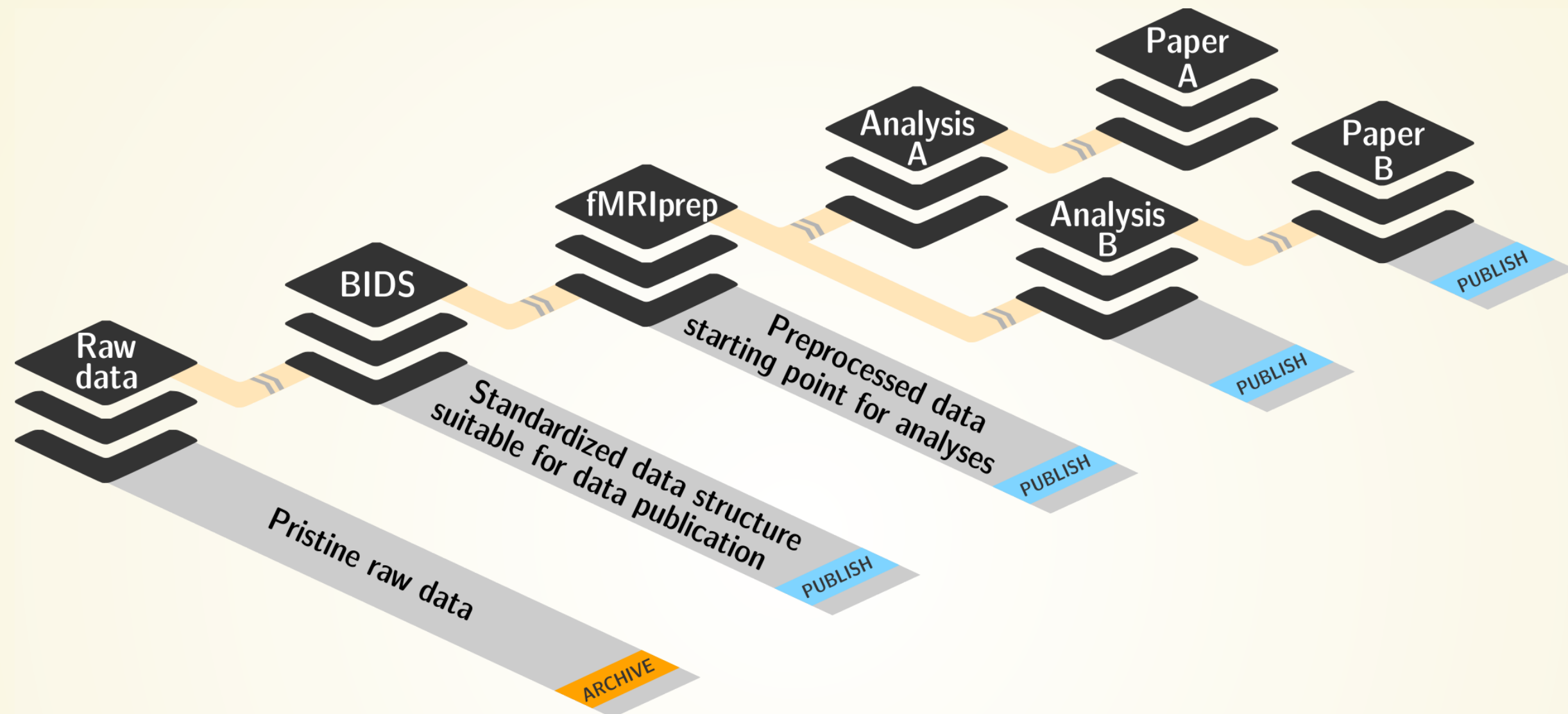
**P1: One thing, one dataset**

**P2: Record where you got it from, and where it is now**

**P3: Record what you did to it, and with what**

Find out more about the YODA principles in [the handbook](#), and more about structuring dataset at [psychoinformatics-de.github.io/rdm-course/02-structuring-data](https://psychoinformatics-de.github.io/rdm-course/02-structuring-data)

# P1: ONE THING, ONE DATASET



- Create **modular** datasets: Whenever a particular collection of files could anyhow be useful in more than one context (e.g. data), put them in their own dataset, and install it as a subdataset.
- Keep everything **structured**: Bundle all components of one analysis into one superdataset, and within this dataset, separate code, data, output, execution environments.
- Keep a dataset **self-contained**, with relative paths in scripts to subdatasets or files. Do not use absolute paths.

