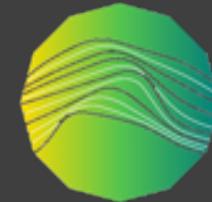




Human Brain Project



EBRAINS

# Provenance tracking with containerized pipelines for reproducible research

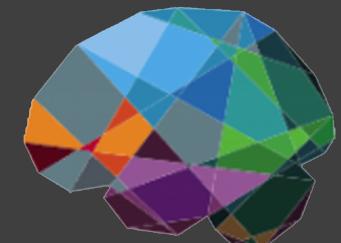


Michael Schirner

Brain Simulation Section (PI: Petra Ritter)

Charité—Universitätsmedizin Berlin

Berlin, July 13th, 2021



THE VIRTUAL BRAIN.

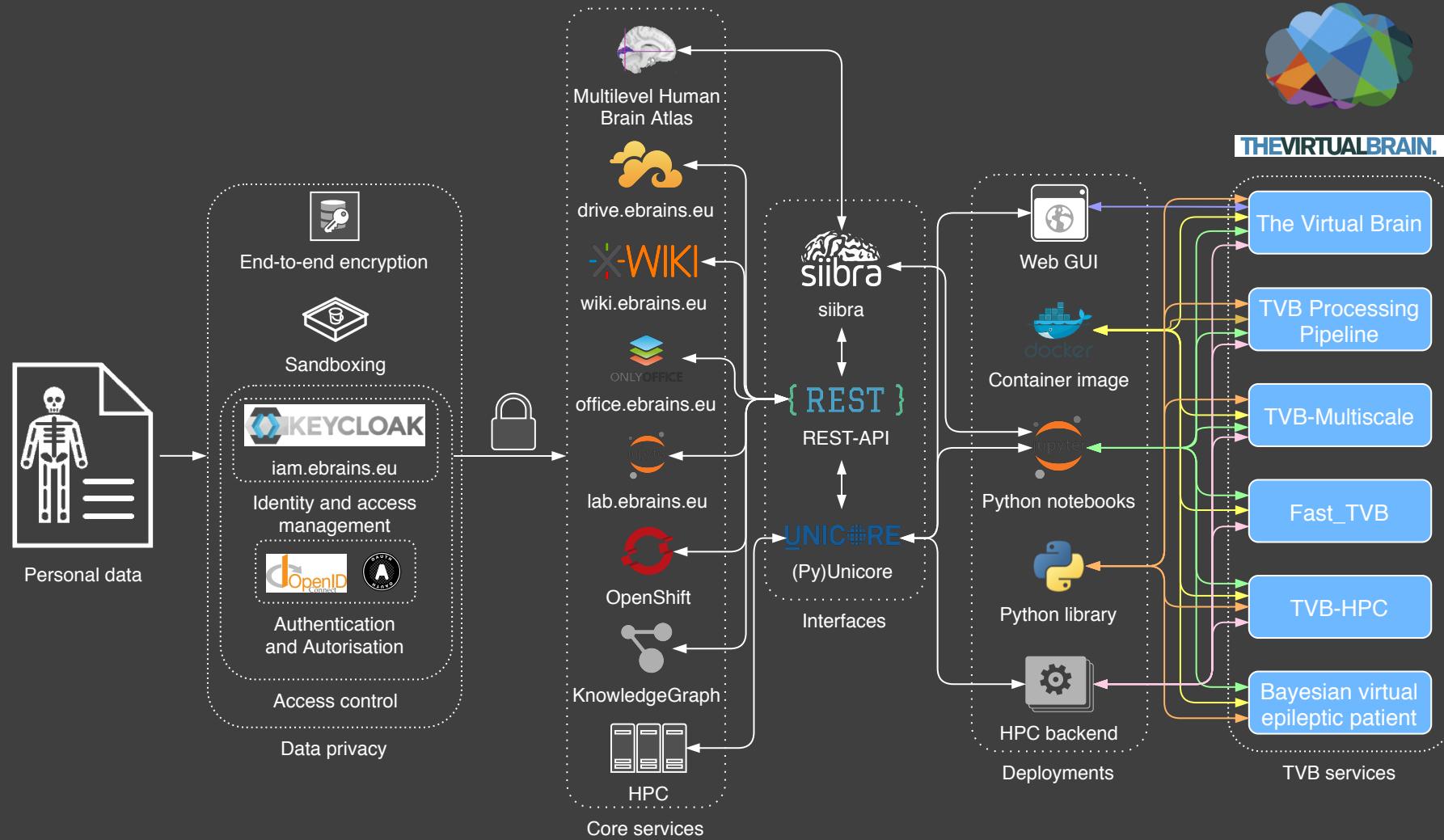
# Outline

1. TVB on EBRAINS Image Processing Container Pipeline
2. Making the pipeline more reproducible and capturing more provenance with DataLad

Acknowledgments: much of the material shown here was created by the DataLad community. Special thanks to

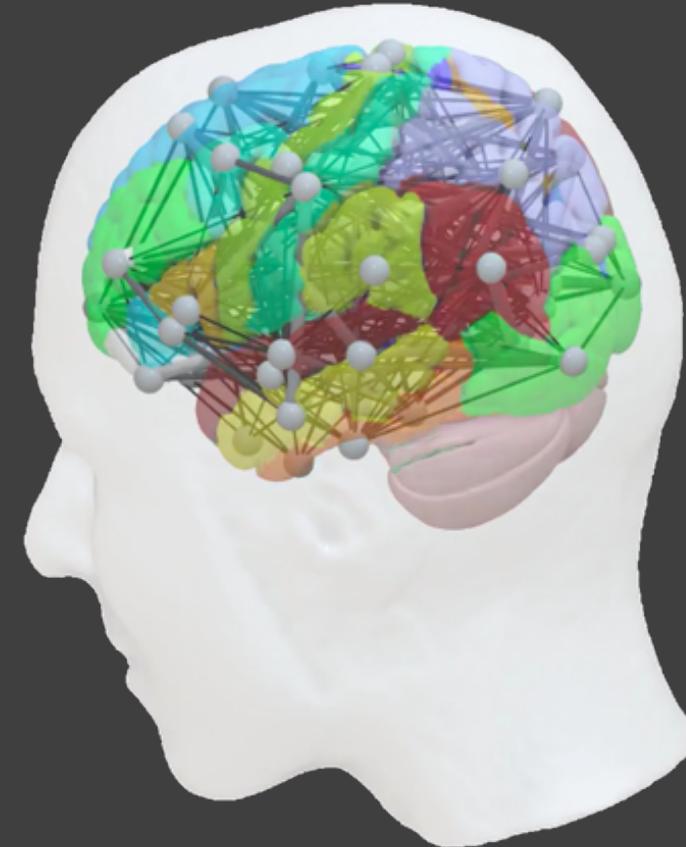
- Michael Hanke
- Adina Wagner.

