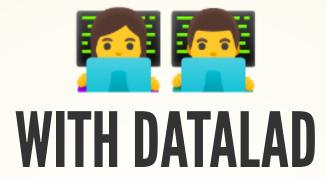
RESEARCH DATA MANAGEMENT



y @AdinaKrik

Adina Wagner Michael Hanke



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

Slides: https://github.com/datalad-handbook/course/

WELCOME!

Approximate workshop schedule

Session 1 (now, 13.30-15.00)

Logistics & Intro 2,

Hands-on Terminal Basics ,

Demo of core functionality 🤵 💻

Session 2 (today, 16.00-18.00)

Hands-on DataLad Basics & Exercises

Session 3 (tomorrow, 11.00-12.30)

Sharing and Collaboration 🧖,

Hands-on Data publication

Session 4 (tomorrow, 13.30-15.00)

Computational reproducibility 🤵 💻,

Outro 👰,

Final QA ?

LOGISTICS AND LINKS

- You can download these slides at https://doi.org/10.5281/zenodo.6827086 (scan the QR code), and you can find their sources at github.com/datalad-handbook/datalad-course
- Some of today's code-along workshop contents are at psychoinformatics-de.github.io/rdm-course
- The workshop will be interactive. If you do not have the software installed on your own system, you can access a Jupyterhub from your browser at datalad-hub.inm7.de
- You can log in to the Juypterhub with a pre-set username (take one out of the jar) and a self-set password. Remember the password for tomorrow!
- A requirements.txt file on Zenodo details the software environment we setup on the Jupyterhub



INTERACTIVITY

- The workshop centers around DataLad (version 0.16 and up) for real-world research data management use cases
- There are no stupid questions; ask anything any time
- Something doesn't look right on your system? Stick a post-it to your screen. We'll take a look together
- We're available outside of sessions, too. Chat about your use cases or questions over a coffee or meal
 - 4 sessions = time for more than a standard introduction.
 - Materials are available online & persistent, we can be flexible & spontaneous if specific topics interest you

AFTER THE WORKSHOP

If you have a question after the workshop, you can reach out for help:

Reach out to to the DataLad team via

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

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
- In the DataLad Handbook
- In the DataLad RDM course
- In the Official API documentation

AUDIENCE RESPONSE SYSTEM

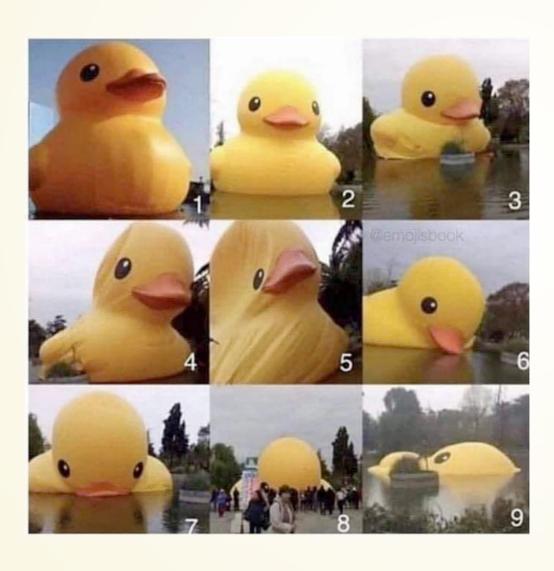
Use your phone to scan the QR code, or open the link in a new browser window

What is a commit hash? http://etc.ch/7YEk

A trending Twitter topic	0
40-character string identifying changes	0
Drugs consumed at a wedding ceremony	0



ON A SCALE OF RUBBER DUCKS...