# WHY MODULARITY?

- 1. Reuse and access management
- 2. Scalability
- 3. Transparency

Original:

```
/dataset
├── sample1
│   └── a001.dat
├── sample2
│   └── a001.dat
└── ...
```

Without modularity, after applied transform (preprocessing, analysis, ...):

```
/dataset
├── sample1
│   ├── ps34t.dat
│   └── a001.dat
├── sample2
│   ├── ps34t.dat
│   └── a001.dat
└── ...
```

Without expert/domain knowledge, no distinction between original and derived data possible.

# WHY MODULARITY?

- 3. Transparency

Original:

```
/raw_dataset
├── sample1
│   └── a001.dat
├── sample2
│   └── a001.dat
└── ...
```

**With modularity** after applied transform (preprocessing, analysis, ...)

```
/derived_dataset
├── sample1
│   └── ps34t.dat
├── sample2
│   └── ps34t.dat
├── ...
└── inputs
    └── raw
        ├── sample1
        │   └── a001.dat
        ├── sample2
        │   └── a001.dat
        └── ...
```

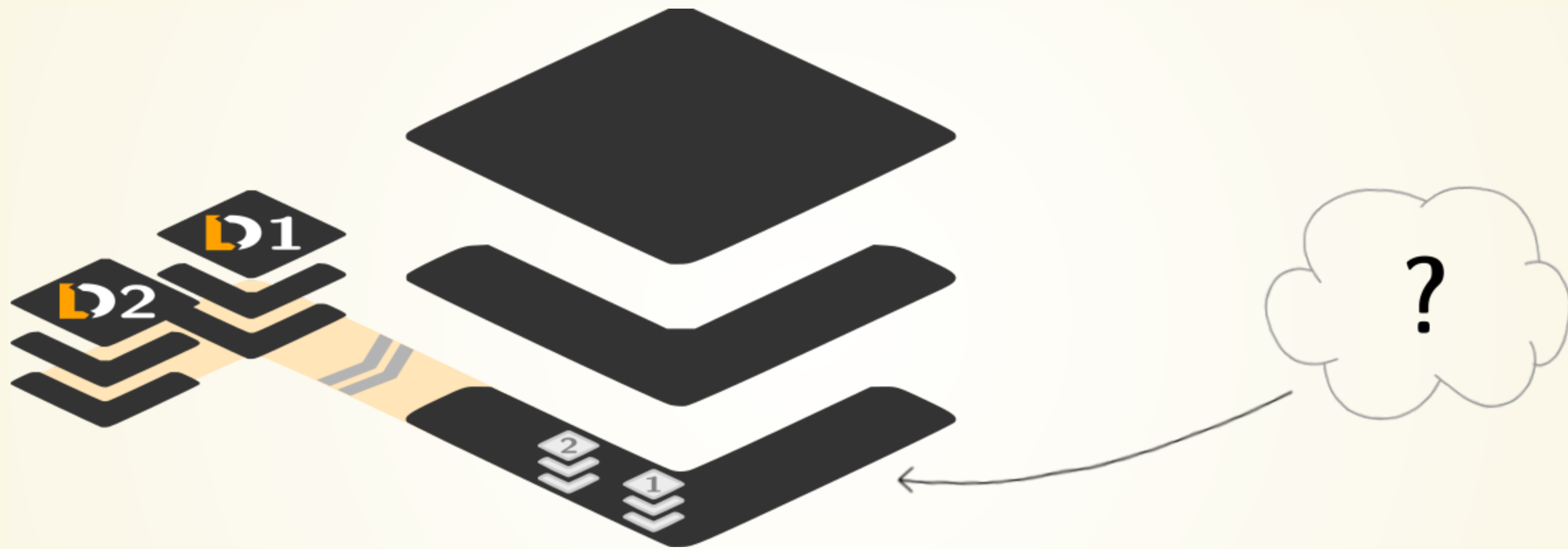
Clearer separation of semantics, through use of pristine version of original dataset within a *new, additional* dataset holding the outputs.