**Thank You!**



# TVB-on-EBRAINS cloud infrastructure

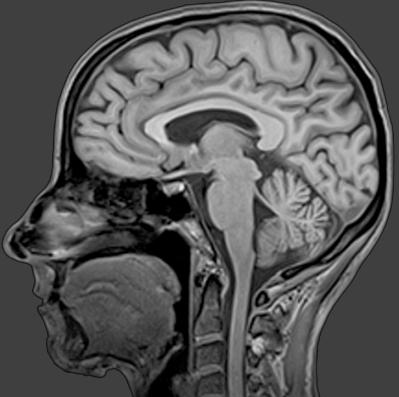
# TVB Image Processing Pipeline



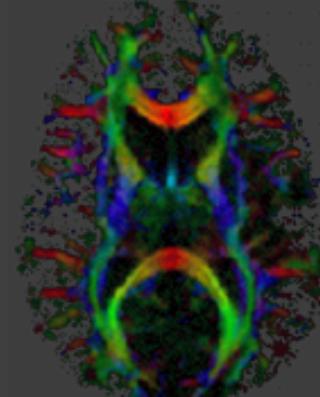
magnetic resonance imaging

brain network model

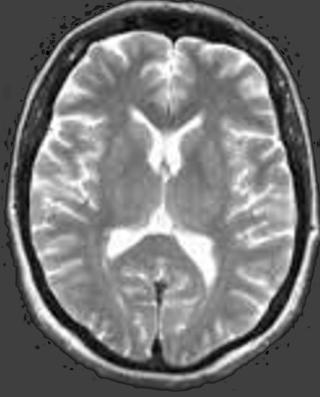
# Pipeline input



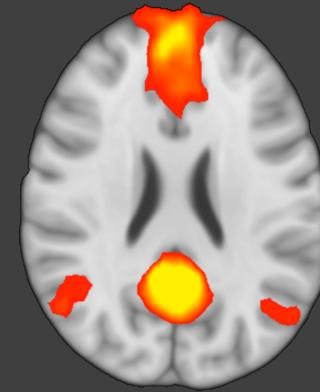
T1-weighted MRI



Diffusion-weighted MRI

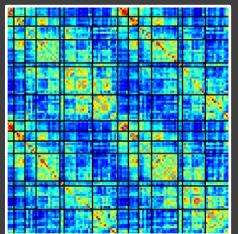


T2-weighted MRI

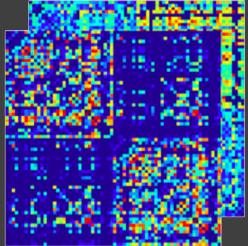


functional MRI

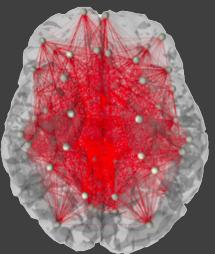
# Pipeline output



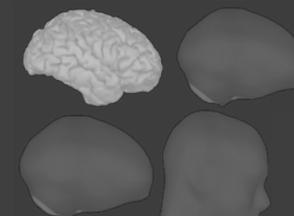
Functional connectivity



Structural connectivity



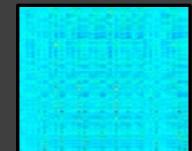
region-wise fMRI



Surface triangulations

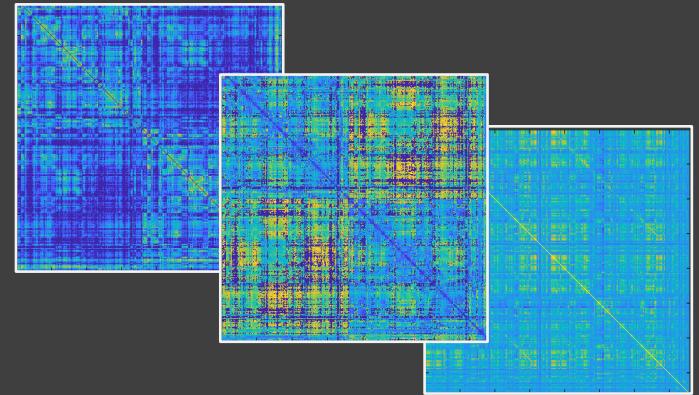
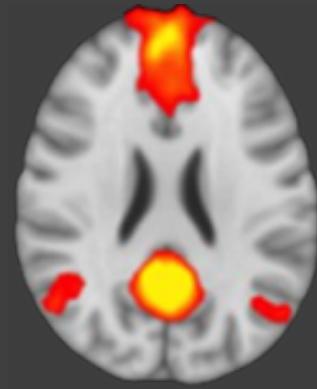


Brain maps



Projection matrices

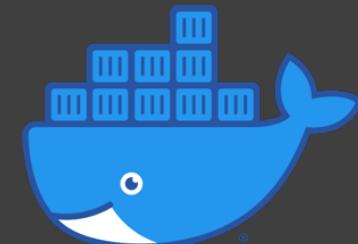
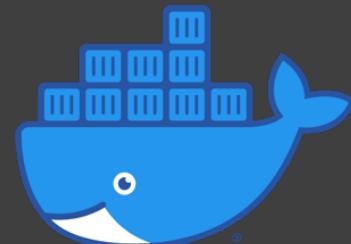
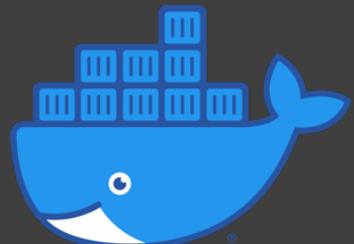
# TVB Image Processing Pipeline