What is a commit hash? http://etc.ch/7YEk

A trending Twitter topic 0

40-character string identifying changes

Drugs consumed at a wedding ceremony

0 votes - 0 participants

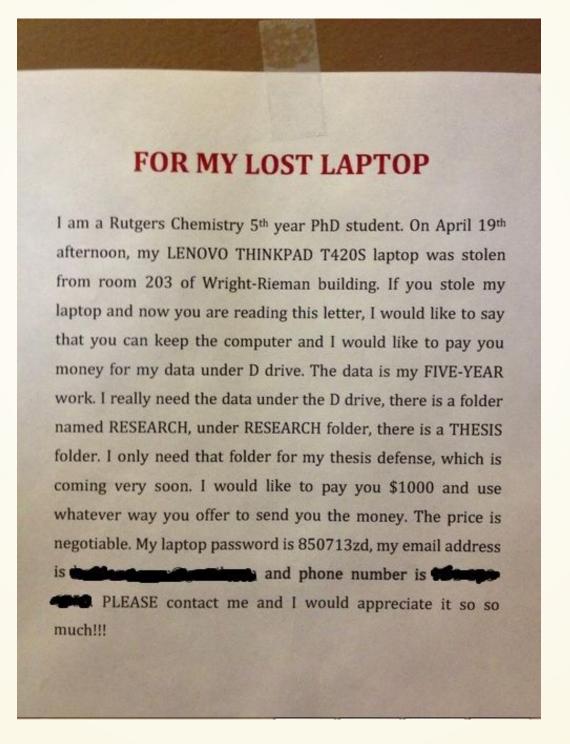


RESEARCH DATA MANAGEMENT

You write a paper & stay up late to generate good-looking figures, but you have to tweak many parameters and display options. The next morning, you have no idea which parameters produced which figures, and which of the figures fit to what you report in the paper.



Your research project produces phenomenal results, but your laptop, the only place that stores the source code for the results, is stolen or breaks



A graduate student complains that a research idea does not work. Their supervisor can't figure out what the student did and how, and the student can't sufficiently explain their approach (data, algorithms, software). Weeks of discussion and mis-communication ensues because the supervisor can't first-hand explore or use the students project.









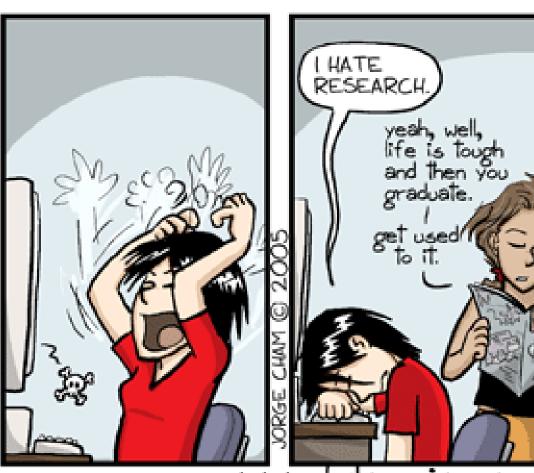
WWW.PHDCOMICS.COM

You wrote a script during your PhD that applied a specific method to a dataset.

Now, with new data and a new project, you try to reuse the script, but forgot how it worked.







www.phdcomics.com

You try to recreate results from another lab's published paper. You base your reimplementation on everything reported in their paper, but the results you obtain look nowhere like the original.



Scriberia

All these problems were paraphrased from Buckheit & Donoho, 1995

Let's do better!



DataLad can help with small or large-scale data management

Free, open source, command line tool & Python API



- A command-line tool, available for all major operating systems (Linux, macOS/OSX, Windows), MIT-licensed
- Build on top of Git and Git-annex
- Allows...
 - ... version-controlling arbitrarily large content version control data and software alongside to code!
 - ... transport mechanisms for sharing and obtaining data consume and collaborate on data (analyses) like software
 - ... (computationally) reproducible data analysis

 Track and share provenance of all digital objects
 - ... and much more
- Completely domain-agnostic