# WHEN TO MODULARIZE?

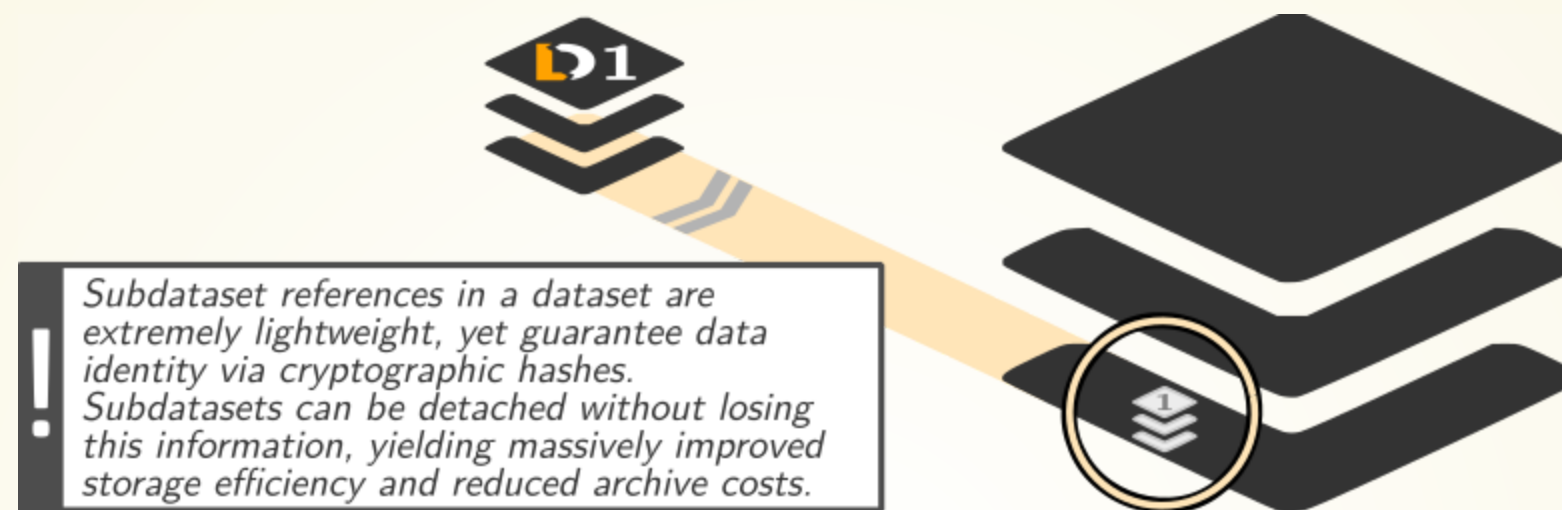
- Target audience is different
  - public vs. private
  - domain specific vs. domain general
- Pace of evolution is different
  - "factual" raw data vs. choices of (pre-)processing
  - completed acquisition vs. ongoing study
- Size impacts I/O and logistics
  - Git can struggle with 1M+ files
  - filesystems (licensing) can struggle with large numbers of inodes
  - More infos: [Go Big or Go Home chapter](#)
- Legal/Access constraints
  - personal vs. anonymized data

## P2: RECORD WHERE YOU GOT IT FROM, AND WHERE IT IS NOW



- **Link** individual datasets to declare data-dependencies (e.g. as subdatasets).
- **Record data's origin** with appropriate commands, for example to record access URLs for individual files obtained from (unstructured) sources "in the cloud".
- Share and **publish** datasets for collaboration or back-up.