`bids/mrtrix3_connectome`

`nipreps/fmriprep`

`thevirtualbrain/tvb-pipeline-converter`



168 lines (155 sloc) | 8.52 KB

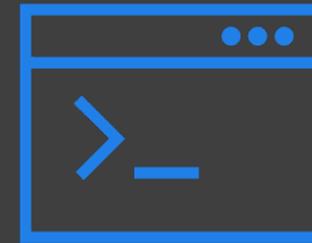
```
1 FROM ubuntu:18.04
2 MAINTAINER Robert E. Smith <robert.smith@florey.edu.au>
3
4 # Core system capabilities required
5 RUN apt-get update && DEBIAN_FRONTEND=noninteractive apt-get
6     bc \
7     build-essential \
8     curl \
9     dc \
10    git \
11    libegl1-mesa-dev \
12    libopenblas-dev \
13    nano \
14    perl-modules-5.26 \
15    python2.7 \
16    python3 \
17    tar \
18    tcsh \
19    tzdata \
20    unzip \
21    wget
22
23 # PPA for newer version of nodejs, which is required for bids
24 RUN curl -sL https://deb.nodesource.com/setup_12.x -o nodesou
25     bash nodesource_setup.sh && \
26     rm -f nodesource_setup.sh && \
27     apt-get install -y nodejs
28
29 # NeuroDebian setup
```

```
137 SUBJECTS_DIR=/opt/freesurfer/subjects \
138 FSLDIR=/opt/fsl \
139 FSLOUTPUTTYPE=NIFTI \
140 FSLMULTIFILEQUIT=TRUE \
141 FSLTCLSH=/opt/fsl/bin/fsltclsh \
142 FSLWISH=/opt/fsl/bin/fslwish \
143 LD_LIBRARY_PATH=/opt/fsl/lib:$LD_LIBRARY_PATH \
144 PATH=/opt/mrtrix3/bin:/usr/lib/ants:/opt/freesurfer/bin:/opt/freesurfer/mn
145 PYTHONPATH=/opt/mrtrix3/lib:$PYTHONPATH
146
147 # MRtrix3 setup
148 # Commitish is 3.0.2 plus relevant hotfix
149 RUN git clone https://github.com/MRtrix3/mrtrix3.git /opt/mrtrix3 && \
150     cd /opt/mrtrix3 && \
151     git checkout 4ab54489f40997f7da1e1915c2adde3373cf6039 && \
152     python3 configure --nogui && \
153     python3 build --persistent --nopaginate && \
154     git describe --tags > /mrtrix3_version && \
155     rm -rf .git/ cmd/ core/ src/ testing/ tmp/ && \
156     cd /
157
158 # Acquire extra MRtrix3 data
159 RUN wget -q "https://osf.io/v8n5g/download" -O /opt/mrtrix3/share/mrtrix3/labe
160     wget -q "https://osf.io/ug2ef/download" -O /opt/mrtrix3/share/mrtrix3/labe
161
162 # Acquire script to be executed
163 COPY mrtrix3_connectome.py /mrtrix3_connectome.py
164 RUN chmod 775 /mrtrix3_connectome.py
165
166 COPY version /version
167
168 ENTRYPOINT ["/mrtrix3_connectome.py"]
```

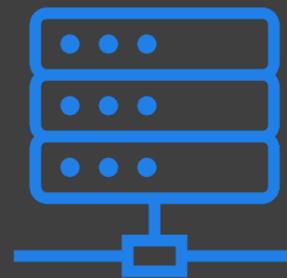
Upload  
encrypted MRI  
to Drive



Configure  
processing with  
GUI



Download  
encrypted results



Forward  
to HPC



DataLad  
commands  
here

Orchestator performs  
de-/encryption & runs  
the containers on HPC

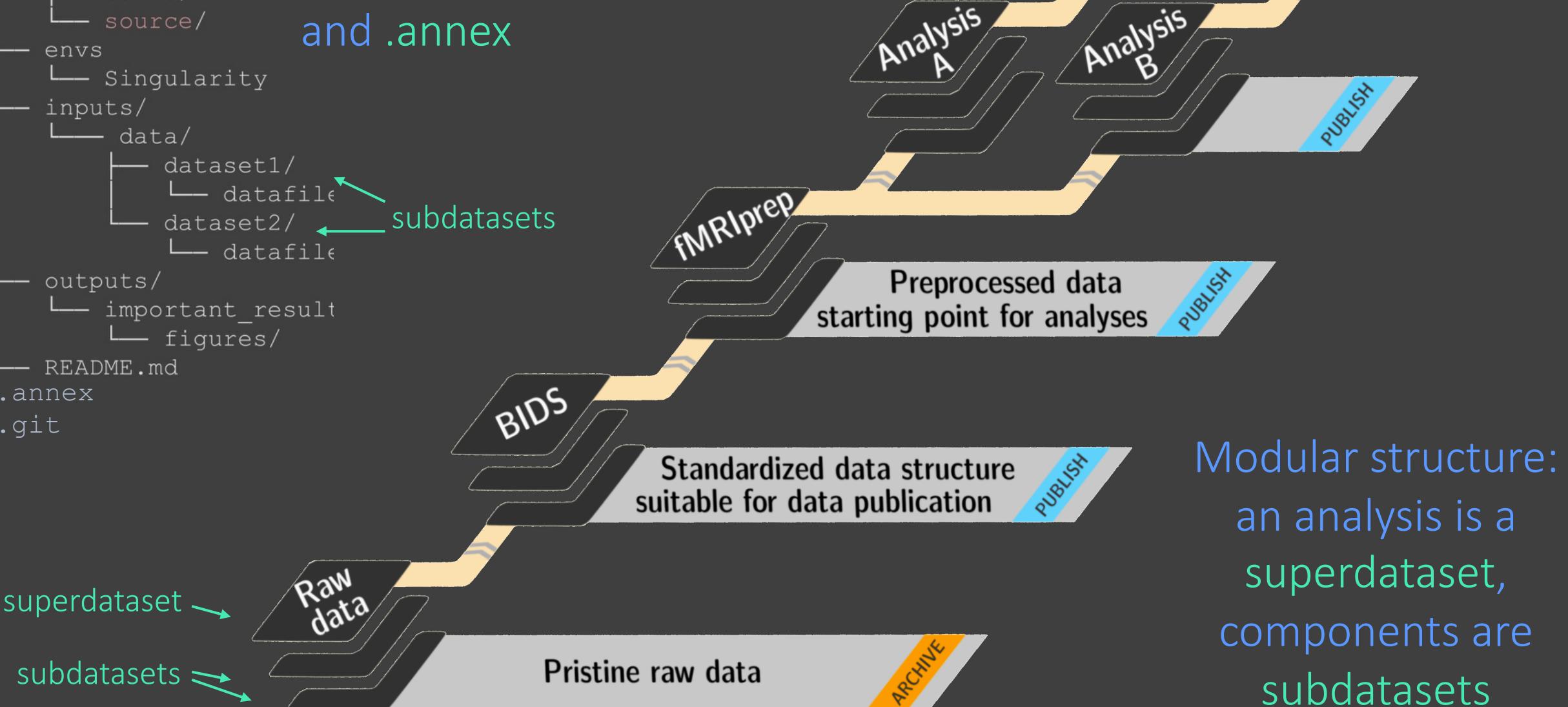


```
superdataset
├── code/
│   └── tests/
│       └── myscript.py
├── docs
│   └── build/
│       └── source/
├── envs
└── Singularity
inputs/
└── data/
    ├── dataset1/
    │   └── datafile
    └── dataset2/
        └── datafile
outputs/
└── important_result
    └── figures/
README.md
.annex
.git
```

