## **DATALAD**

### DECENTRALIZED MANAGEMENT OF DIGITAL OBJECTS FOR OPEN SCIENCE

Dr. Adina Wagner



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



Slides: DOI 10.5281/zenodo.15193934 (Scan the QR code) files.inm7.de/adina/talks/html/nhr

## **ACKNOWLEDGEMENTS**

# DataLad software & ecosystem

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

#### **Funders**







SPONSORED BY THE

BMBF 01GQ1411











#### **Collaborators**





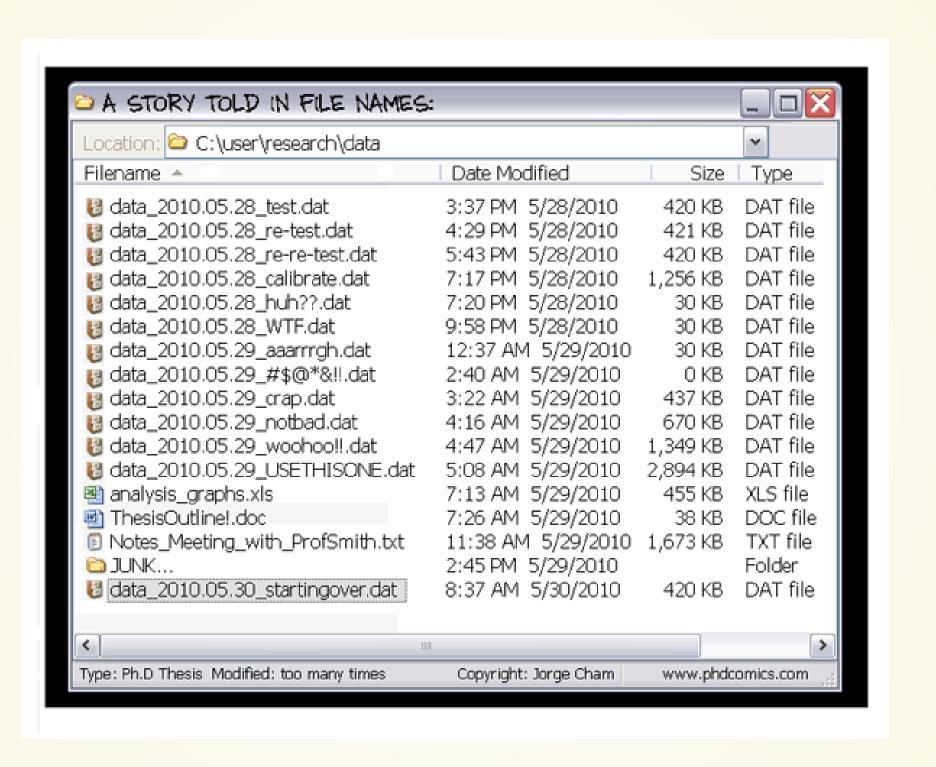


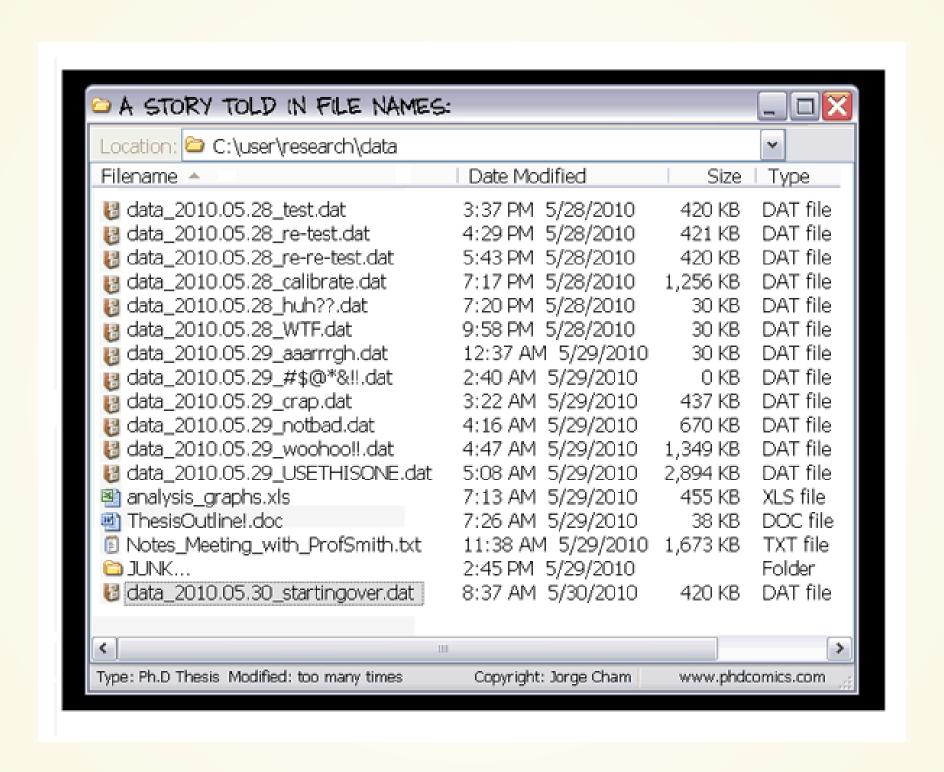


eBRAIN Health









#### <u>Data</u> changes

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

Mar 2019

#### ABCD Data Release 2.0 available now on the NIMH Data Archive!

The second annual curated ABCD data release 2.0 is available now on the NIMH Data Archive. Data Release 2.0 includes baseline data on the full participant cohort, ages 9-10 years. This release contains much of the same type of data as in Data Release 1.1, and also includes genotypic data for the first time. Smokescreen genotyping array data are available on almost 11,000 participants. These include common variations, as well as variations associated with addiction, smoking behavior and nicotine metabolism.

ABCD Study data will be released annually. The next data release will be in early summer 2020 and will include the first longitudinal data from the 6-month and 1-year follow-up assessments.

#### Spring 2019

#### **Issues Identified with Data Release 2.0**

Problems were identified with the data in 11 imagingrelated files (5 diffusion tensor imaging (DTI) files and 6 fMRI Monetary Incentive Delay (MID) task files). See below for more detail. Corrected files will be re-uploaded to the NIMH Data Archive and available soon to users.

DTI subcortical and cortical ROIs and white fiber tracts (Part 1)

DTI subcortical and cortical ROIs and white matter fiber tracts (Part 4)

dMRI DTI Destrieux Parcellations Part 1

dMRI DTI Destrieux Parcellations Part 2

dMRI DTI Full Part 2

Task fMRI MID Average SEM Destriex Parcellations Part 2
Task fMRI MID Run 1 Beta Weights Destrieux Parcellations
Part 2

Task fMRI MID Run 2 Beta Weights Destrieux Parcellations Part 2

Task fMRI MID Average Beta Weights Destrieux Parcellations Part 2

Task fMRI MID Run 1 SEM Destrieux Parcellations Part 2
Task fMRI MID Run 2 SEM Destrieux Parcellations Part 2

An error was also discovered in imaging data collected from Siemens scanners between September 2017 and December 2017 where structural images are flipped left-right. These data will be updated in a patch release later this year.

July 2019

#### ABCD Data Release 2.0.1 available now on the NIMH Data Archive

Due to reporting compilation and processing errors in 2.0 Data Release, a 2.0.1 Fix Release has been issued. Please ensure curated data (datasheets, minimally processed data) from the original Data Release 2.0 are replaced with data from 2.0.1 Fix Release. The following release notes were updated to reflect these changes:

- NDA 2.0.1 Release Notes ABCD README FIRST
- NDA 2.0.1 Release Notes Imaging Instruments
- NDA 2.0.1 Changes and Known Issues Fix Release 2.0.1
- NDA 2.0.1 Diffusion Magnetic Resonance Imaging
- NDA 2.0.1 Task-Based Functional Magnetic Resonance Imaging
- NDA 2.0.1 Mental Health
- NDA 2.0.1 Genetics

#### Dec 2019

#### Issues Identified with Data Release 2.0.1

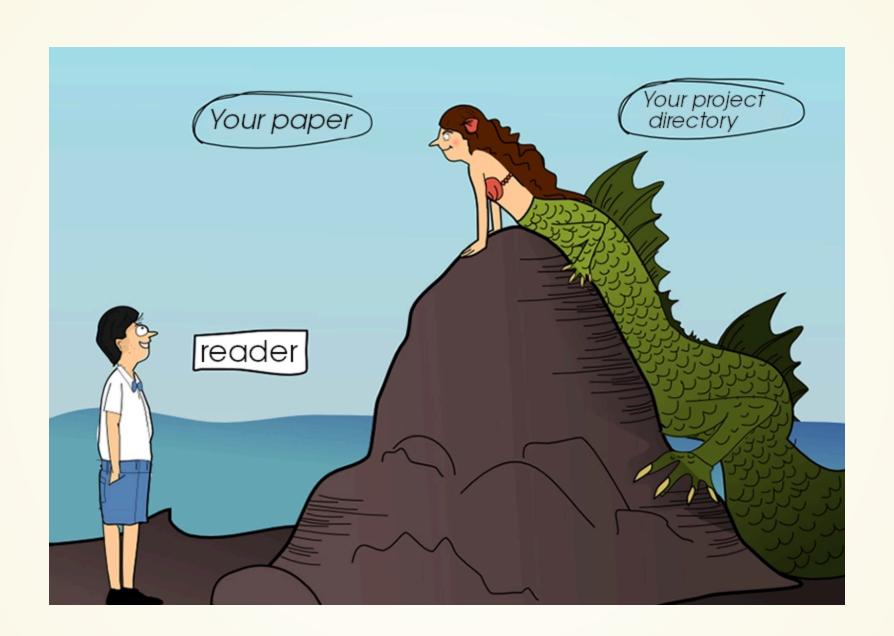
Issues have been identified with ABCD Data Release 2.0.1 that impact some of the neuroimaging and genetics data. See below for details.

#### fMRI:

Due to incorrect post-processing, all task and resting-state fMRI data obtained on Phillips scanners should be excluded from all analyses. This includes tabulated and minimally processed data. This issue <u>does not</u> <u>affect</u> other modalities (sMRI, dMRI) or raw DICOM data (Fast Track). Philips fMRI data have been removed from DEAP and RDS. They will be corrected and made available with Data Release 3.0.

Philips fMRI *minimally processed* and *tabulated data* can be excluded by using the ABCD MRI Info instrument (abcd\_mri01) and exclude by: mri\_info\_manufacturer = "Philips Medical Systems". An R script is available at https://github.com/ABCD-STUDY/fMRI-cleanup to remove Philips fMRI data from tabulated data.

We recommend that all users re-run fMRI related analyses without the Philips fMRI data. For published studies, we encourage users to confirm findings without the Philips fMRI data and provide a corrigendum with the updated results.





- Domain-agnostic command-line tool (+ graphical user interface), built on top of Git & Git-annex
- Open source (MIT) research software developed since 2013
- Available for all major operating systems
- Major features:

Version-controlling arbitrarily large content

Version control data & software alongside to code!