# DATASET LINKAGE



```
$ datalad clone --dataset . http://example.com/ds inputs/rawdata
```

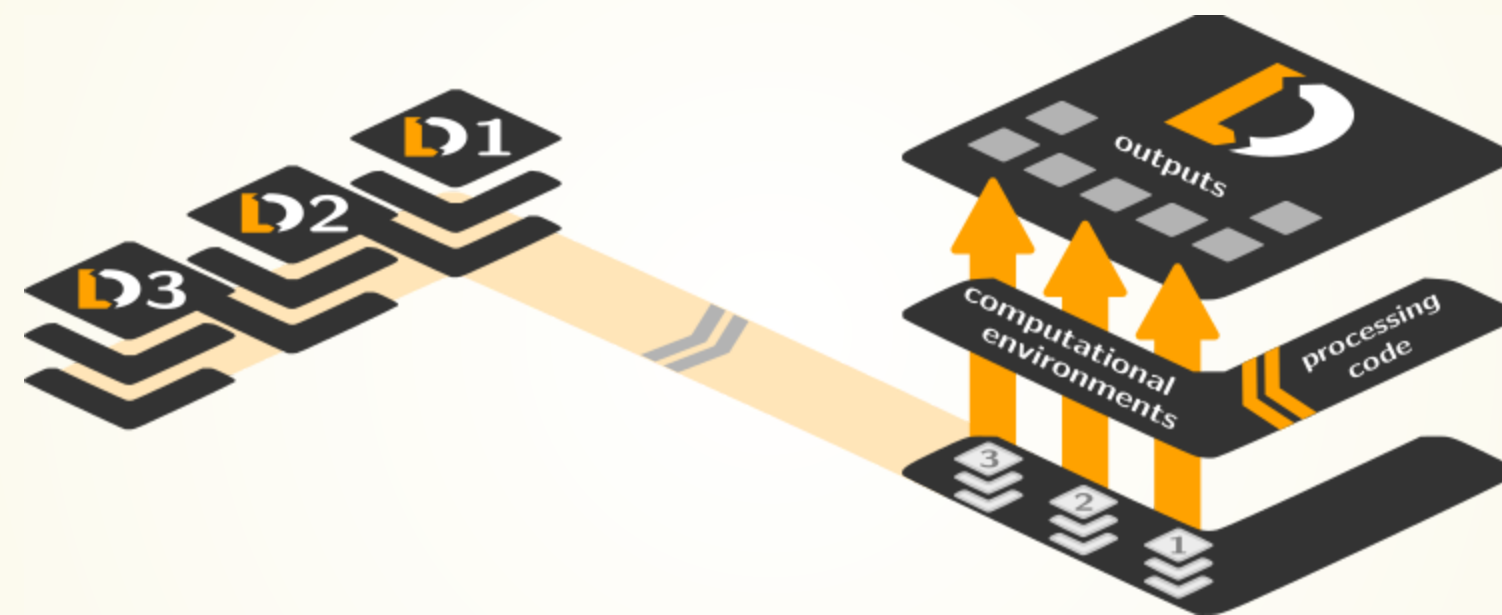
copy

```
$ git diff HEAD~1
diff --git a/.gitmodules b/.gitmodules
new file mode 100644
index 0000000..c3370ba
--- /dev/null
+++ b/.gitmodules
@@ -0,0 +1,3 @@
+[submodule "inputs/rawdata"]
+    path = inputs/rawdata
+    url = http://example.com/importantds
diff --git a/inputs/rawdata b/inputs/rawdata
new file mode 160000
index 0000000..fabf852
--- /dev/null
+++ b/inputs/rawdata
@@ -0,0 +1 @@
+Subproject commit fabf8521130a13986bd6493cb33a70e580ce8572
```

copy

Each (sub)dataset is a separately, but jointly version-controlled entity. If none of its data is retrieved, subdatasets are an extremely **lightweight** data dependency and yet **actionable** (`datalad get` retrieves contents on demand)

## P3: RECORD WHAT YOU DID TO IT, AND WITH WHAT



- Collect and store **provenance** of all contents of a dataset that you create
- "Which script produced which output?", "From which data?", "In which **software environment**?" ... Record it in an ideally machine-readable way with **datalad (containers-)run**

## TAKE HOME MESSAGES

### **Data deserves version control**

It changes and evolves just like code, and exhaustive tracking lays a foundation for reproducibility

### **Reproducible science relies on good data management**

But effort pays off: Increased transparency, better reproducibility, easier accessibility, efficiency through automation and collaboration, streamlined procedures for synchronizing and updating your work, ...

### **DataLad can help with some things**

Have access to more data than you have disk space

Who needs short-term memory when you can have automatic provenance capture?