superdataset

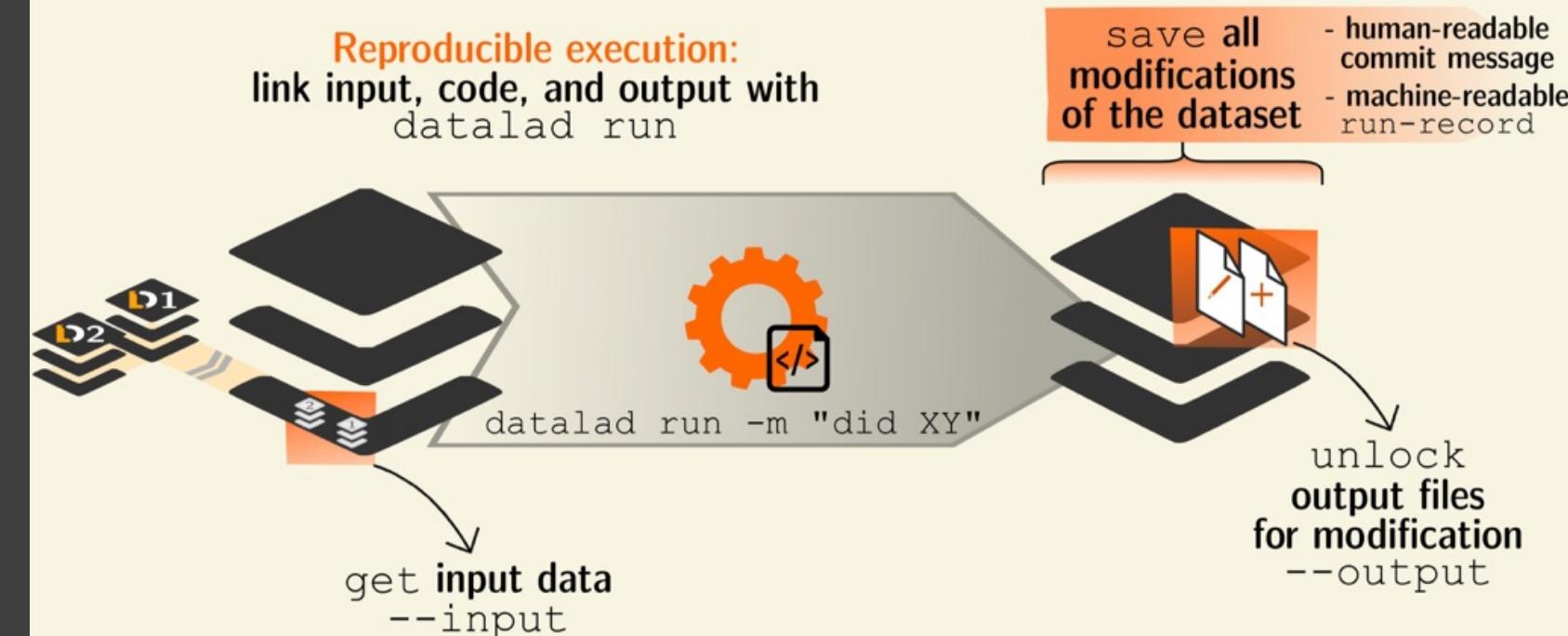
subdatasets

A DataLad dataset is a normal folder plus .git and .annex

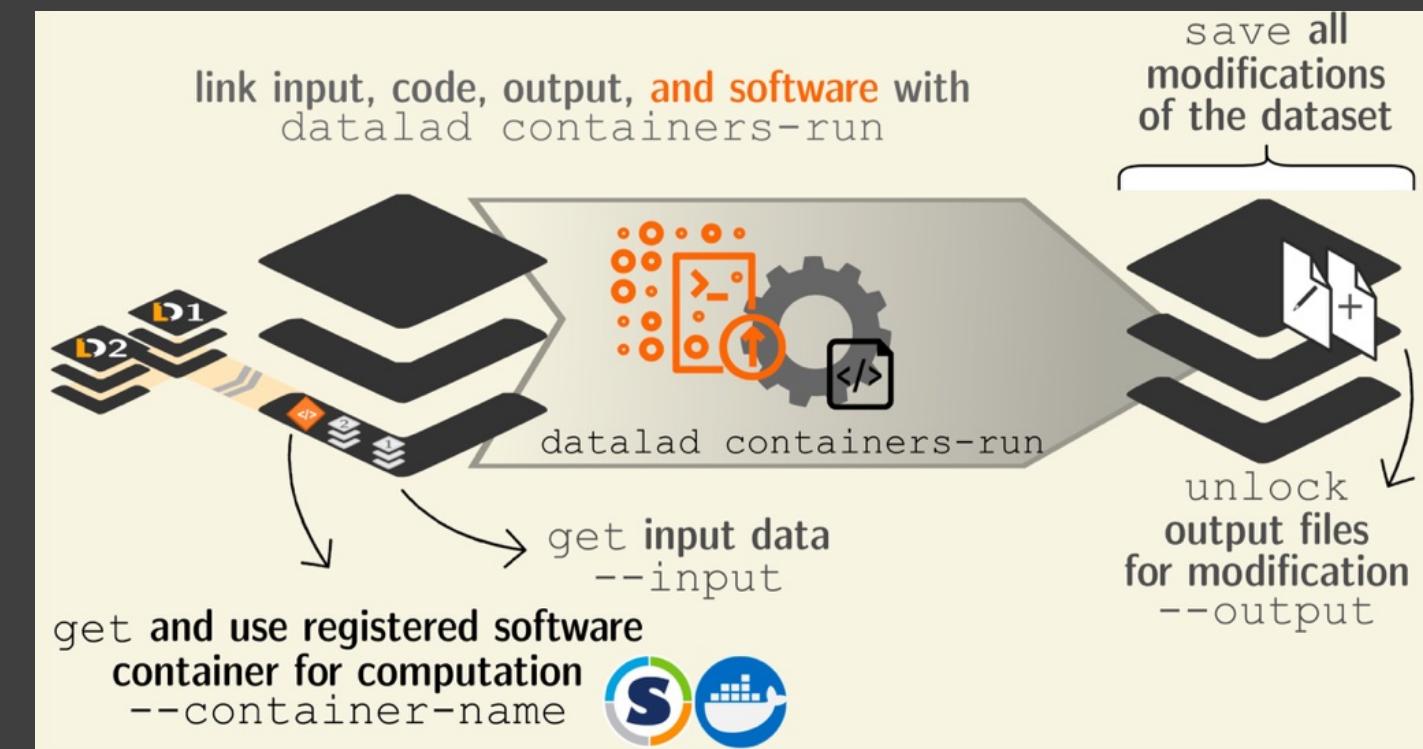


Modular structure:  
an analysis is a superdataset,  
components are subdatasets

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



Containerization captures software and environment as provenance



```
$ datalad run -m "analyze iris data with classification analysis" \
--input "input/iris.csv" \ retrieved on demand
--output "prediction_report.csv" \
--output "pairwise_relationships.png" \
"python3 code/script.py"
[INFO    ] Making sure inputs are available (this may take some time)
get(ok): input/iris.csv (file) [from web...]
[INFO    ] == Command start (output follows) =====
[INFO    ] == Command exit (modification check follows) =====
add(ok): pairwise_relationships.png (file)
add(ok): prediction_report.csv (file)
save(ok): . (dataset)
action summary:
add (ok: 2)
get (notneeded: 2, ok: 1)
save (notneeded: 1, ok: 1)
```

Execute the analysis script  
in a computationally  
reproducible manner



```
$ git log  
commit df2dae9b5af184a0c463708acf8356b877c511a8 (HEAD -> master)  
Author: Adina Wagner adina.wagner@t-online.de  
Date: Tue Dec 1 11:58:18 2020 +0100
```

[DATALAD RUNCMD] analyze iris data with classification analysis

==== Do not change lines below ===

```
{  
    "chain": [],  
    "cmd": "python3 code/script.py",  
    "dsid": "9ffdbfcd-f4af-429a-b64a-0c81b48b7f62",  
    "exit": 0,  
    "extra_inputs": [],  
    "inputs": [  
        "input/iris.csv"  
    ],  
    "outputs": [  
        "prediction_report.csv",  
        "pairwise_relationships.png"  
    ],  
    "pwd": ".."  
}
```

^^^ Do not change lines above ^^^

datalad (containers)-run  
produces a machine-readable  
provenance record, identified by a  
Git commit hash

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/gh-asim/revision_2
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:ElectronicTea
2020-03-18 09:58 +0100	Michael Hanke	o Last detections
2020-03-18 09:59 +0100	Adina Wagner	o name parameter in caption

- the dataset logs everything that was done, by whom & when
- reset dataset to historic state
- add, remove or revisit individual steps