**Transport mechanisms for sharing & obtaining data** 

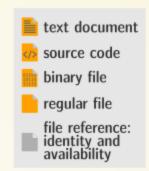
Consume & collaborate on data (analyses) like software

(Computationally) reproducible data analysis

Track and share provenance of all digital objects

(... and much more)



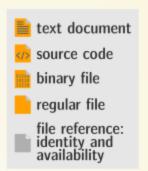


A DataLad dataset is a joint Git/git-annex repository that can version control any file

```
# turn any directory into a dataset
# with version control # file co
% datalad create <directory> % datalad
```

```
# save a new state of a dataset with
# file content of any size
% datalad save
```

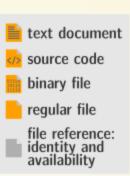




Which data (at which version), with which code, running with what parameterization in which computational environment, to generate what?

```
# execute any command and capture its output
# while recording all input versions too
% datalad run --input ... --output ... <command>
```

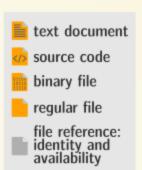




Decentral data transport to Git hosting, local or remote infrastructure, or external hosting services

```
# transfer data and metadata to other sites and services
# with fine-grained access control for dataset components
% datalad push --to <site-or-service>
```





Outcomes can be validated. This enables audits, promotes accountability, and streamlines automated "upgrades" of outputs

```
# obtain dataset (initially only identity,
# availability, and provenance metadata)
% datalad clone <url>
```

```
# immediately actionable provenance records
# full abstraction of input data retrieval
% datalad rerun <commit|tag|range>
```



Datasets can be (re-)used as modular components in larger contexts — propagating their traits. They are verifiable, portable, self-contained data structures

```
# declare a dependency on another dataset and
# re-use it a particular state in a new context
% datalad clone -d <superdataset> <url> <path-in-dataset>
```









```
o [DATALAD RUNCMD] add non-defaced anatomical images
                                    o [DATALAD RUNCMD] reconvert DICOMs without defacing
                                    o [master] {origin/HEAD} {origin/master} {origin/synced/master} [DATALAD] dataset aggregate metadata update
018-05-11 09:23 +0200 Michael Hanke
                                    o Enable DataLad metadata extractors
018-05-11 09:19 +0200 Michael Hanke
                                    o [DATALAD] new dataset
018-05-11 09:17 +0200 Michael Hanke
                                    o [DATALAD] Set default backend for all files to be MD5E
018-05-11 09:17 +0200 Michael Hanke
                                    o <v1.5> Update changelog for 1.5
o BF: Re-import respiratory trace after bug fix in converter (fixes gh-11)
018-01-19 14:09 +0100 Michael Hanke
                                    o Fix type in physio log converter (fixes gh-11)
018-01-14 18:59 +0100 Michael Hanke
                                    o : Report per-stimulus events (fixes gh-6)
917-01-10 10:10 +0100 Michael Hanke
                                    o Ada LDG-compatible stimuli/ directory (with symlinks)
016-12-10 20:18 +0100 Michael Hanke
                                    o Minor tweaks to gaze overlay script
016-11-15 07:04 +0100 Michael Hanke
                                    o Add "TaskName" meta data field for compliance with BIDS 1.0.0
016-10-30 11:03 +0100 Michael Hanke
                                    o Add task-* physio.json files
016-09-21 08:33 +0200 Michael Hanke
                                    o BF: Fix task label in file names of contracting retmap run.
016-09-21 08:23 +0200 Michael Hanke
016-08-04 13:14 +0200 Michael Hanke
                                    • Update changelog
                                    o Add cut position information to allow for timing verification of generated stimulus files
016-08-03 22:22 +0200 Michael Hanke
                                    o {origin/ } Mention openfmri as download source
016-05-27 17:35 +0200 Michael Hanke
016-04-04 09:31 +0200 Michael Hanke
                                     O Update publication links
                                    o Disable invalid test
016-03-31 11:26 +0200 Michael Hanke
```





```
o [DATALAD RUNCMD] add non-defaced anatomical images
                                     o [DATALAD RUNCMD] reconvert DICOMs without defacing
                                     o [master] {origin/HEAD} {origin/master} {origin/synced/master} [DATALAD] dataset aggregate metadata update
018-05-11 09:23 +0200 Michael Hanke
                                     o Enable DataLad metadata extractors
018-05-11 09:19 +0200 Michael Hanke
                                     o [DATALAD] new dataset
018-05-11 09:17 +0200 Michael Hanke
                                     o [DATALAD] Set default backend for all files to be MD5E
018-05-11 09:17 +0200 Michael Hanke
                                     o <v1 > Update changelog for 1.5
                                     o BF: Re-import respiratory trace after bug fix in converter (fixes gh-11)
018-01-19 14:09 +0100 Michael Hanke
                                     o Fix type in physio log converter (fixes gh-11)
018-01-14 18:59 +0100 Michael Hanke
                                     o : Report per-stimulus events (fixes gh-6)
917-01-10 10:10 +0100 Michael Hanke
                                     o Add LTDS-compatible stimuli/ directory (with symlinks)
016-12-10 20:18 +0100 Michael Hanke
                                     o Minor tweaks to gaze overlay script
016-11-15 07:04 +0100 Michael Hanke
                                     o Add "TaskName" meta data field for compliance with BIDS 1.0.0
016-10-30 11:03 +0100 Michael Hanke
                                     o Add task-* physio.json files
016-09-21 08:33 +0200 Michael Hanke
                                     o BF: Fix task label in file names of contracting retmap run.
016-09-21 08:23 +0200 Michael Hanke
016-08-04 13:14 +0200 Michael Hanke
                                     Update changelog
                                     o Add cut position information to allow for timing verification of generated stimulus files
016-08-03 22:22 +0200 Michael Hanke
                                     o {origin/ } Mention openfmri as download source
016-05-27 17:35 +0200 Michael Hanke
016-04-04 09:31 +0200 Michael Hanke
                                     O Update publication links
                                     o Disable invalid test
016-03-31 11:26 +0200 Michael Hanke
```





```
o [DATALAD RUNCMD] add non-defaced commit 6da25fb6fee2c698d35f52066698b6f94850f4d2
   -03-13 10:46 +0100 Adina Wagner
  20-03-13 10:29 +0100 Adina Wagner
                                       o [DATALAD RUNCMD] reconvert DICOM
                                       o [master] {origin/HEAD} {origin/m
                                                                                       Michael Hanke <michael.hanke@gmail.com>
                                                                            \uthor:
                                       o Enable DataLad metadata extracto
                                                                           AuthorDate: Fri Jan 19 14:09:53 2018 +0100
                                       [DATALAD] new dataset
                                       O [DATALAD] Set default backend for
                                                                           CommitDate: Fri Jan 19 14:11:23 2018 +0100
                                       o <v1.5> Update changelog for 1.5
                                       o BF: Re-import respiratory trace
                                                                               BF: Re-import respiratory trace after bug fix in converter (fixes gh
2018-01-19 14:09 +0100 Michael Hanke
                                       o Fix type in physio log converter
                                       o ENH: Report per-stimulus events
                                                                              .er task-movielocalizer run-1 recording-cardresp physio.tsv.gz | 2 +-
                                       o Add BIDS-compatible stimuli/ dir
                                       • Minor tweaks to gaze overlay scr
                                       o Add "TaskName" meta data field f
016-09-21 08:33 +0200 Michael Hanke
                                       o Add task-* physio.json files
                                       o BF: Fix task label in file names
                                       • Update changelog
                                       o Add cut position information to
016-05-27 17:35 +0200 Michael Hanke
                                       o {origin/ } Mention openfmri as d
                                       O Update publication links
 016-04-04 09:31 +0200 Michael Hanke
                                       o Disable invalid test
mainl 6da25fb6fee2c698d35f52066698b6f94850f4d2 - commit 10 of 79
                                                                          |[diff] 6da25fb6fee2c698d35f52066698b6f94850f4d2 - line 1 of 2391
```

- Datasets have an annex to track files without placing their content into Git
- Rather than content, identity (hash) and location information is put into Git:

- Datasets have an annex to track files without placing their content into Git
- Rather than content, identity (hash) and location information is put into Git:
  - Where the filesystem allows it, annexed files are symlinks:

```
$ ls -1 sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz copy lrwxrwxrwx 1 adina adina 142 Jul 22 19:45 sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz -> ../../.git/annex/objects/kZ/K5/MD5E-s24180157--aeb0e5f2e2d5fe4ade97117a8cc5232f.nii.gz/MD5E-s241-aeb0e5f2e2d5fe4ade97117a8cc5232f.nii.gz
```

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

- Datasets have an annex to track files without placing their content into Git
- Rather than content, identity (hash) and location information is put into Git:
  - Where the filesystem allows it, annexed files are symlinks:

```
$ ls -1 sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz copy lrwxrwxrwx 1 adina adina 142 Jul 22 19:45 sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz -> ../../.git/annex/objects/kZ/K5/MD5E-s24180157--aeb0e5f2e2d5fe4ade97117a8cc5232f.nii.gz/MD5E-s241-aeb0e5f2e2d5fe4ade97117a8cc5232f.nii.gz
```

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

The symlink reveals: This internal data organization based on identity hash

```
$ md5sum sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz
aeb0e5f2e2d5fe4ade97117a8cc5232f sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz
```

- Datasets have an annex to track files without placing their content into Git
- Rather than content, identity (hash) and location information is put into Git:
  - Where the filesystem allows it, annexed files are symlinks:

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

The symlink reveals: This internal data organization based on identity hash

```
$ md5sum sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz
aeb0e5f2e2d5fe4ade97117a8cc5232f sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz
```

 The (tiny) symlink instead of the (potentially large) file content is committed version controlling precise file identity without checking contents into Git

```
diff --git a/sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz b/sub-02/func/sub-02_task-oneback_run-01_bold.nii.
new file mode 120000
index 00000000..398e7f1
--- /dev/null
+++ b/sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz
@@ -0,0 +1 @@
+../../.git/annex/objects/kZ/K5/MD5E-s24180157--aeb0e5f2e2d5fe4ade97117a8cc5232f.nii.gz/MD5E-s24180157--aeb0e5f2e2
```

- Datasets have an annex to track files without placing their content into Git
- Rather than content, identity (hash) and location information is put into Git:
  - Where the filesystem allows it, annexed files are symlinks:

```
$ ls -1 sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz copy lrwxrwxrwx 1 adina adina 142 Jul 22 19:45 sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz -> ../../.git/annex/objects/kZ/K5/MD5E-s24180157--aeb0e5f2e2d5fe4ade97117a8cc5232f.nii.gz/MD5E-s241-aeb0e5f2e2d5fe4ade97117a8cc5232f.nii.gz
```

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

The symlink reveals: This internal data organization based on identity hash

 The (tiny) symlink instead of the (potentially large) file content is committed version controlling precise file identity without checking contents into Git

```
diff --git a/sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz b/sub-02/func/sub-02_task-oneback_run-01_bold.nii.
new file mode 120000
index 00000000..398e7f1
--- /dev/null
+++ b/sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz
@@ -0,0 +1 @@
+../../.git/annex/objects/kZ/K5/MD5E-s24180157--aeb0e5f2e2d5fe4ade97117a8cc5232f.nii.gz/MD5E-s24180157--aeb0e5f2e2
```