Link versioned data to your analysis at no disk-space cost

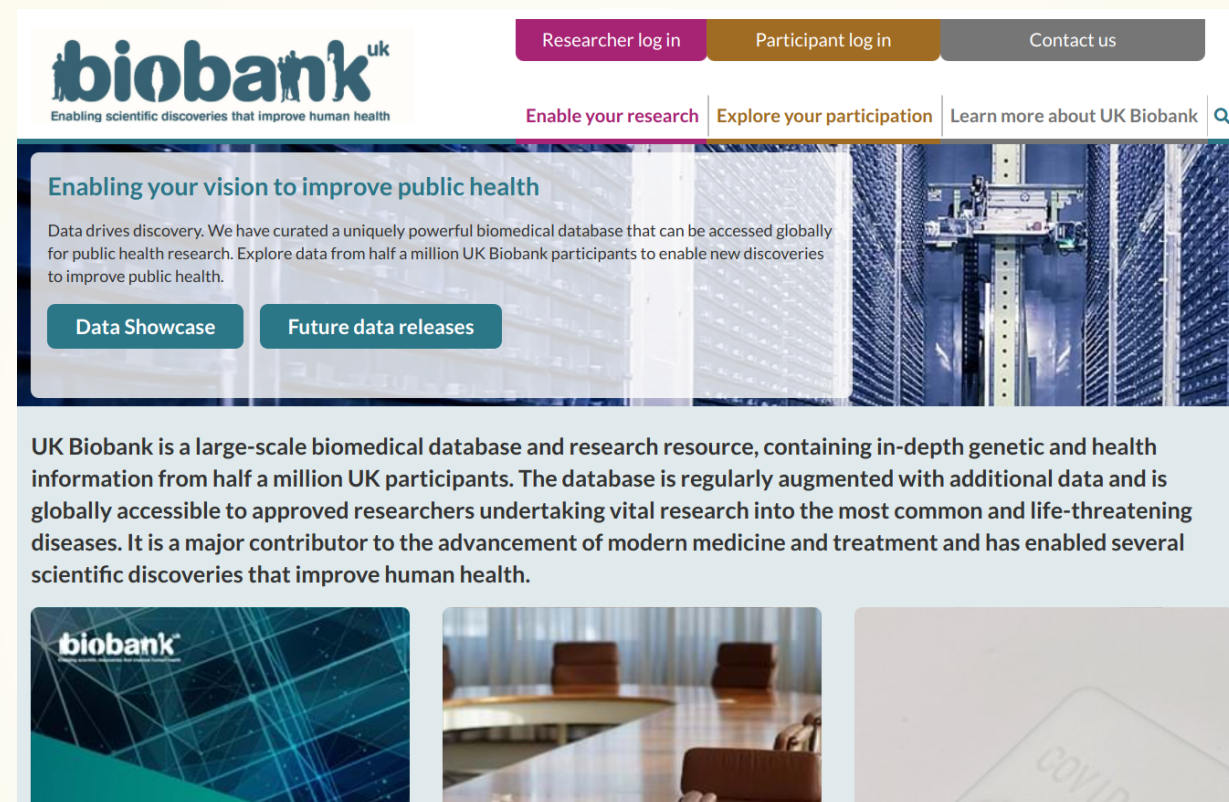
...

# SCALABILITY



# FAIRLY BIG: SCALING UP

Objective: Process the UK Biobank (imaging data)