```
$ git log  
commit df2dae9b5af184a0c463708acf8356b877c511a8 (HEAD -> master)  
Author: Adina Wagner adina.wagner@t-online.de  
Date:   Tue Dec 1 11:58:18 2020 +0100
```

```
$ datalad rerun df2dae9b5af1  
datalad rerun df2dae9b5af18  
[INFO    ] run commit df2dae9; (analyze iris data...)  
[INFO    ] Making sure inputs are available (this may take some time)  
unlock(ok): pairwise_relationships.png (file)  
unlock(ok): prediction_report.csv (file)  
[INFO    ] == Command start (output follows) =====  
[INFO    ] == Command exit (modification check follows) =====  
add(ok): pairwise_relationships.png (file)  
add(ok): prediction_report.csv (file)  
action summary:
```

```
  add (ok: 2)  
  get (notneeded: 3)  
  save (notneeded: 2)  
  unlock (ok: 2)
```

- rerun recorded computations via commit hash
- helps to reproduce and verify a result (via content hash)
- granularity can be freely decided
- recompute: not everything needs to be stored or transported

# Setting up a reproducible, version-controlled, provenance-tracked workflow

## Install DataLad

With Miniconda DataLad can be installed with just three commands

```
ssh daint # or the supercomputer you are working on
module load daint-mc # supercomputer-specific module
module load cray-python/3.8.5.0 # load Python environment

# The easiest way to install DataLad with all
# dependencies on a supercomputer without root
# permissions is by using conda
wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
bash Miniconda3-latest-Linux-x86_64.sh
# acknowledge license, initialize Miniconda3, close and
# re-open shell
conda install -c conda-forgedatalad
```

## Install DataLad extensions

Two DataLad extensions are required:

`datalad-container` and

`datalad-neuroimaging`

container  
support

Metadata  
extraction

Install them with pip like this:

```
pip install datalad-neuroimaging datalad-container
```

## Configure DataLad

Setup your git identity. Specify your git username and email, which will be used to track changes in your DataLad data set: changes you make are associated with your name and email address.

```
cd ~  
git config --global --add user.name "MichaelSchirner"  
git config --global --add user.email m.schirner@fu-berlin.de
```