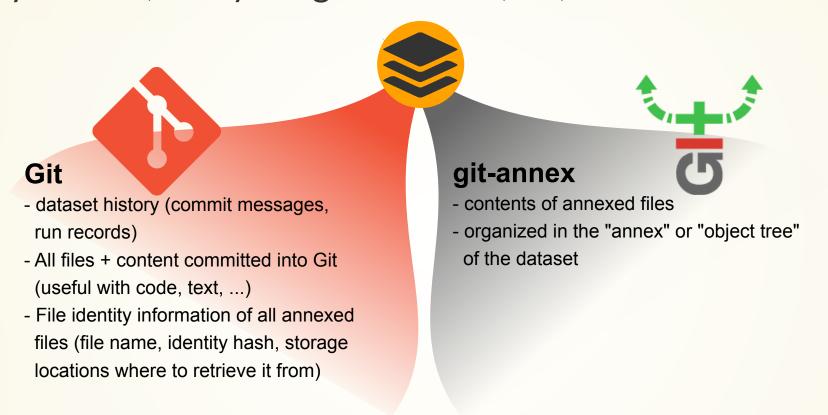
 File availability information is stored to record a decentral network of file content. A file can exist in multiple different locations.

```
$ git annex whereis sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz
whereis sub-02/func/sub-02_task-oneback_run-01_bold.nii.gz (2 copies)
   8c3680dd-6165-4749-adaa-c742232bc317 -- git@8242caf9acd8:/data/repos/adswa/bidsdata.git [gir fff8fdbc-3185-4b78-bd12-718717588442 -- adina@muninn:~/bids-data [here]
ok
```

## 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



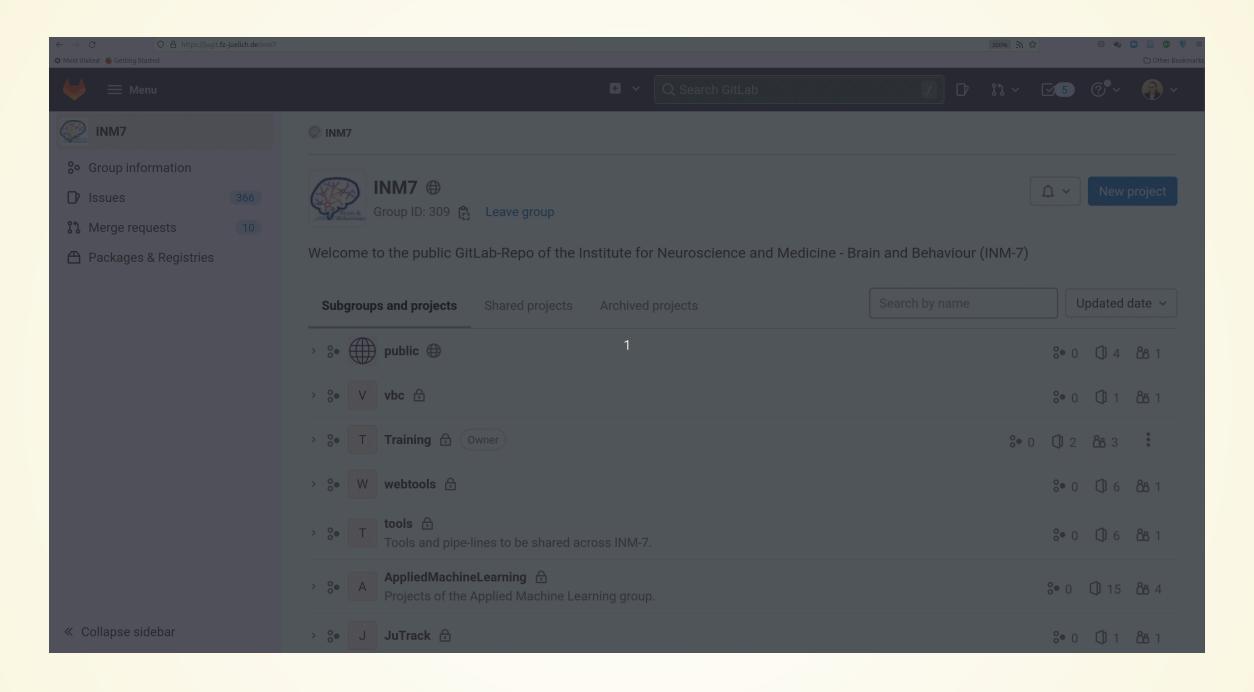
Git	git-annex
handles <b>small</b> files well (text, code)	handles <b>all</b> types and sizes of files well
file contents are in the Git history and will be shared upon git/datalad push	file contents are in the annex. Not necessarily shared
Shared with every dataset clone	Can be kept private on a per-file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

# (RAW) DATA MISMANAGEMENT

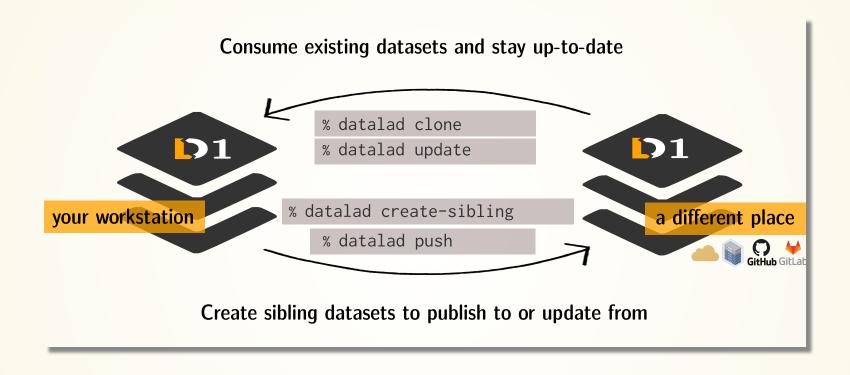
- Multiple large datasets are available on a compute cluster
- Each researcher creates their own copies of data
- Multiple different derivatives and results are computed from it
- Data, copies of data, half-baked data transformations, results, and old versions of results are kept - undocumented

# SHARE DATA LIKE SOURCE CODE

## SHARE DATA LIKE SOURCE CODE



# SHARE DATA LIKE SOURCE CODE



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

 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 .
```

 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:

 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:

 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 copy
get(ok): /tmp/studyforrest-data-phase2/sub-01/ses-movie/func/
sub-01_ses-movie_task-movie_run-1_bold.nii.gz (file) [from mddatasrc...]
```

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

#### TRANSPORT LOGISTICS: LOTS OF DATA, 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 copy
get(ok): /tmp/studyforrest-data-phase2/sub-01/ses-movie/func/
sub-01_ses-movie_task-movie_run-1_bold.nii.gz (file) [from mddatasrc...]
```

Have access to more data on 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 .
```

Drop file content that is not needed:

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): /[...]/sub-01\_ses-movie\_task-movie\_run-1\_bold.nii.gz (file)



Drop file content that is not needed:

Only "meta data" stays behind, and files can be re-obtained on demand. This allows for disk-space-aware computing workflows:

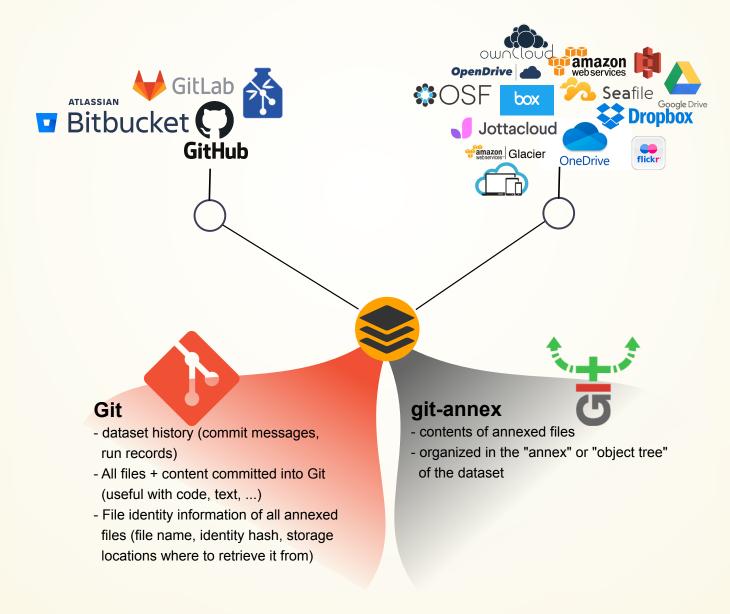
Drop file content that is not needed:

Only "meta data" stays behind, and files can be re-obtained on demand. This allows for disk-space-aware computing workflows:

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

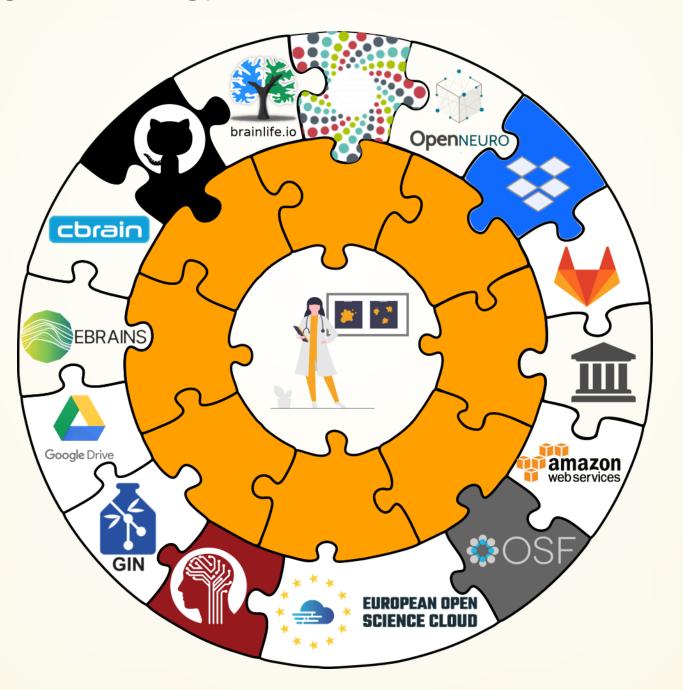
## **PUBLISHING DATASETS**

Publish datasets, their annexed contents, or both to infrastructure of your choice

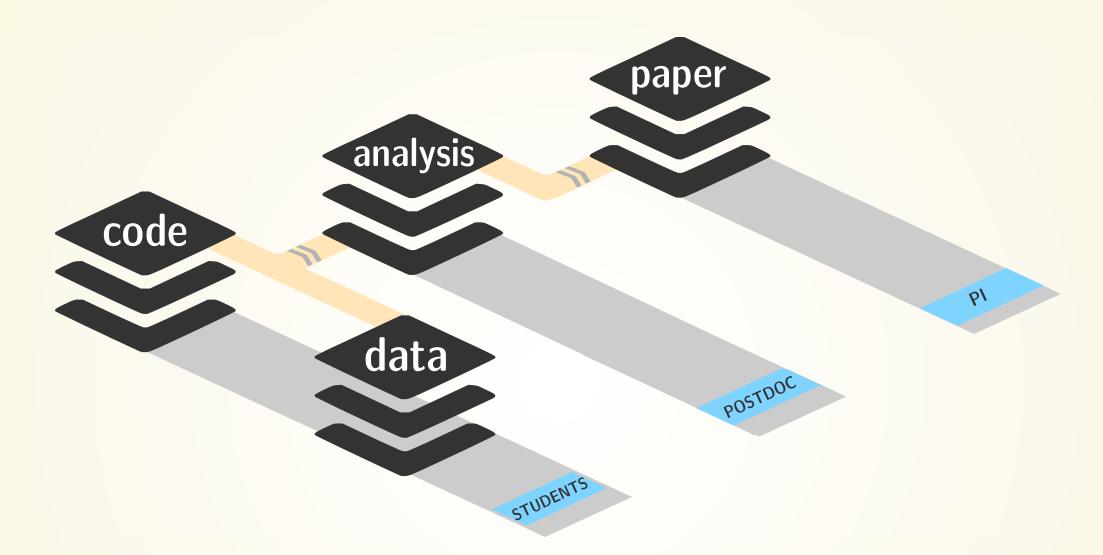


# **INTEROPERABILITY**

 DataLad is built to maximize interoperability and streamline routines across hosting and storage technology