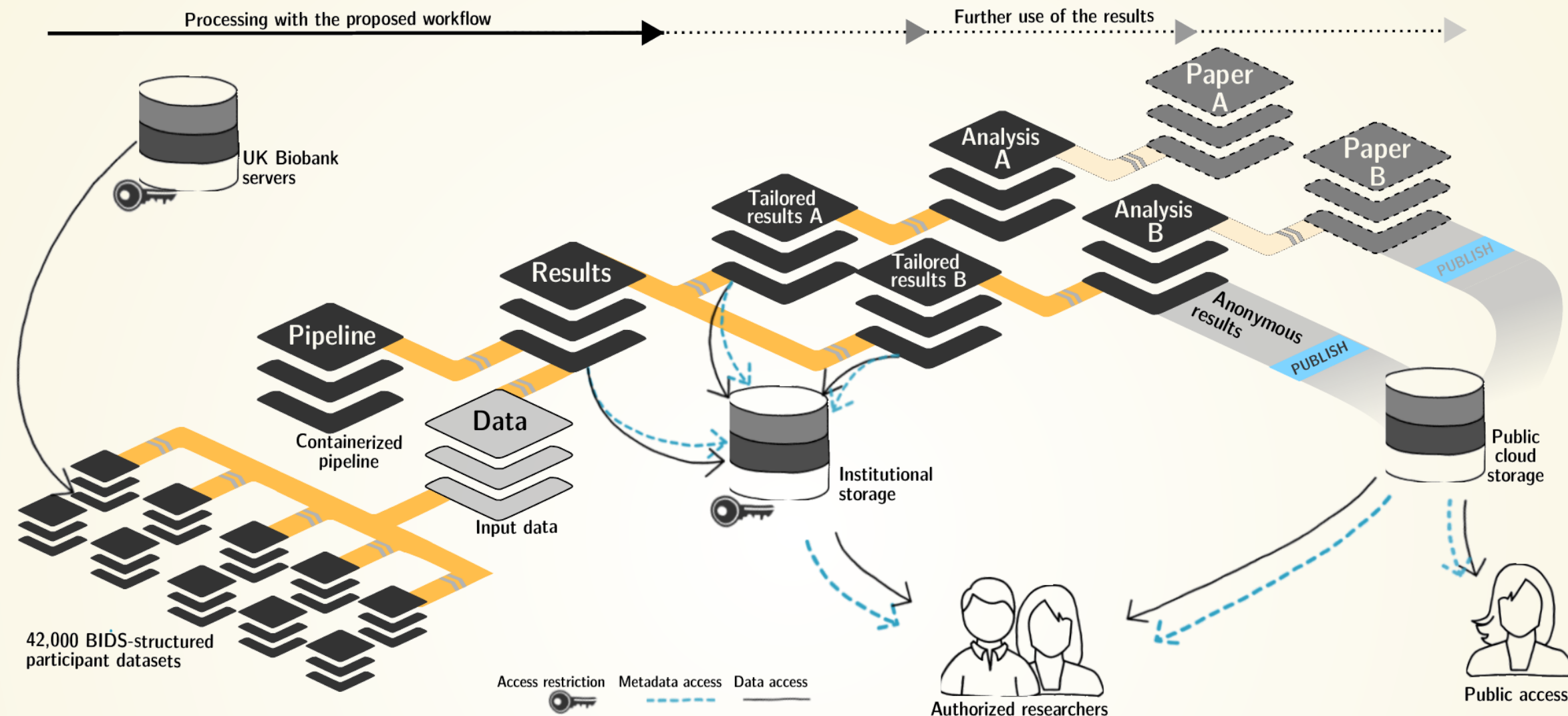


- 76 TB in 43 million files in total
- 42,715 participants contributed personal health data
- Strict DUA
- Custom binary-only downloader
- Most data records offered as (unversioned) ZIP files

# CHALLENGES

- Process data such that
  - Results are computationally reproducible (without the original compute infrastructure)
  - There is complete linkage from results to an individual data record download
  - It scales with the amount of available compute resources
- Data processing pipeline
  - Compiled MATLAB blob
  - 1h processing time per image, with 41k images to process
  - 1.2 M output files (30 output files per input file)
  - 1.2 TB total size of outputs