# Create an input dataset



If your data is not yet a DataLad dataset, you can transform it into one with the following commands.

```
# create a dataset in an existing directory
$ datalad create -f .
# save its contents
$ datalad save . -m "Import all data"
```



For large datasets (e.g. HCP, UK Biobank),  
we create subdataset for each subject

# Create a new DataLad dataset

First we create a new empty dataset and then add existing containers to it.

```
$ datalad create pipeline  
$ cd pipeline
```

## Add container(s)

Add a software-container to the dataset using *datalad containers-add* from the datalad-container extension. The *--url* parameter can be a local path to your container image, or a URL to a container hub such as Dockerhub or Singularity Hub.

```
$ datalad containers-add cat --url<path/or/url/to/image> \  
--call-fmt "singularity run -B {{pwd}} --cleanenv {img} {cmd}"
```



Re-usable  
container  
pipeline  
dataset

## Create empty dataset and then clone and add the container as subdataset

Here we create an empty dataset and then clone and register our needed container(s) as a subdataset ([documentation](#)). Arguments like `--bind` that are intended for singularity rather than the underlying command should be specified with `--call-fmt` when calling `containers-add`. It's also fine to edit the cmdexec value in `.datalad/config` after the fact.

```
# Create a source dataset for all analysis components
source_ds="pipeline"
datalad create $source_ds
cd $source_ds

# clone the container-dataset as a subdataset.
containername='tvbpipe-fmriprep'
containerstore="/scratch/sn3000/bp000225/tvbpipe-fmriprep"
datalad clone -d . ${containerstore} code/pipeline

# Register the container in the top-level dataset.
# If necessary, configure your own container call in the
# --call-fmt argument. If your container does not need a
# custom call format, remove the --call-fmt flag and its
# options below.
datalad containers-add \
--call-fmt 'singularity run -B {{pwd}} --cleanenv {img} {cmd}' \
-i /scratch/sn3000/bp000225/${source_ds}/code/pipeline/.datalad/environments/${containername}
$containername

# amend the previous commit with a nicer commit message
git commit --amend -m 'Register pipeline dataset'
```

Analysis  
superdataset to  
specify all  
dependencies of  
an analysis

## Import custom code

Here we can copy custom code into the dataset. For example, FreeSurfer (part of fmriprep, MRtrix3\_connectome) needs a license file (free but must be individually obtained), which is now copied from the home folder into the data set.

```
# committing changes after copying license file
cp ~/license.txt code/license.txt
datalad save -m "Add Freesurfer license file"
```

## Quickly test-run the container on the login node

Before we create SLURM batch scripts we will give the container a test run to make sure we set up everything correctly so far.

```
# load Singularity module again in case there was an
# environment change
module load singularity

# the subject we test
subid="sub-CON03"

# create workdir for fmriprep inside the dataset to
# simplify singularity call PWD will be available in the
# container
mkdir -p .git/tmp/wdir

# Replace the command with an fmriprep parametrization
# that fits your analysis
datalad containers-run \
  -m "Compute ${subid}" \
  -n tvbpipe-fmriprep \
  --explicit \
  -o fmriprep/${subid} \
  -i inputs/data/${subid}/ses-postop/ \
  -i code/license.txt \
  "inputs/data/${subid} . participant --participant-label $subid \
   --anat-only -w .git/tmp/wdir --fs-no-reconall --skip-bids-validation \
   --fs-license-file code/license.txt"
```

# Limitations

- not every feature of DataLad may work „out of the box“
- sometimes customization is required

## ↗ Find-out-more: Fine-tuning: Enable re-running

If you want to make sure that your dataset is set up in a way that you have the ability to rerun a computation quickly, the following fMRIprep-specific consideration is important: If fMRIprep finds preexisting results, it will fail to run. Therefore, all outputs of a job need to be removed before the job is started<sup>[3]</sup>. We can simply add an attempt to do this in the script (it wouldn't do any harm if there is nothing to be removed):

```
(cd fmriprep && rm -rf logs "$subid" "$subid.html" dataset_description.json  
(cd freesurfer && rm -rf fsaverage "$subid")
```

```
# pybids (inside fmriprep) gets angry when it sees dangling symlinks  
# of .json files -- wipe them out, spare only those that belong to  
# the participant we want to process in this job  
find sourcedata -mindepth 2 -name '*.json' -a ! -wholename "$1"/*
```

meanwhile fixed in  
PyBIDS thanks to  
community  
interaction  
<https://github.com/bids-standard/pybids/issues/631>

# Summary / take home messages

- The presented workflow is quite complex
- Nevertheless, it boils down to a handful of DataLad commands
- It's not difficult to re-combine the commands for different use cases
- Containers capture all software and environment
- DataLad captures how containers are executed including all **varying parts** (e.g. ad-hoc configuration, live user input, etc.)
- Enables **everyone** (e.g. reviewers, trainees) to **re-run** individual steps of the workflow on a **fine-grained level** on **another machine** with a **simple command**

THANK

YOU

# Questions to the audience

Q1: Was the presentation understandable?  
What was not clear?

Q2: Is this type of large-scale processing  
relevant for your project?

Q3: Who in the audience considers using  
such a technology?