ACKNOWLEDGEMENTS

Software

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

Illustrations

The Turing Way project & Scriberia















Collaborators





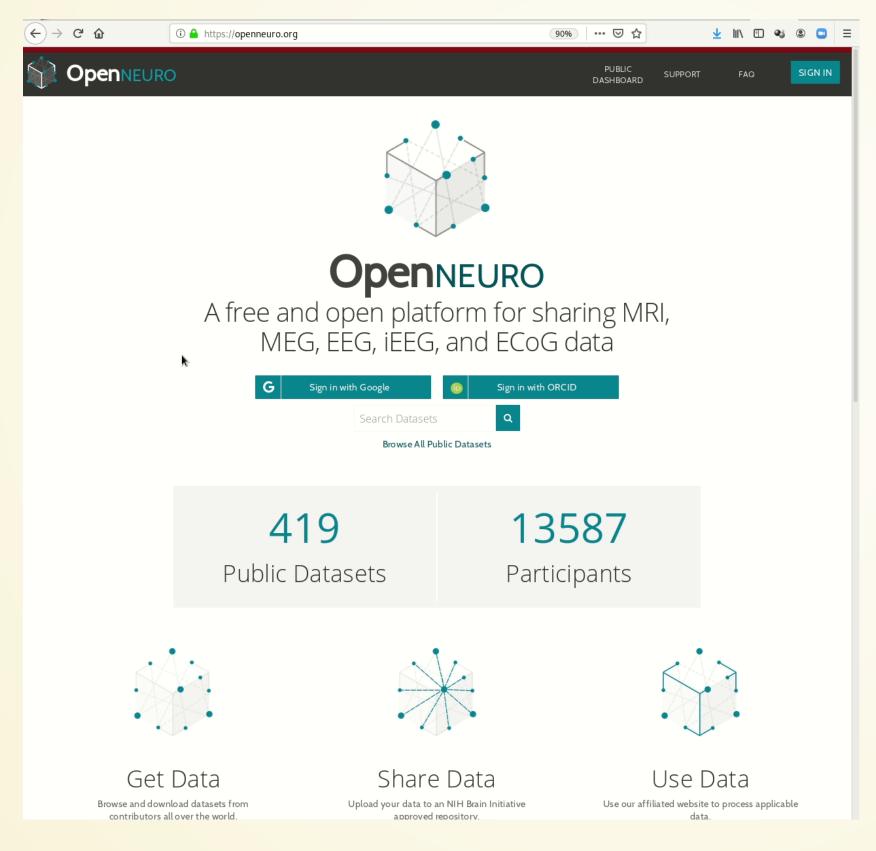




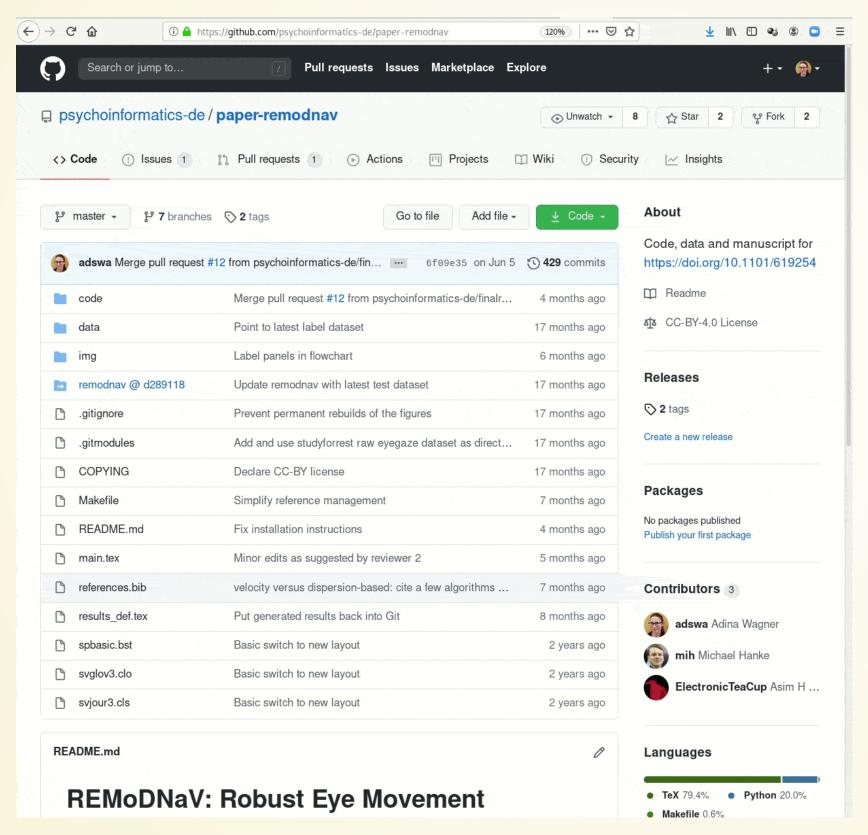




