RDM IN SYSTEMS NEUROSCIENCE:

CHALLENGES AND BEST PRACTICES





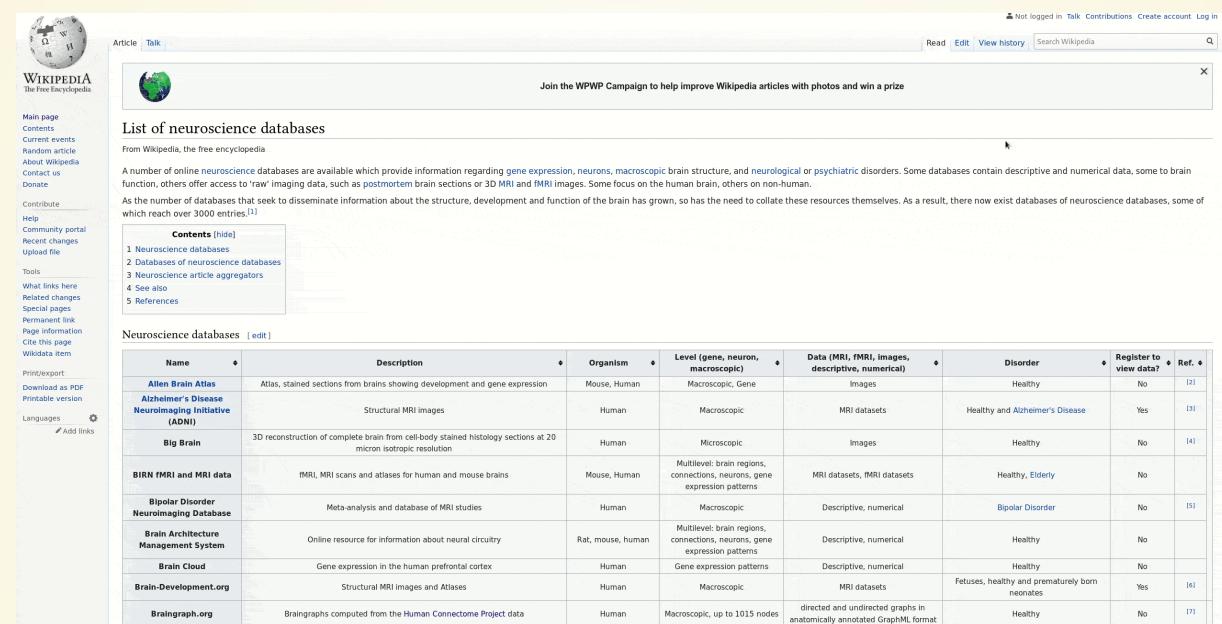
Institute of Neuroscience and Medicine (INM-7), Research Center Jülich



Institute for Experimental Psychology, HHU Düsseldorf

Slides: files.inm7.de/adina/talks/html/fdm_nrw_hhu.html DOI; 10.5281/zenodo.10122803