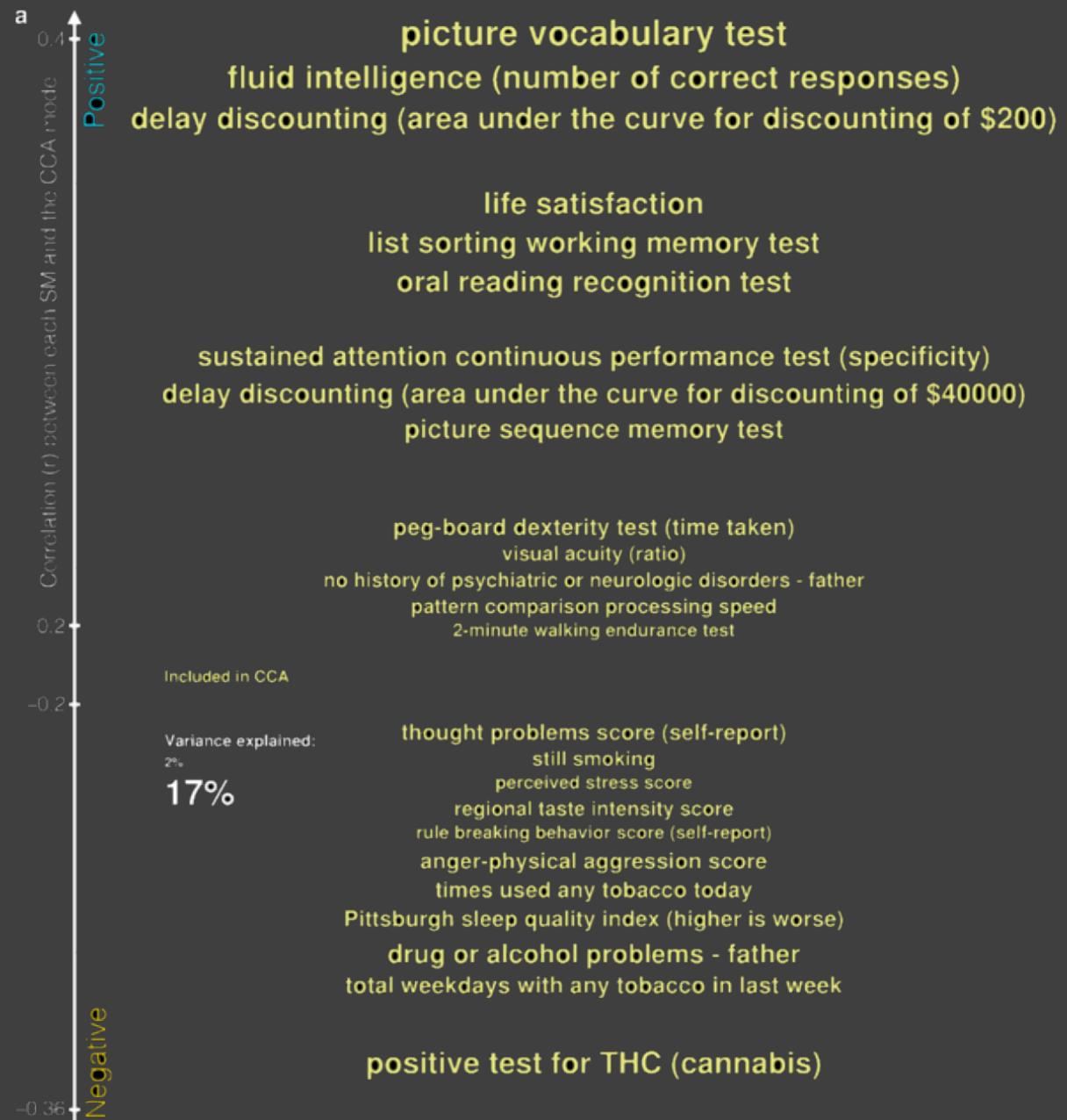
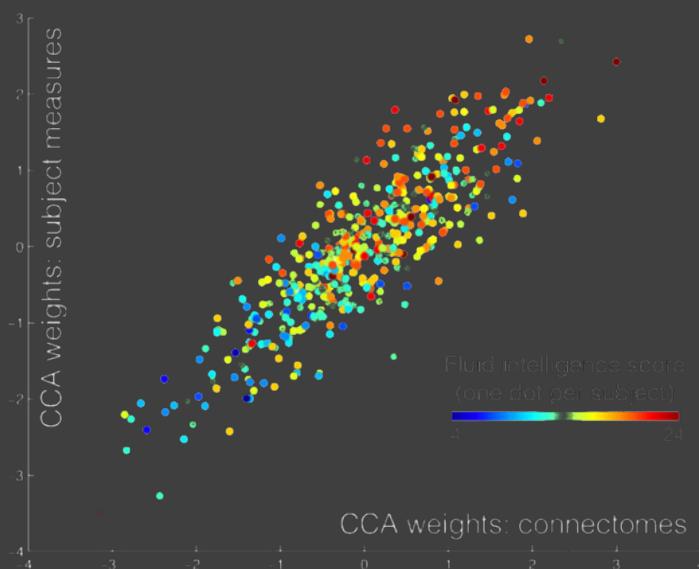
Q4: Provenance tracking: on a scale from  
hand-written lab book to DataLad – where  
do you locate yourself?

# fMRI contains highly personal information

A positive-negative mode of population covariation links brain connectivity, demographics and behavior

Stephen M Smith<sup>1</sup>, Thomas E Nichols<sup>2</sup>, Diego Vidaurre<sup>3</sup>, Anderson M Winkler<sup>1</sup>, Timothy E J Behrens<sup>1</sup>, Matthew F Glasser<sup>4</sup>, Kamil Ugurbil<sup>5</sup>, Deanna M Barch<sup>4</sup>, David C Van Essen<sup>4</sup> & Karla L Miller<sup>1</sup>

We investigated the relationship between individual subjects' functional connectomes and 280 behavioral and demographic measures in a single holistic multivariate analysis relating imaging to non-imaging data from 461 subjects in the Human Connectome Project. We identified one strong mode of population co-variation: subjects were predominantly spread along a single 'positive-negative' axis linking lifestyle, demographic and psychometric measures to each other and to a specific pattern of brain connectivity.