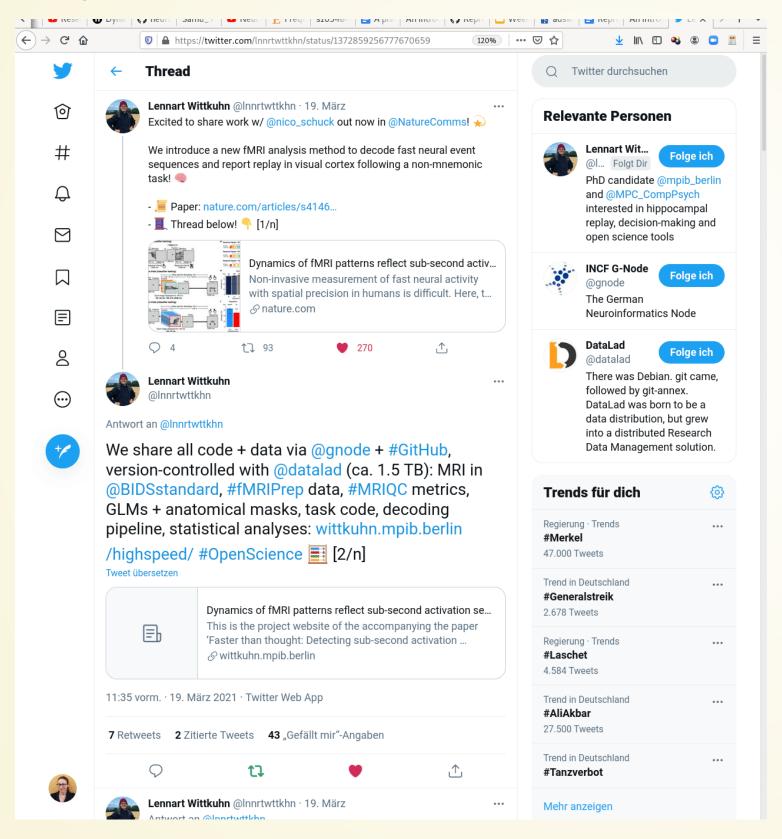
 Behind-the-scenes infrastructure component for data transport and versioning (e.g., used by OpenNeuro, brainlife.io, the Canadian Open Neuroscience Platform (CONP), CBRAIN)