A growing culture of open data



A large and growing amount of open source software











































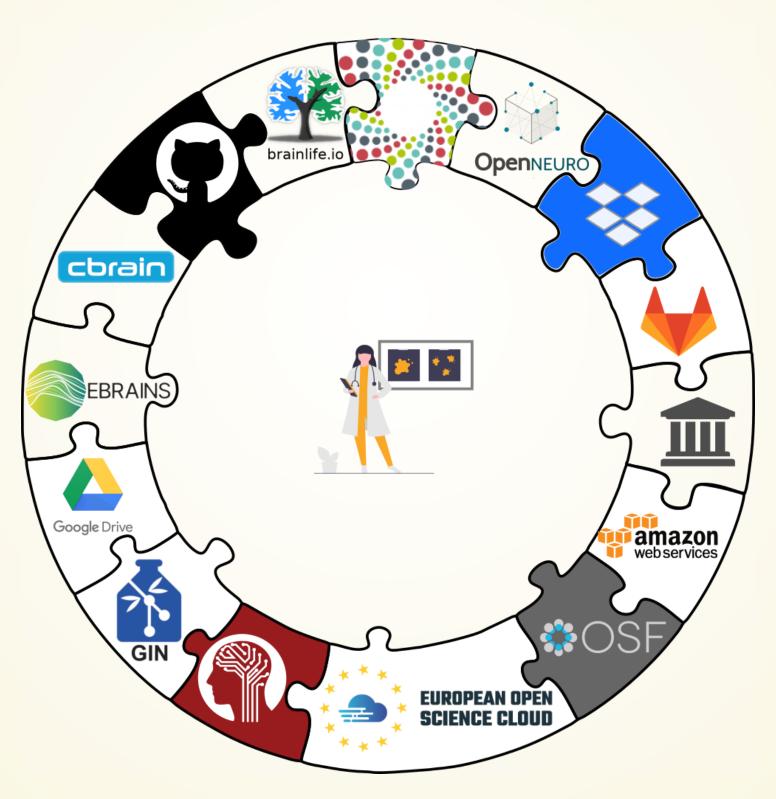




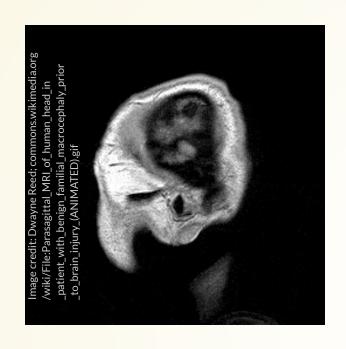


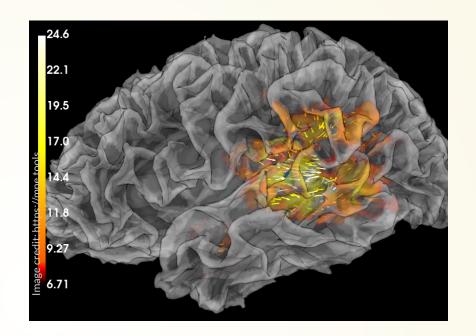
... and many more!

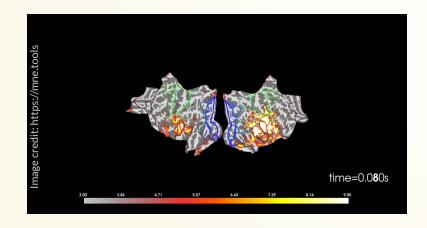
Many readily available, often free, sometimes FOSS, services for data storage and collaboration

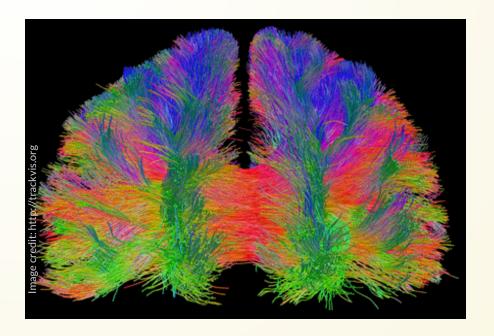


Work on fascinating questions with fascinating data









THE GOOD NEWS, THE BAD NEWS

Interesting

... that's subject to GDPR, making storage, analysis, and sharing

data

difficult

Data analysis

Reproducibility is threatened by intransparent, multi-stepped

analyses¹ & unstable results across software versions². Dataset

sizes exceed computational capabilities (e.g., HCP project:

~100TB, 1.2k people).

Data sharing

Heterogenous distribution and updating, many scientists lack

data management skills

Collaboration

Few interoperable workflows across institutes, rather: isolated

solutions

¹ Botvinik-Nezer et al., 2020: Variability in the analysis of a single neuroimaging dataset by many teams

² Kiar et al., 2020: Comparing perturbation models for evaluating stability of neuroimaging pipelines



- Domain-agnostic data management tool (command-line + graphical user interface), built on top of Git & Git-annex
- 10+ year open source project (100+ contributors), available for all major OS
- Born from rethinking data:
 - Just like code, data is not static.
 - Just like code, data is subject to collaboration. Stream-lined workflows for sharing and collaborating should be possible, mirroring those in software development.
 - Provenance of data is essential for reproducible, trustworthy, and FAIR science
 - Flexibility and interoperability with existing tools is the key to sustainability and ease of use



- Domain-agnostic command-line tool (+ graphical user interface), built on top of Git & Git-annex
- 10+ year open source project (100+ contributors), available for all major OS
- Major features:

Version-controlling arbitrarily large content

Version control data & software alongside to code!