§

Art. 24 GDPR

## Responsibility of the controller

*“...shall implement appropriate technical and organisational measures to ensure and...demonstrate that processing is performed in accordance with this Regulation”*



**Controller determines alone or jointly with others the purposes and means of processing**

**Responsible for protecting the personal data!**

# Data privacy risks



*Intercept during transit*



*Break-in on device*



*User impersonation*

**Problem 1:** health data shall be secured by default

*Solution: health data shall be always encrypted except during processing*

**Problem 2:** how to secure unencrypted data during processing?

*Solution: Sandboxes, temporary mount file systems, authenticated data containers*



# Exposed (anonymized)

Tracked

## VERSION HISTORY

## FILE NAMES

patient\_record/

photo.jpg

insurance.json

visit/

1.json

2.json

imaging/

brain/20031.img

thorax/20201.img

## METADATA

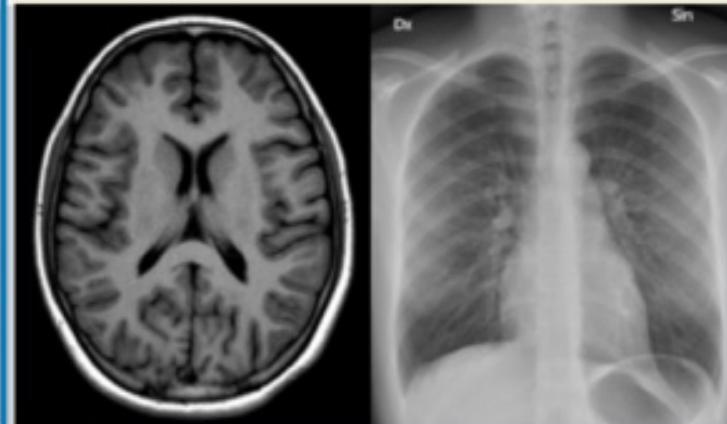
```
{"type": "patient", "id": "122bf-342-422", ANON IDS  
"imaging": ["brain", "thorax"], "sex": "male",  
"diagnosis": ["F51", "G12", "G52"], "age": 52...}
```



PRISTINE INFORMATION

```
{"account": "293924",  
"id": "a399b30033d..."}
```

```
{"date": "2020-04-21",  
"diagnosis": "...",  
"notes": "..." ...}
```



Authorized/local access only











# Data

Analysis	Data	
	Same	Different
Same	Reproducible	Replicable
Different	Robust	Generalisable

# Open Scholarship: what science can learn from (open source) software development

- transparency & accessibility: everybody can engage
- version control: unique persistent identifiers for digital objects
- containerization (dependency management): code, data and environment in one robust package
- attribution: “commits” for fine-grained tracking of contributions
- issue tracking: peer-review by everyone and all the time
- unit testing: provenance; from raw data to figure
- continuous integration: consistent deployment of high-quality artefacts at all stages of development





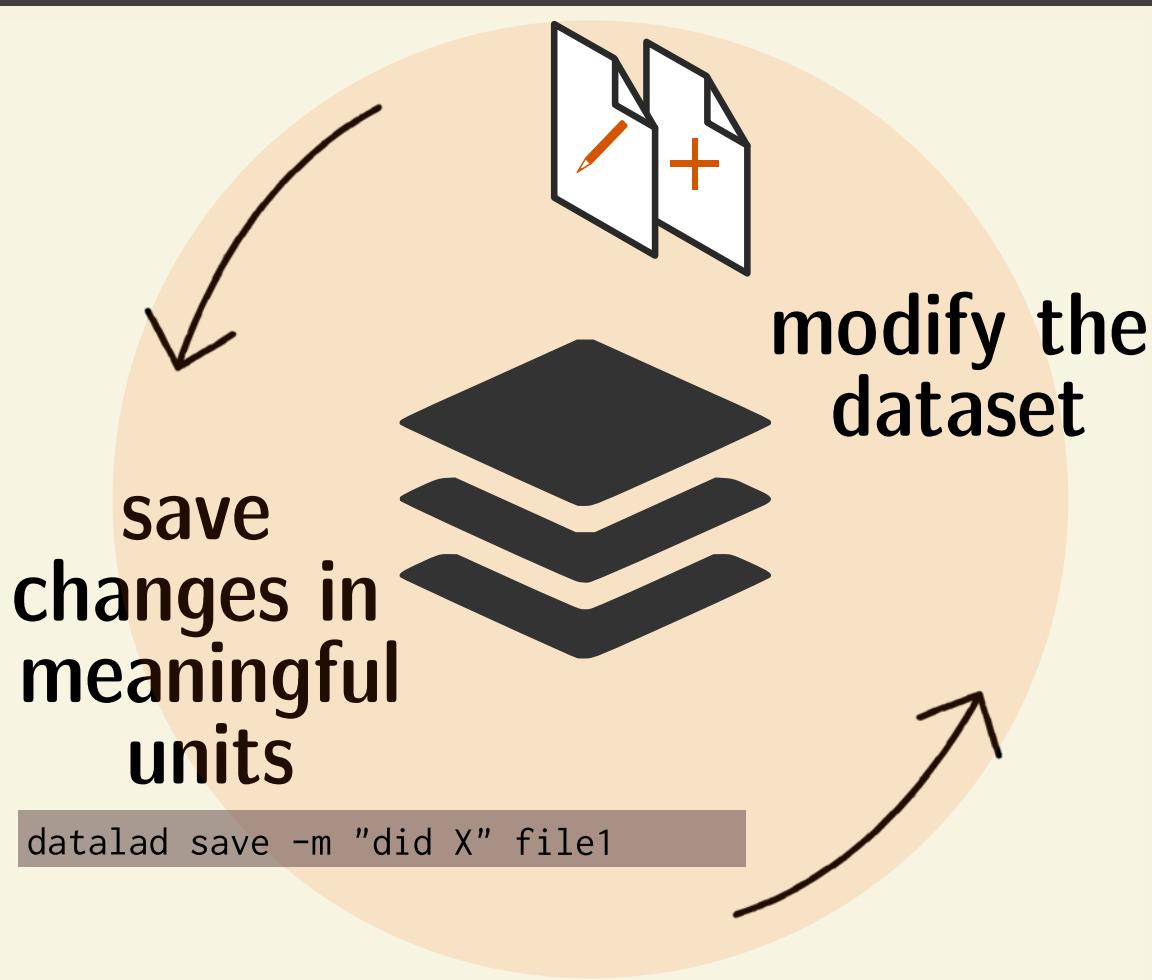
- “Git for data”
- based on **Git** and **git-annex**
- command-line tool for Linux, OS X, Windows
- version control arbitrarily large files
- provenance tracking
- share and obtain data
- reproducible workflows
- domain-agnostic



- DataLad datasets
  - are directories managed by DataLad
  - are Git repositories
  - can be nested: linked subdirectories
  - can be created or installed
  - can branch and merge



## Local version control



`datalad create`

Create empty dataset

`datalad save`

record state of dataset to history incl. commit message

`datalad download-url`

obtain web content and record origin

`datalad status`

report tracking state