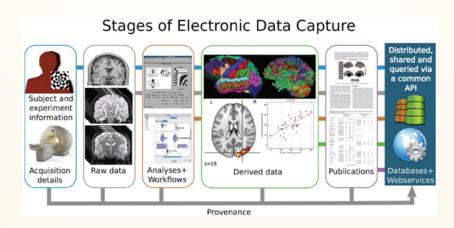
#### **MODULARITY**



- Typical workflow in science
  - Prior works (algorithm development, empirical data, etc.) are combined to produce novel results with to goal of a publication
  - Aggregation across time and contributors
  - Aiming for (but often failing) to be reproducible

#### VERSION CONTROL BEYOND SINGLE REPOSITORIES

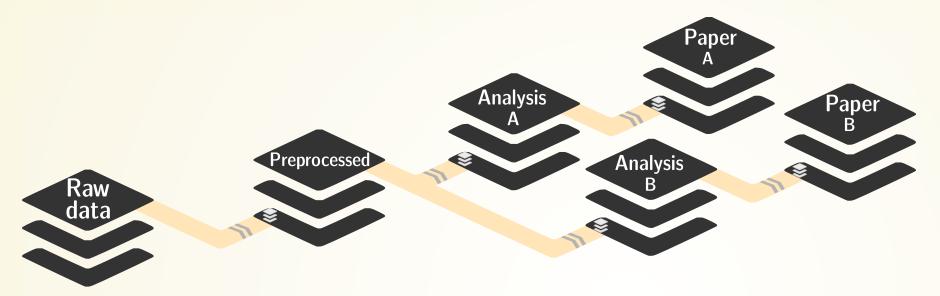
- Why are multiple repositories needed (in science)?
  - Size impacts I/O and logistics
    - Git can struggle with 1M+ files or 100k+ commits
    - Filesystems (licensing) can struggle with large numbers of inodes
  - Target audience is different
    - Public vs. private or personal vs. anonymized data
  - Pace of evolution or access patterns are different
    - "Factual" raw data vs. choices of (pre-)processing
    - Completed acquisition vs. ongoing study



 A single repository is not enough, but Git/Git-annex are not optimized for such use cases

### DATASET NESTING

Seamless nesting mechanisms:

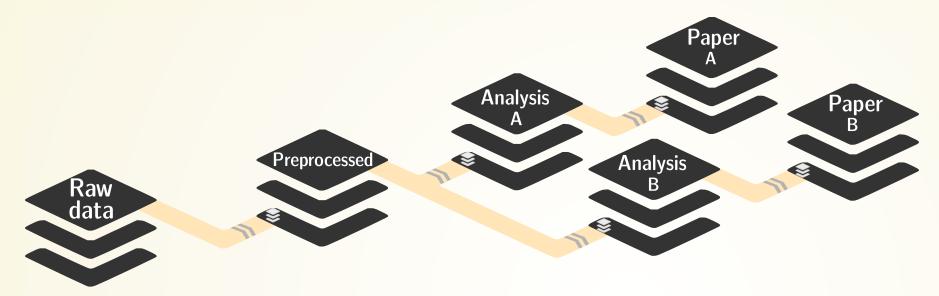


Nest modular datasets to create a linked hierarchy of datasets, and enable recursive operations throughout the hierarchy

- hierarchies of datasets in super-/sub-dataset relationships
- based on Git submodules, but more seamless: Mono-repo feel thanks to recursive operations

#### DATASET NESTING

Seamless nesting mechanisms:



Nest modular datasets to create a linked hierarchy of datasets, and enable recursive operations throughout the hierarchy

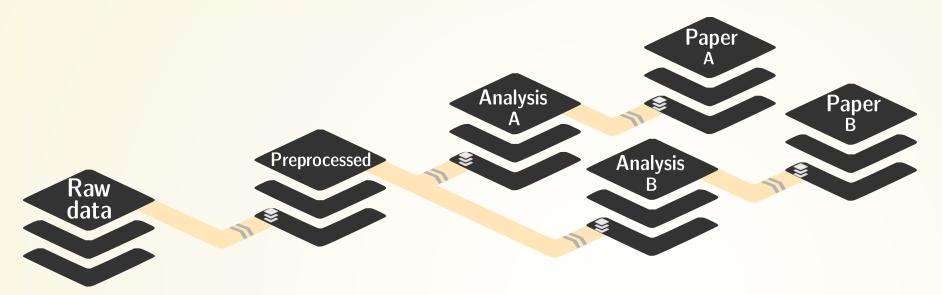
- hierarchies of datasets in super-/sub-dataset relationships
- based on Git submodules, but more seamless: Mono-repo feel thanks to recursive operations
- Overcomes scaling issues with large amounts of files

```
adina@bulk1 in /ds/hcp/super on git:master) datalad status --annex -r 15530572 annex'd files (77.9 TB recorded total size) nothing to save, working tree clean
```

(github.com/datalad-datasets/human-connectome-project-openaccess)

#### DATASET NESTING

Seamless nesting mechanisms:



Nest modular datasets to create a linked hierarchy of datasets, and enable recursive operations throughout the hierarchy

- hierarchies of datasets in super-/sub-dataset relationships
- based on Git submodules, but more seamless: Mono-repo feel thanks to recursive operations
- Overcomes scaling issues with large amounts of files

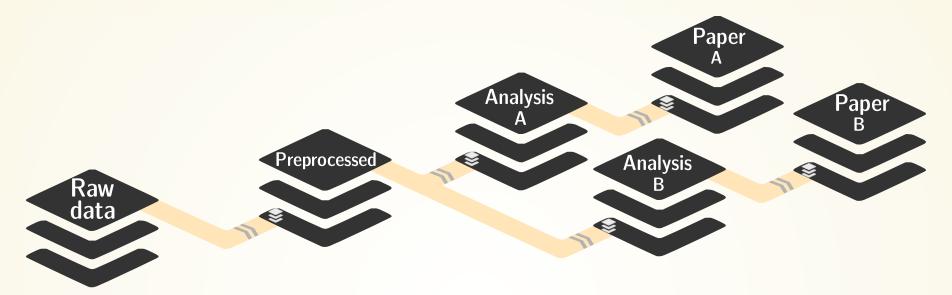
```
adina@bulk1 in /ds/hcp/super on git:master) datalad status --annex -r 15530572 annex'd files (77.9 TB recorded total size) nothing to save, working tree clean
```

(github.com/datalad-datasets/human-connectome-project-openaccess)

 Modularizes research components for transparency, reuse, and access management

#### INTUITIVE DATA ANALYSIS STRUCTURE

You can link datasets together in superdataset-subdataset hierarchies:



Nest modular datasets to create a linked hierarchy of datasets, and enable recursive operations throughout the hierarchy

```
copy

2 # we can install analysis input data as a subdataset to the dataset

3 $ datalad clone -d . https://github.com/datalad-handbook/iris_data.git input/

4 [INFO ] Scanning for unlocked files (this may take some time)

5 [INFO ] Remote origin not usable by git-annex; setting annex-ignore

6 install(ok): input (dataset)

7 add(ok): input (file)

8 add(ok): .gitmodules (file)

9 save(ok): . (dataset)

10 action summary:

11 add (ok: 2)

12 install (ok: 1)

13 save (ok: 1)
```

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

"Shit, why buttons did I click and in which order did I use all those tools?"

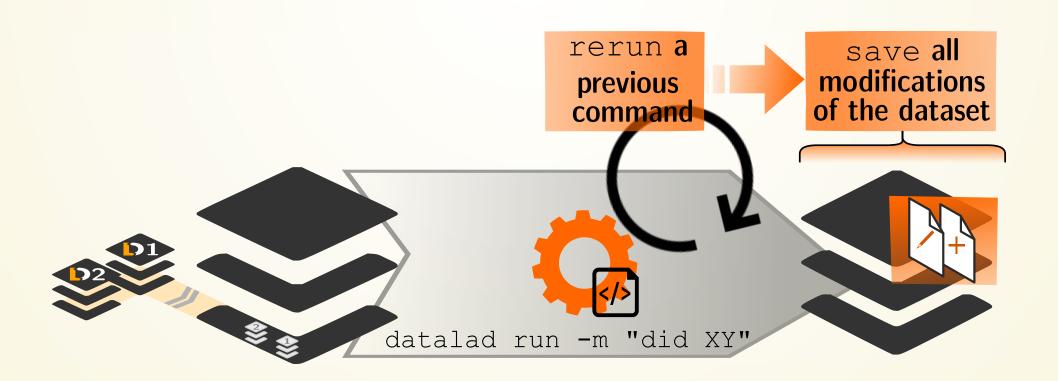


datalad run wraps around anything expressed in a command line call and saves the dataset modifications resulting from the execution.



datalad run wraps around anything expressed in a command line call and saves the dataset modifications resulting from the execution.

datalad rerun repeats captured executions. If the outcomes differ, it saves a new state of them.



datalad run wraps around anything expressed in a command line call and saves the dataset modifications resulting from the execution.

datalad rerun repeats captured executions. If the outcomes differ, it saves a new state of them.

datalad containers-run executes command line calls inside a tracked software container and saves the dataset modifications resulting from the execution.



```
сору
"python3 code/extract lc timeseries.py"
```

```
1 $ datalad containers-run \
                                                                                           сору
   --message "Time series extraction from Locus Coeruleus"
   --container-name nilearn \
   --input 'mri/* bold.nii' \
   --output 'sub-*/LC timeseries run-*.csv' \
   "python3 code/extract lc timeseries.py"
```

```
сору
-- Git commit --
    commit 5a7565a640ff6de67e07292a26bf272f1ee4b00e
               Adina Wagner adina.wagner@t-online.de
    AuthorDate: Mon Nov 11 16:15:08 2019 +0100
                Adina Wagner adina.wagner@t-online.de
    Commit:
    CommitDate: Mon Nov 11 16:15:08 2019 +0100
    [DATALAD RUNCMD] Time series extraction from Locus Coeruleus
    === Do not change lines below ===
     "cmd": "singularity exec --bind {pwd} .datalad/environments/nilearn.simg bash..",
     "dsid": "92ea1faa-632a-11e8-af29-a0369f7c647e",
     "inputs": [
     "mri/*.bold.nii.gz",
     ".datalad/environments/nilearn.simg"
     ],
     "outputs": ["sub-*/LC timeseries run-*.csv"],
    ^^^ Do not change lines above ^^^
 sub-01/LC timeseries run-1.csv | 1 +
```

```
1 $ datalad rerun 5a7565a640ff6de67

2 [INFO ] run commit 5a7565a640ff6de67; (Time series extraction from Locus Coeruleus)

3 [INFO ] Making sure inputs are available (this may take some time)

4 get(ok): mri/sub-01_bold.nii (file)

5 get(ok): mri/sub-02_bold.nii (file)

6 [...]

7 [INFO ] == Command start (output follows) =====

8 [INFO ] == Command exit (modification check follows) =====

9 add(ok): sub-01/LC_timeseries_run-*.csv(file)

10 add(ok): sub-02/LC_timeseries_run-*.csv (file)

11 [...]

12 action summary:

13 add (ok: 30)

14 get (ok: 30)

15 save (ok: 2)

16 unlock (ok: 30)
```

Scientific building blocks are not static.