Transport mechanisms for sharing, updating & obtaining data

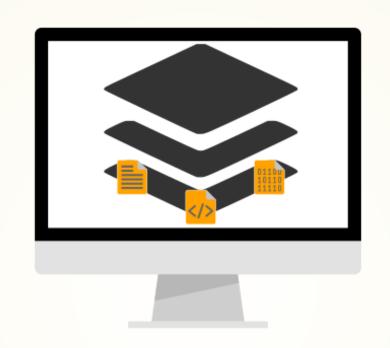
Consume & collaborate on data (analyses) like software

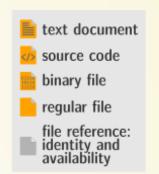
(Computationally) reproducible data analysis

Track and share provenance of all digital objects

(... and much more)

EXHAUSTIVE TRACKING OF RESEARCH COMPONENTS

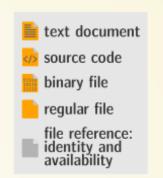




Well-structured datasets (using community standards), and portable computational environments — and their evolution — are the precondition for reproducibility

CAPTURE COMPUTATIONAL PROVENANCE





Which data was needed at which version, as input into which code, running with what parameterization in which computional environment, to generate an outcome?

```
# execute any command and capture its output
# while recording all input versions too

% datalad run --input ... --output ... <command>
```

EXHAUSTIVE CAPTURE ENABLES PORTABILITY



Precise identification of data and computational environments combined with provenance records form a comprehensive and portable data structure, capturing all aspects of an investigation.

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

text document source code

binary file

regular file

file reference: identity and availability

REPRODUCIBILITY STRENGTHENS TRUST



text document
source code
binary file
regular file
file reference:
identity and
availability

Outcomes of computational transformations can be validated by authorized 3rd-parties. 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>
```

ULTIMATE GOAL: (RE-)USABILITY



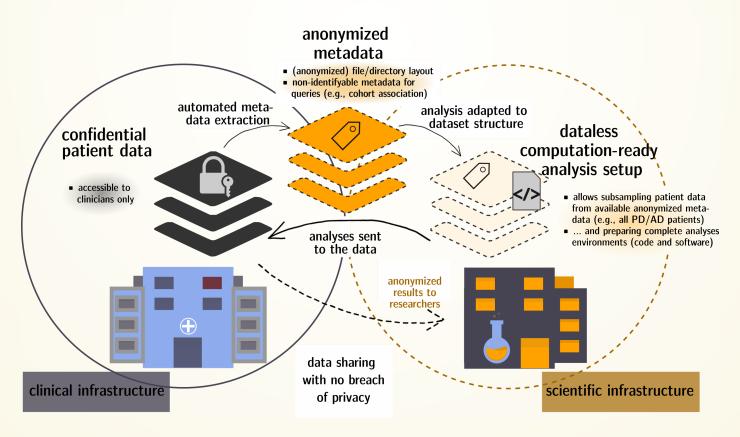
Verifiable, portable, self-contained data structures that track all aspects of an investigation exhaustively can be (re-)used as modular components in larger contexts — propagating their traits