# FAIRLY BIG SETUP



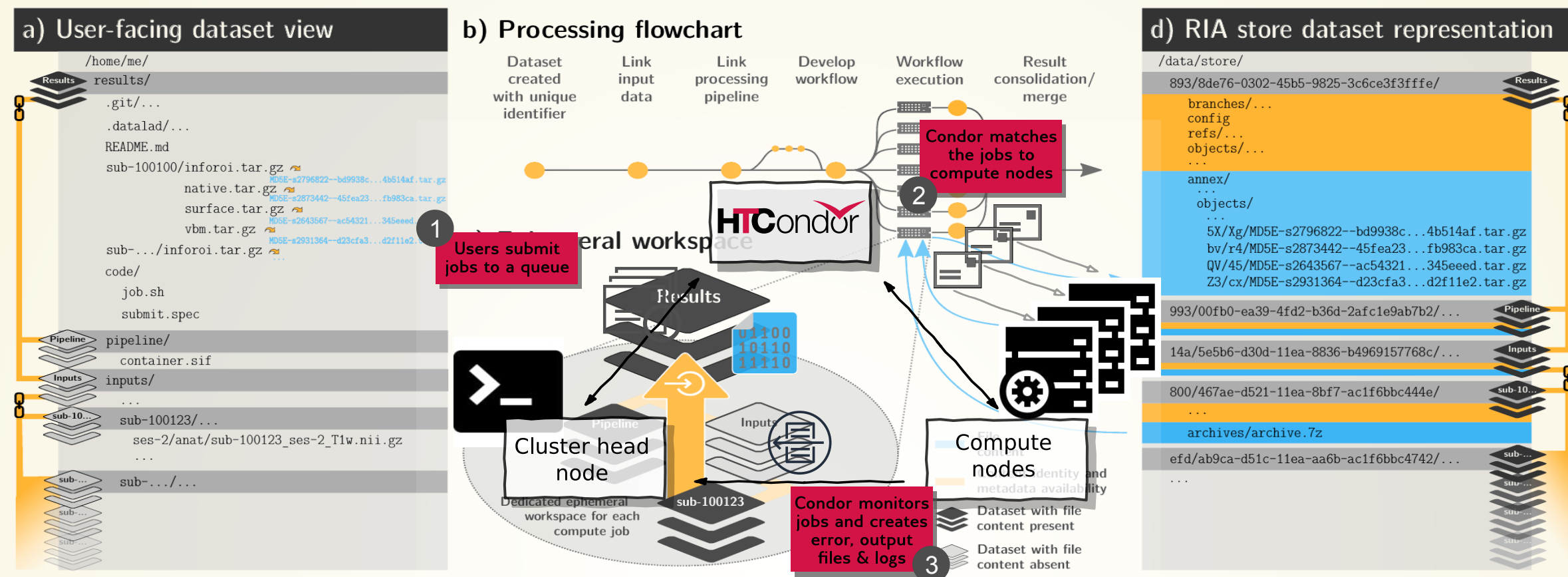
## Exhaustive tracking

- **data-lad-ukbiobank** extension downloads, transforms & track the evolution of the complete data release in DataLad datasets
- Native and BIDSified data layout (at no additional disk space usage)
- Structured in 42k individual datasets, combined to one superdataset
- Containerized pipeline in a software container
- Link input data & computational pipeline as dependencies

Wagner, Waite, Wierzbica et al. (2021). FAIRly big: A framework for computationally reproducible processing of large-scale data.