Scientific building blocks are not static.
Version control beyond text

Scientific building blocks are not static.
Version control beyond text

Science is build from modular units.

Scientific building blocks are not static.
Version control beyond text
Science is build from modular units.
Nesting

Scientific building blocks are not static.

Version control beyond text

Science is build from modular units.

Nesting

Science is exploratory, iterative, multi-stepped, and complex.

Scientific building blocks are not static.

Version control beyond text

Science is build from modular units.

Nesting

Science is exploratory, iterative, multi-stepped, and complex.

Provenance

Scientific building blocks are not static.

Version control beyond text

Science is build from modular units.

Nesting

Science is exploratory, iterative, multi-stepped, and complex.

Provenance

Science is collaborative.

Scientific building blocks are not static.

Version control beyond text

Science is build from modular units.

Nesting

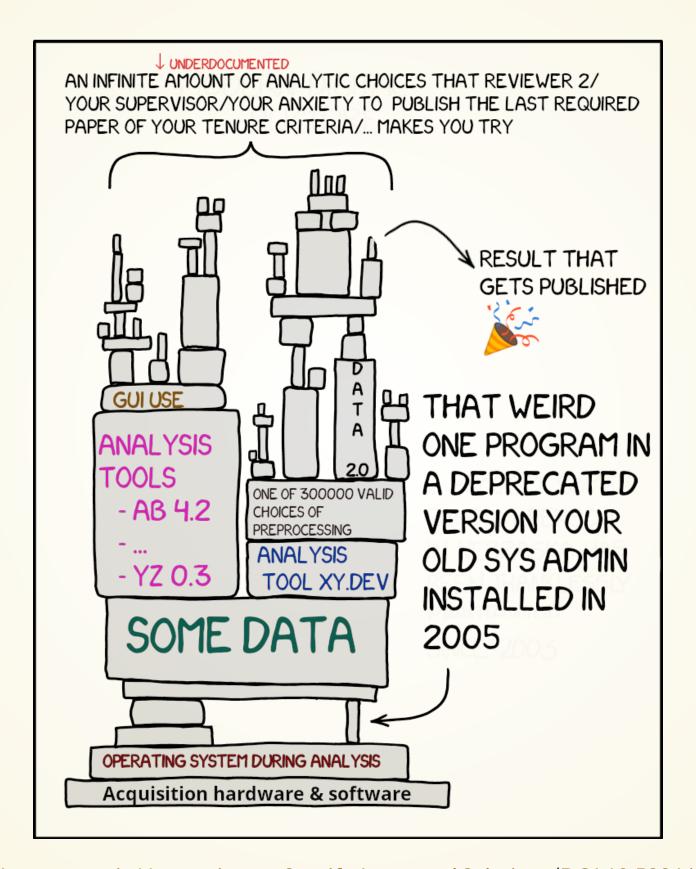
Science is exploratory, iterative, multi-stepped, and complex.

Provenance

Science is collaborative.

**Transport logistics** 

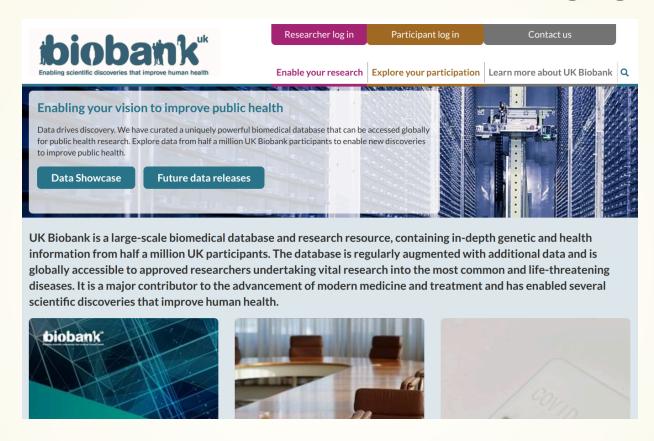
## RESEARCH DATA MANAGEMENT IS TIED TO REPRODUCIBILITY



Reproducibility Management in Neuroscience - Specific Issues and Solutions (DOI 10.5281/zenodo.4285927)

#### FAIRLY BIG: SCALING UP

Objective: Process the UK Biobank (imaging data)

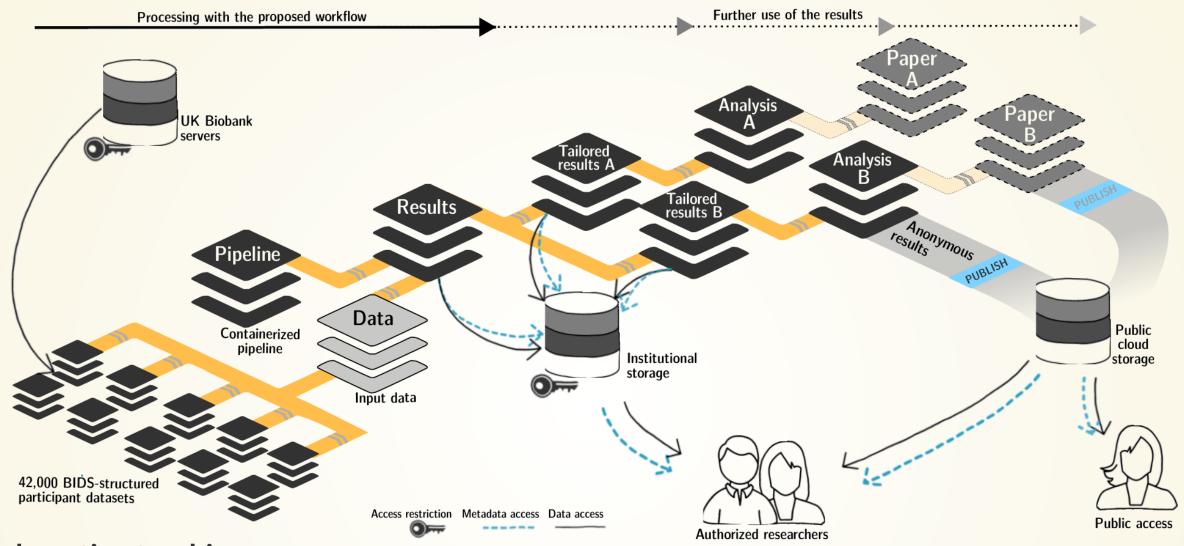


- 76 TB in 43 million files in total
- 42,715 participants contributed personal health data
- Strict DUA
- Custom binary-only downloader
- Most data records offered as (unversioned) ZIP files

#### **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

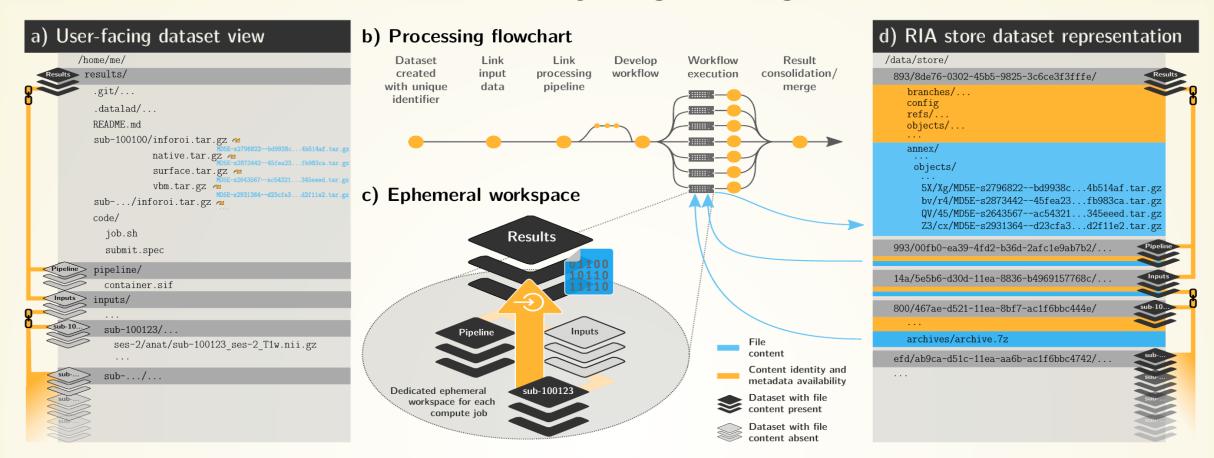
## FAIRLY BIG SETUP



### **Exhaustive tracking**

- datalad-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

## FAIRLY BIG WORKFLOW



#### portability

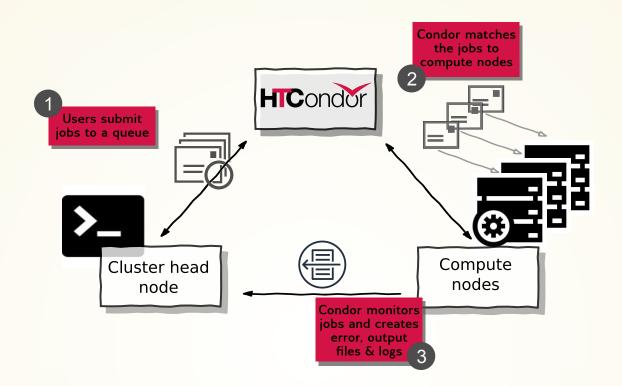
- Parallel processing: 1 job = 1 subject (number of concurrent jobs capped at the capacity of the compute cluster)
- Each job is computed in a ephemeral (short-lived) dataset clone, results are pushed back: Ensure exhaustive tracking & portability during computation
- Content-agnostic persistent (encrypted)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