```
# 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>
```

TRACKING & SHARING - WHERE IS THE PRIVACY?

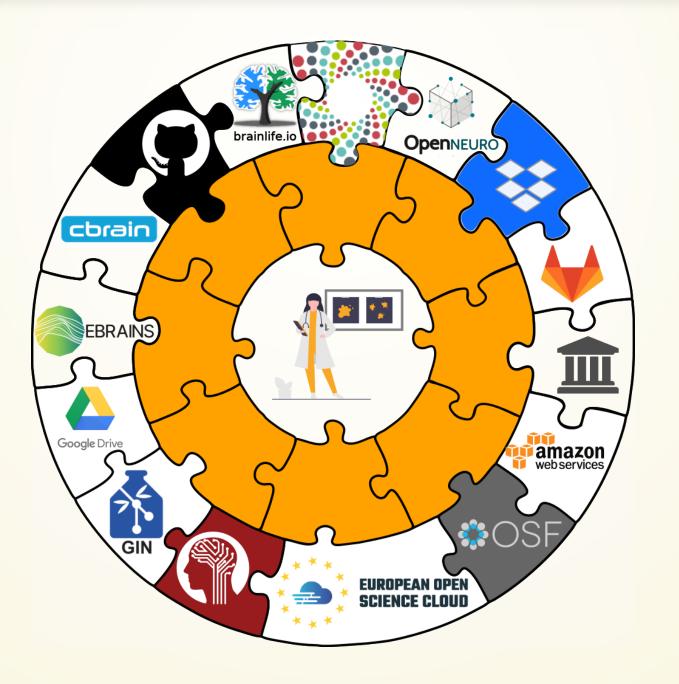
- Datasets have an optional annex for (large or sensitive) data.
- Rather than file content, identity (hash) and location information is tracked.
 Users have fine-grained control over transport and access; encryption is possible.
- Datasets can separate data access from meta data access



FINE-GRAINED FILE TRANSPORT

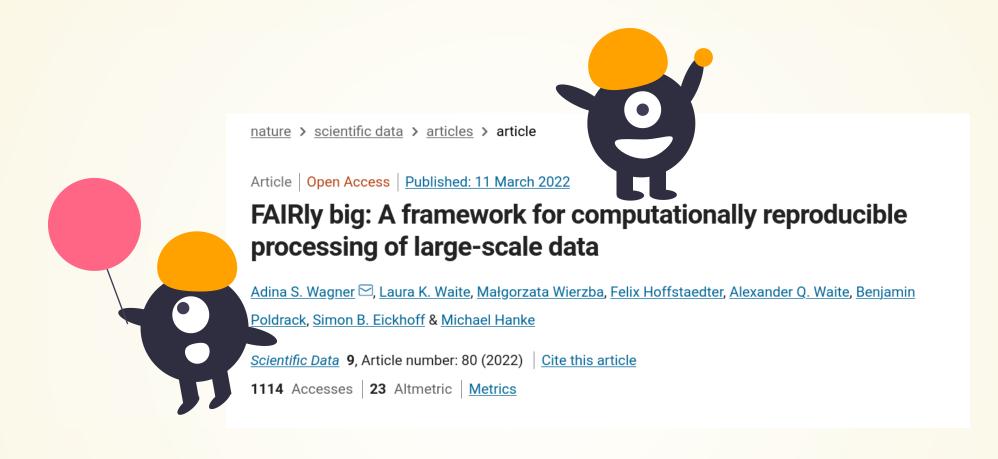
Cloned datasets are lean, because file content is retrieved on demand:

```
1 $ datalad clone \
2    https://github.com/datalad-datasets/human-connectome-project-openaccess.git
3 $ cd human-connectome-project-openaccess && du -sh
4    5.0M .
5 $ datalad get HCP1200/102513/T1w/T1w_acpc_dc.nii.gz
6    get(ok): HCP1200/102513/T1w/T1w_acpc_dc.nii.gz (file) [from datalad...]
7 $ datalad drop HCP1200/102513/T1w/T1w_acpc_dc.nii.gz
8    drop(ok)
```



BIG DATA

FAIRLY BIG: SCALING UP



Objective: Process the UK Biobank (imaging data) with 76 TB in 43 million files from 42,715 participants

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

WHERE CAN I FIND OUT MORE?

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 (github.com/datalad/datalad)

Reach out to the user community with

A question on neurostars.org with a datalad tag

Find more user tutorials or workshop recordings

- On DataLad's YouTube channel (www.youtube.com/channel/datalad)
- In the DataLad Handbook (handbook.datalad.org)
- In the DataLad RDM course (psychoinformatics-de.github.io/rdm-course)
- In the Official API documentation (docs.datalad.org)

{OPEN,TRANSPARENT,REPRODUCIBLE} SCIENCE

- Treat data like software: obtain, version, share, and update data
- Simplified data management, disk-space aware storage & computing
- Transparent and reproducible science: link code, data, software, and execution in a human- and machine-readable way
- Collaborate: Generic workflows, interoperabality with established tools & services



ACKNOWLEDGEMENTS

Software

- Joey Hess (git-annex)
- The DataLad team & contributors

THANKS!

Questions?



Funders





BMBF 01GQ1905









Collaborators