## FAIRLY BIG WORKFLOW

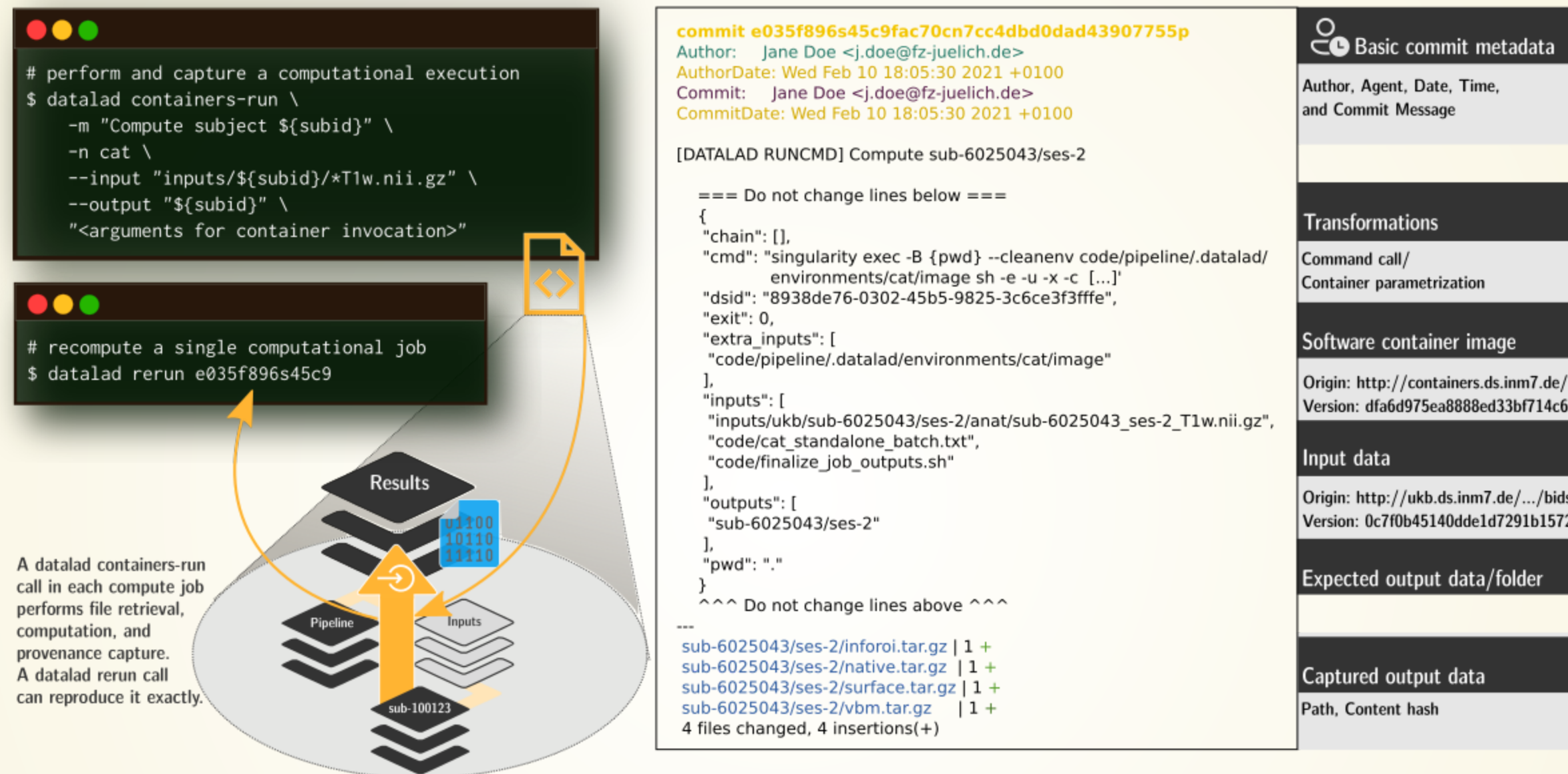


## portability

- Parallel processing: 1 job = 1 subject (number of concurrent jobs capped at the capacity of the compute cluster)
- Each job is computed in a ephemeral (short-lived) dataset clone, results are pushed back: Ensure exhaustive tracking & portability during computation
- Content-agnostic persistent (encrypted) storage (minimizing storage and inodes)
- Common data representation in secure environments

Wagner, Waite, Wierzba et al. (2021). FAIRly big: A framework for computationally reproducible processing of large-scale data.

# FAIRLY BIG PROVENANCE CAPTURE



## Provenance

- Every single pipeline execution is tracked
- Execution in ephemeral workspaces ensures results individually reproducible without HPC access

Wagner, Waite, Wierzba et al. (2021). FAIRly big: A framework for computationally reproducible processing of large-scale data.

# FAIRLY BIG MOVIE

- Two computations on clusters of different scale (small cluster, supercomputer).  
Full video: <https://youtube.com/datalad>
- Two full (re-)computations, programmatically comparable, verifiable, reproducible -- on any system with data access



# THANK YOU FOR YOUR ATTENTION!



Slides: [DOI 10.5281/zenodo.7627723](https://doi.org/10.5281/zenodo.7627723) (Scan the QR code)



Women neuroscientists are underrepresented in neuroscience. You can use the Repository for Women in Neuroscience to find and recommend neuroscientists for conferences, symposia or collaborations, and help making neuroscience more open & divers.

## COMMAND SUMMARIES

## SUMMARY - LOCAL VERSION CONTROL

**datalad create** creates an empty dataset.

Configurations (**-c yoda**, **-c text2git**) add useful structure and/or configurations.

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.

## SUMMARY - DATASET CONSUMPTION & NESTING

**`data1ad clone`** installs a dataset.

It can be installed “on its own”: Specify the source (url, path, ...) 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 annexed files, `data1ad get` downloads file content on demand.

**Datasets preserve their history.**

The superdataset records only the version state of the subdataset.

## SUMMARY - REPRODUCIBLE EXECUTION

**data1ad run** records a command and its impact on the dataset.

All dataset modifications are saved - use it in a clean dataset.

Data/directories specified as **--input** are retrieved prior to command execution.

Use one flag per input.

Data/directories specified as **--output** will be unlocked for modifications prior to a rerun of the command.

Its optional to specify, but helpful for recomputations.

**data1ad containers-run** can be used to capture the software environment as provenance.

Its ensures computations are ran in the desired software set up. Supports Docker and Singularity containers

**data1ad rerun** can automatically re-execute run-records later.

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