## 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

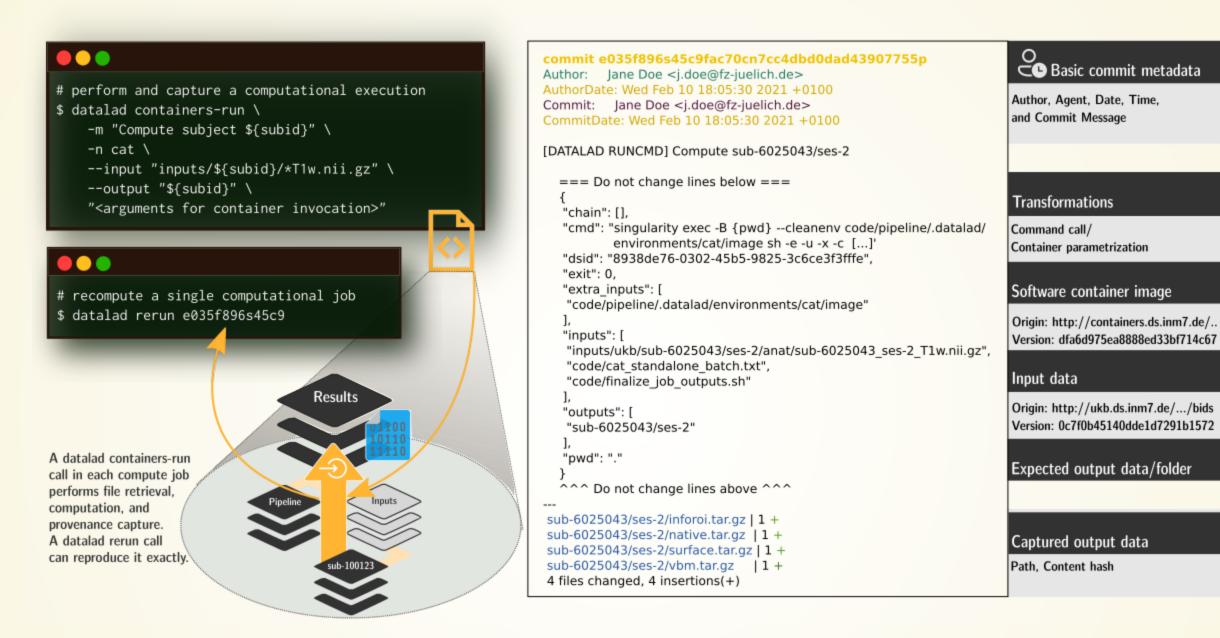
## FAIRLY BIG WORKFLOW



#### 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

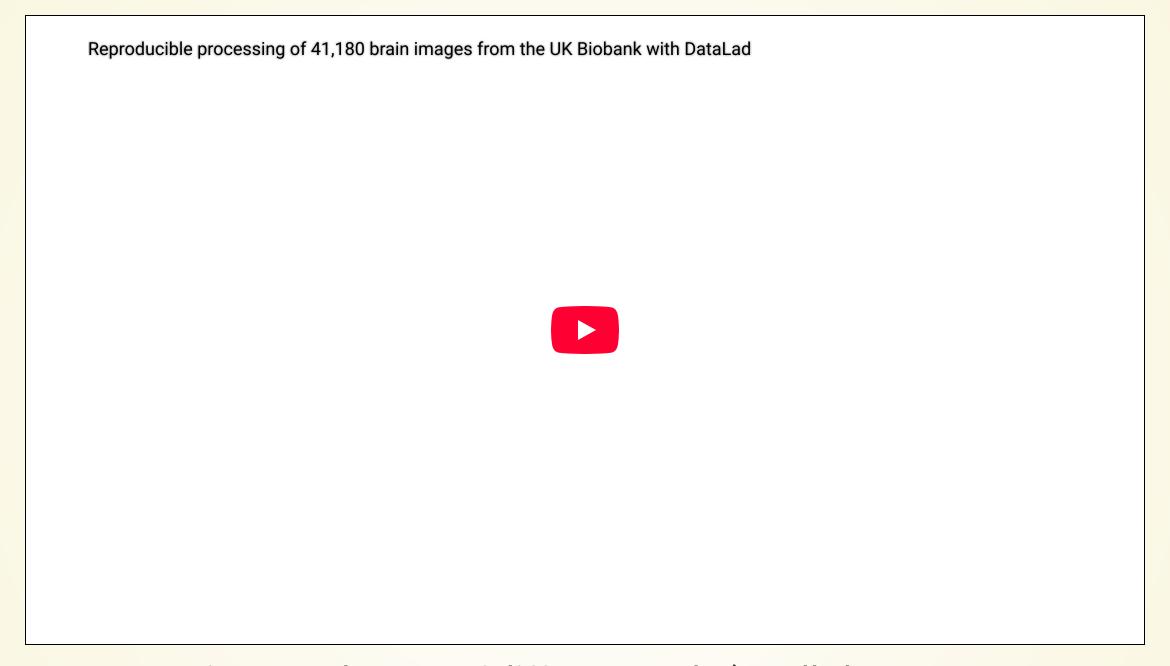
## FAIRLY BIG PROVENANCE CAPTURE



#### **Provenance**

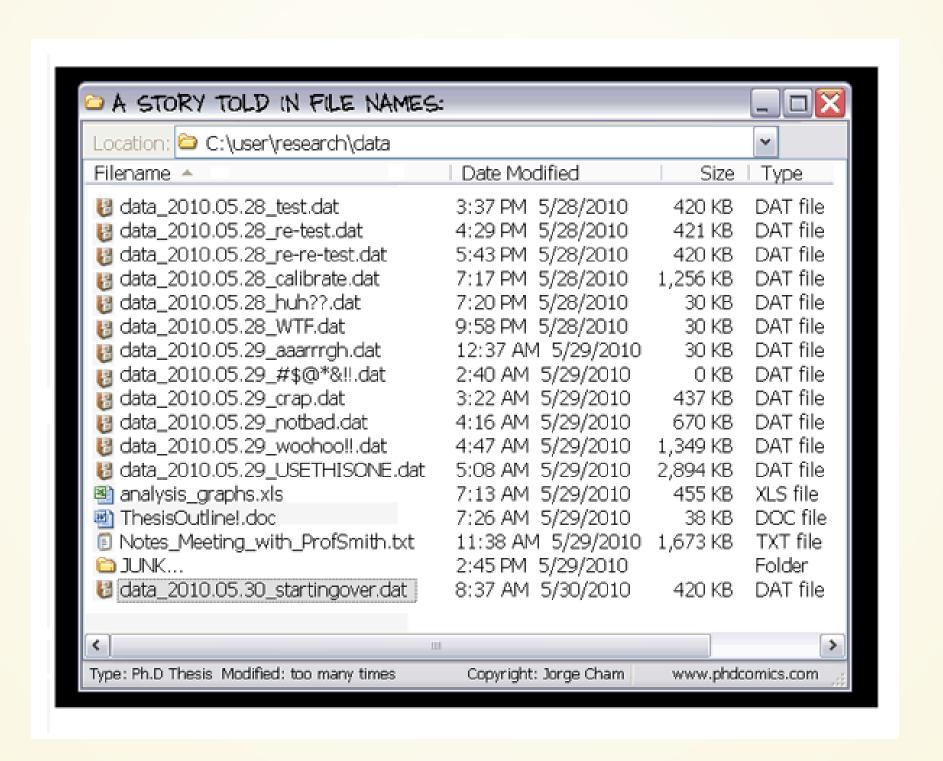
- Every single pipeline execution is tracked
- Execution in ephemeral workspaces ensures results individually reproducible without HPC access

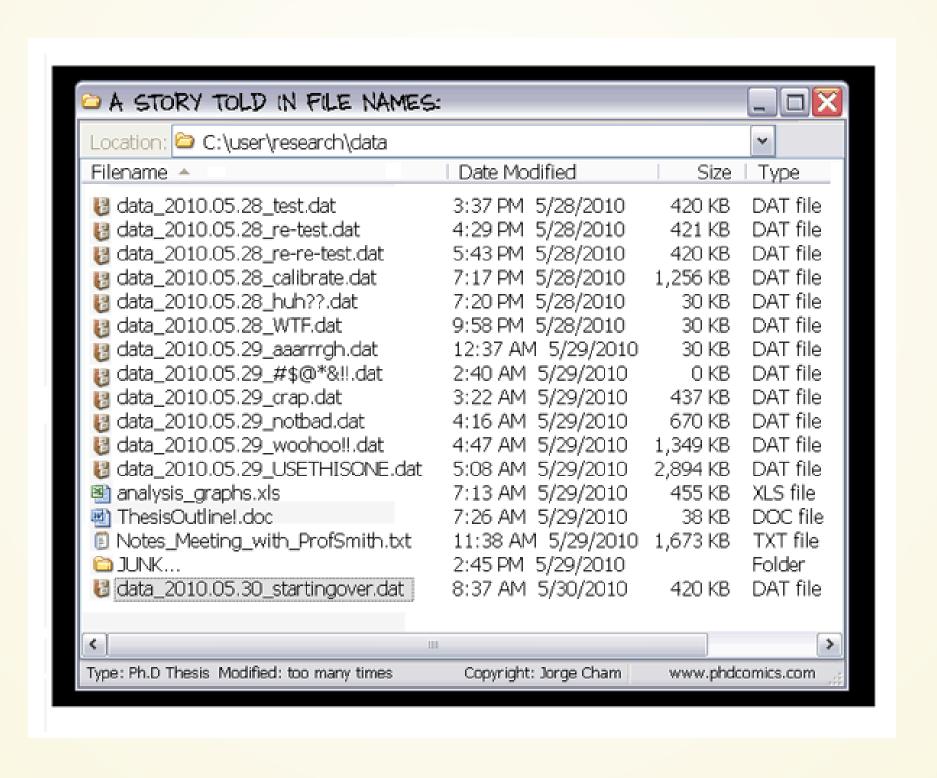
## FAIRLY BIG MOVIE



- 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

## **CURRENT AND FUTURE DEVELOPMENTS**





Mar 2019

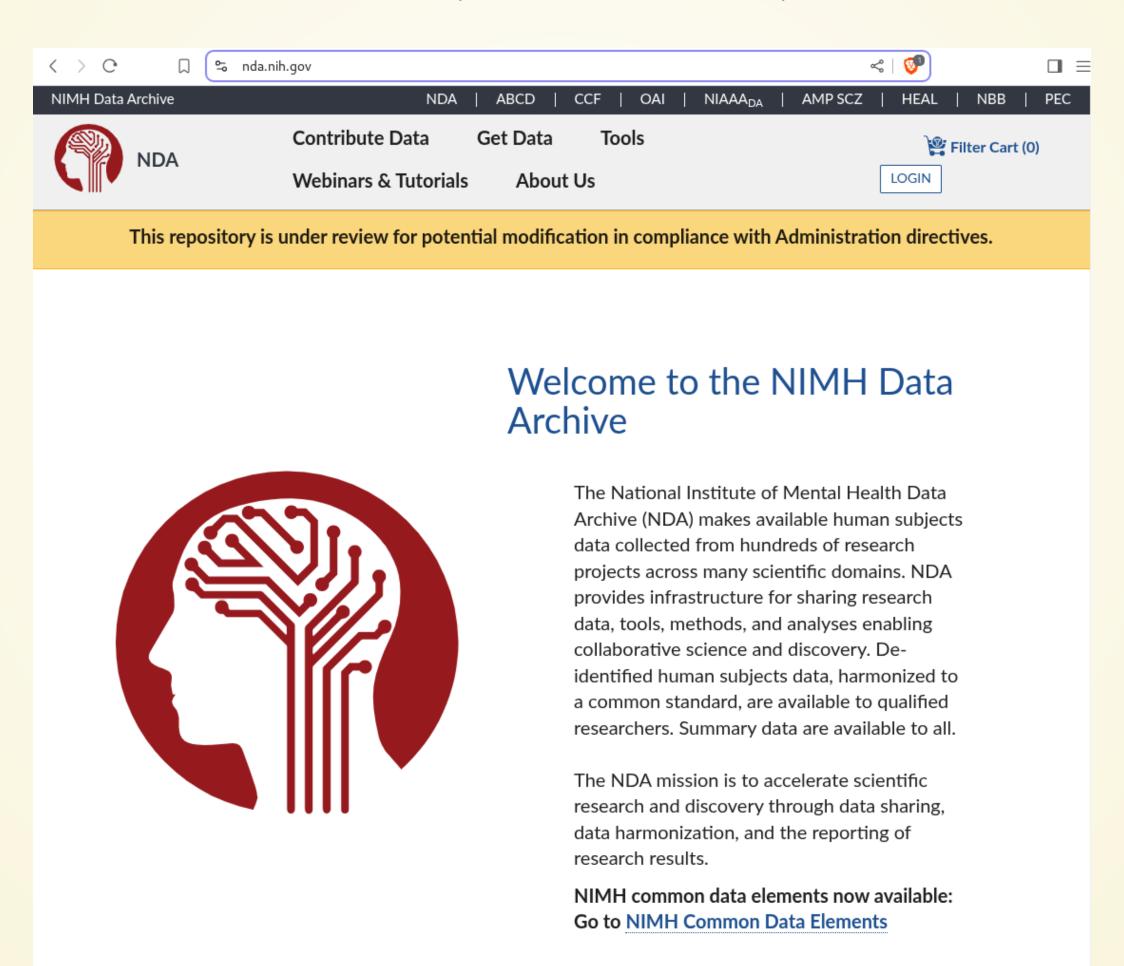
### ABCD Data Release 2.0 available now on the NIMH Data Archive!

The second annual curated ABCD data release 2.0 is available now on the NIMH Data Archive. Data Release 2.0 includes baseline data on the full participant cohort, ages 9-10 years. This release contains much of the same type of data as in Data Release 1.1, and also includes genotypic data for the first time. Smokescreen genotyping array data are available on almost 11,000 participants. These include common variations, as well as variations associated with addiction, smoking behavior and nicotine metabolism.

ABCD Study data will be released annually. The next data release will be in early summer 2020 and will include the first longitudinal data from the 6-month and 1-year follow-up assessments.

#### Mar 2025

www.404media.co/nih-archives-repositories-marked-for-review-for-potential-modification



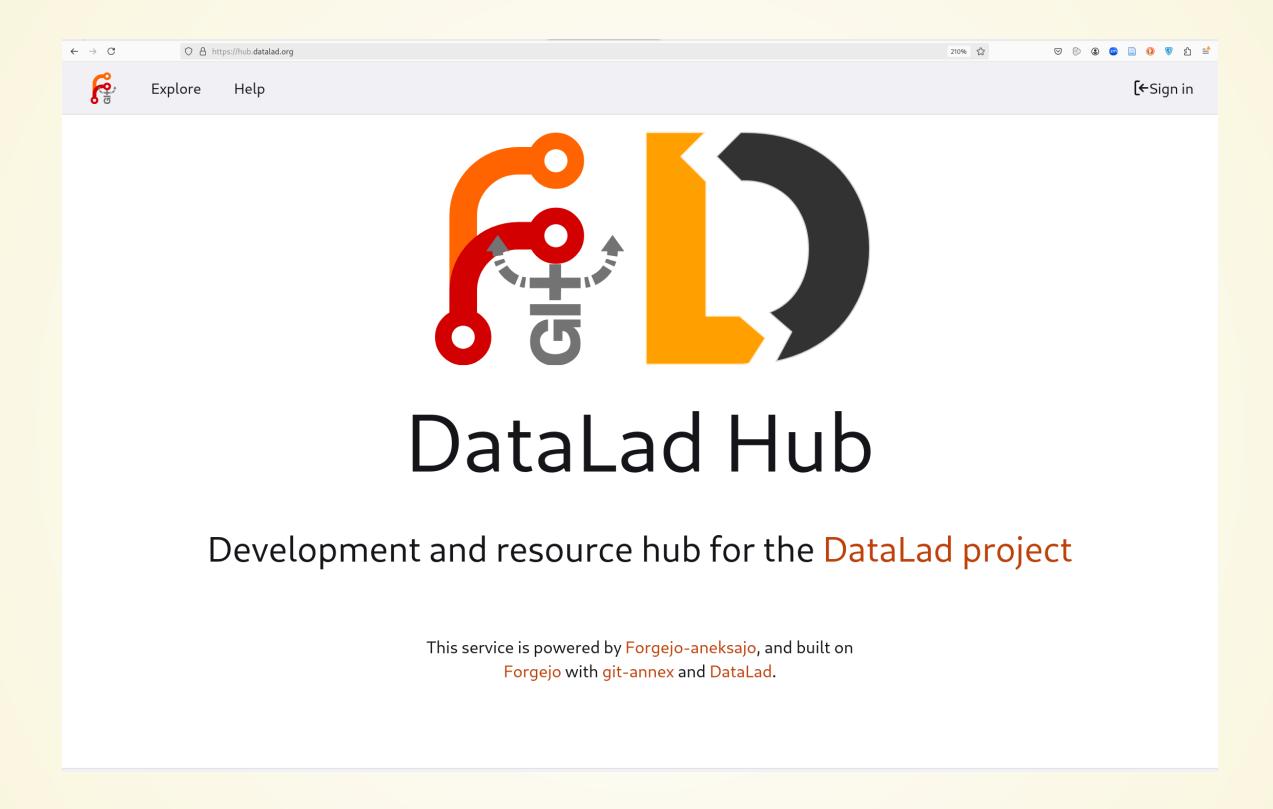
## FREEDOM? CHOSE DECENTRALIZATION

- Infrastructure is ephemeral:
  - Change of institutional contracts
  - Change of affiliations
  - Geopolitical developments?
- DataLad datasets are portable
  - Effortless migrations to different Git or data hosting
  - Versioning allows for integrity checks

Delineation and advantages of decentral versus central RDM: Hanke et al., (2021). In defense of decentralized research data management

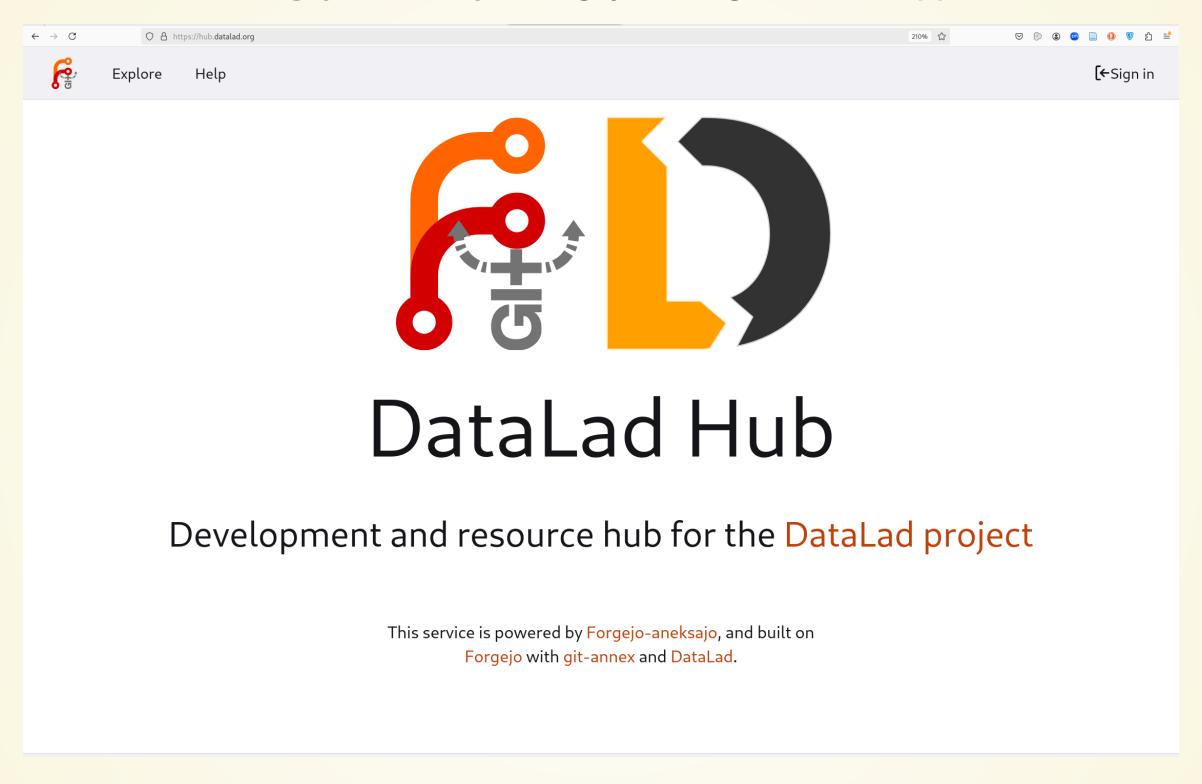
## GOING SELF-HOSTED WITH FORGEJO-ANEKSAJO

• Forgejo (forgejo.org): Fork of Gitea



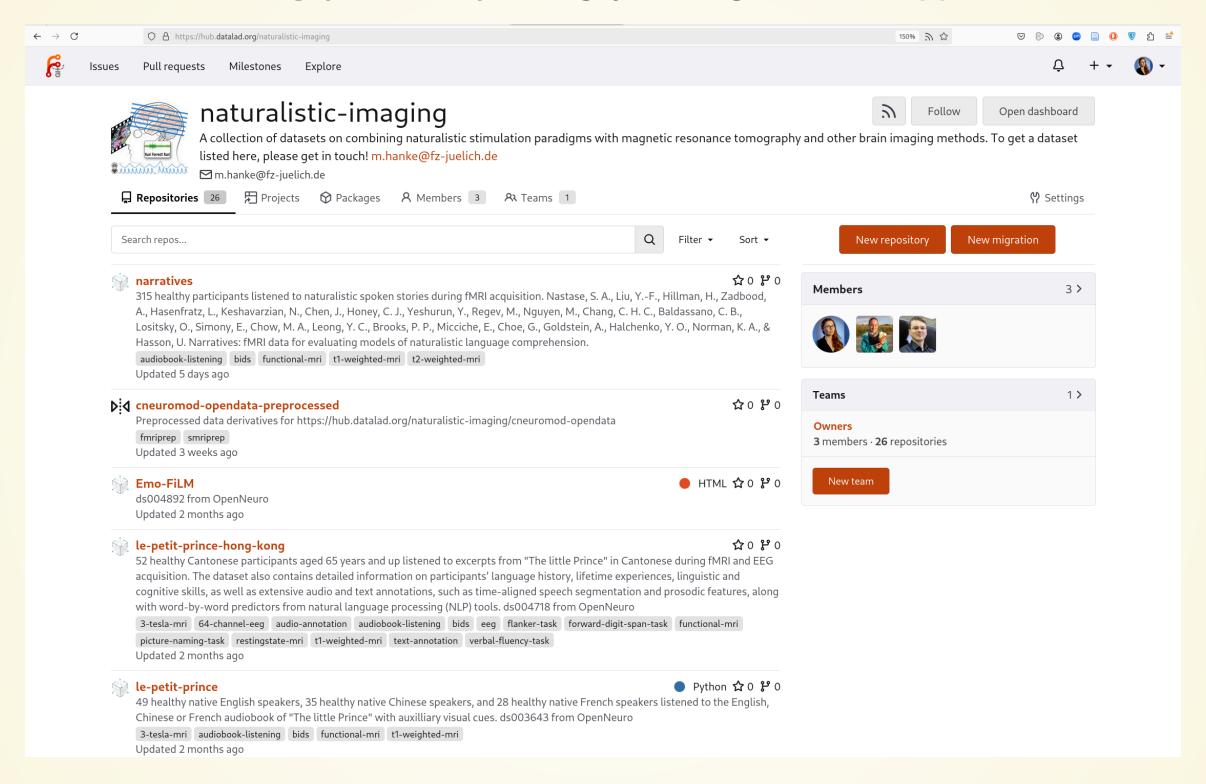
## GOING SELF-HOSTED WITH FORGEJO-ANEKSAJO

- Forgejo (forgejo.org): Fork of Gitea
- Forgejo-aneksajo: Forgejo with git-annex support

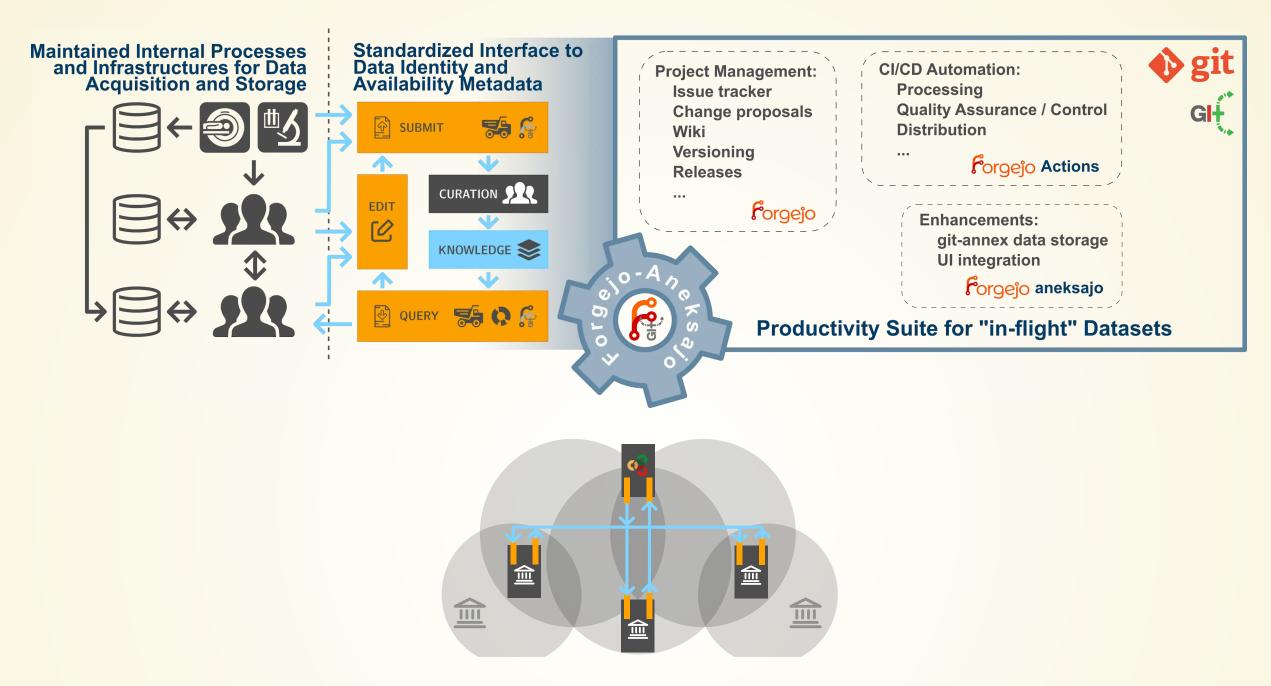


## GOING SELF-HOSTED WITH FORGEJO-ANEKSAJO

- Forgejo (forgejo.org): Fork of Gitea
- Forgejo-aneksajo: Forgejo with git-annex support



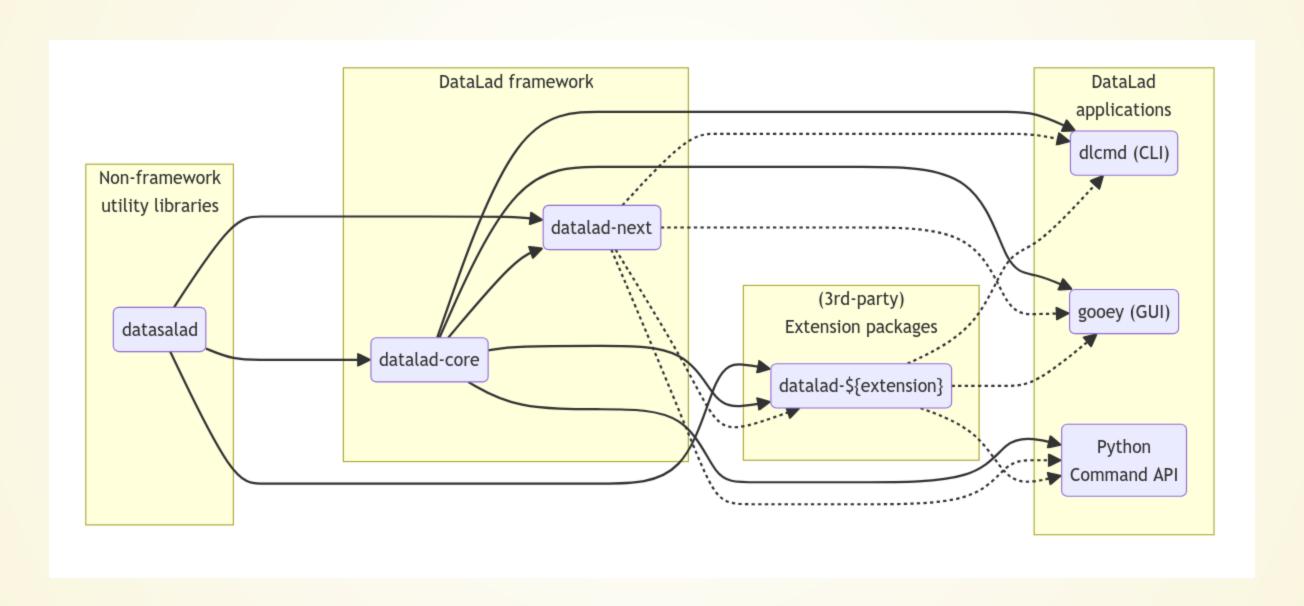
## FULL-STACK RDM FOR INDEPENDENT, INTEROPERABLE COLLABORATORS



scale-free organization: consortium, institution, lab, researcher

- maximum contributor benefit
- solutions, e.g., atris.fz-juelich.de, hub.trr379.de
- minimum contributor cost
- self-hostable, independently governed
   self-contained contributor scopes, not inheriting complexity of others

## DEVELOPMENT ROADMAP



# JOIN US!

#### Distribits 2025

- International conference on technologies for distributed data management
- 2 day conference plus single-day Hackathon
- @ Haus der Universität Düsseldorf
- Registration open until May 1st

3.-25. Oct 2025 Düsseldor

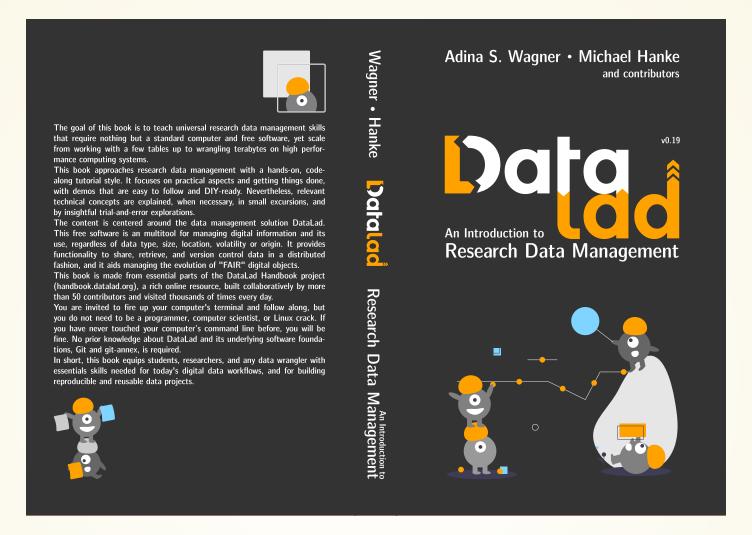
Distribits.Live

201551110115 26d166007fdfd10422025520069df060692667

## DATALAD CONTACT AND MORE INFORMATION

Website + Demos	http://datalad.org		
Documentation	http://handbook.datalad.org		
Talks and tutorials	https://youtube.com/datalad		
Development	http://github.com/datalad		
Support	https://matrix.to/#/#datalad:matrix.org		
Open data	http://datasets.datalad.org		
Mastodon	@datalad@fosstodon.org		

## EXTENSIVE DOCUMENTATION AND TRAINING MATERIALS



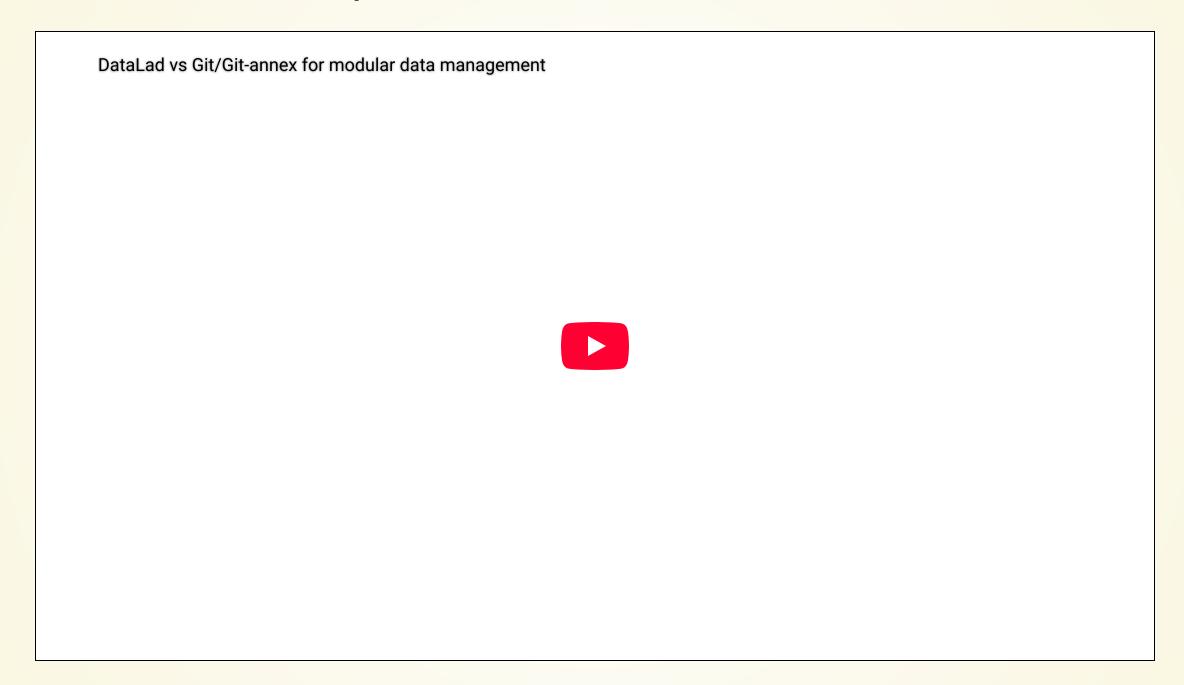
https://handbook.datalad.org (or ISBN 979-8857037973)

- educational materials on technologies targeting researchers, not developers (executable paper, student surpervisor workflow, ...)
- handbook on concepts, workflows, and use cases
- weekly public (virtual) office hour

# THANKS!



# TALK IS CHEAP, SHOW ME THE CODE: GIT VS. DATALAD



https://www.youtube.com/watch?v=Yrg6DgOcbPE