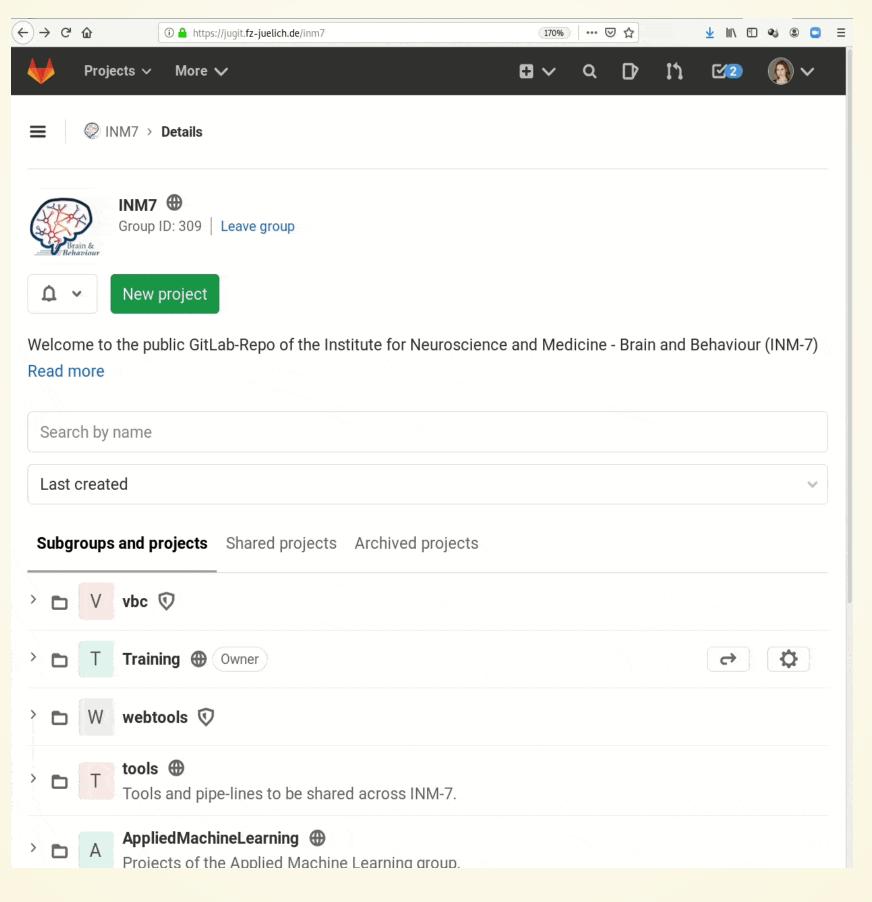
 Creating and sharing reproducible, open science: Sharing data, software, code, and provenance



 Creating and sharing reproducible, open science: Sharing data, software, code, and provenance



Central data management and archival system



Scalable computing framework for reproducible science

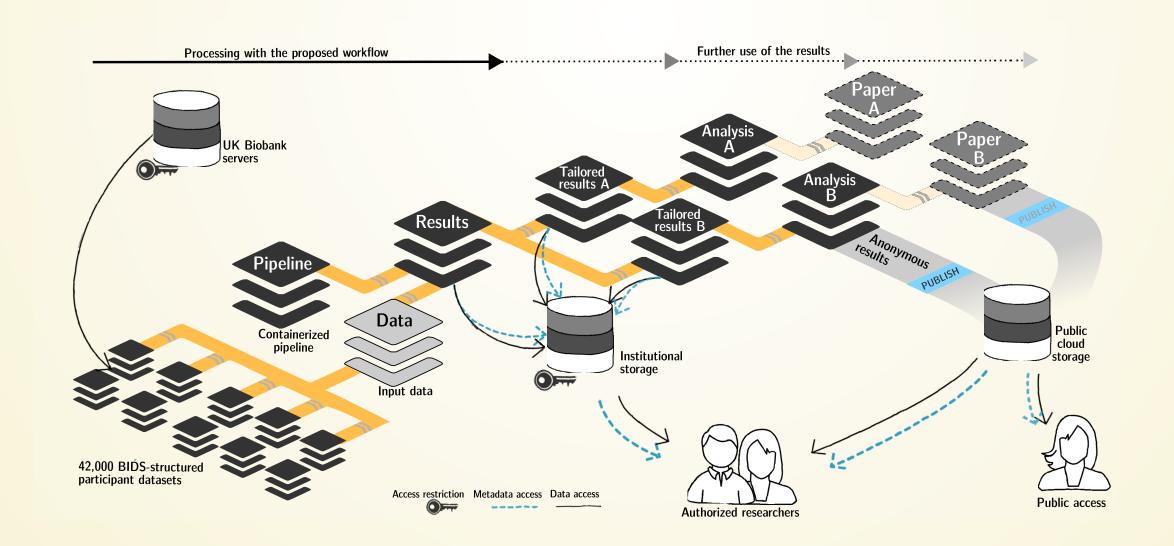
Article Open Access | Published: 11 March 2022

FAIRly big: A framework for computationally reproducible processing of large-scale data

Adina S. Wagner ☑, Laura K. Waite, Małgorzata Wierzba, Felix Hoffstaedter, Alexander Q. Waite, Benjamin Poldrack, Simon B. Eickhoff & Michael Hanke

Scientific Data 9, Article number: 80 (2022) Cite this article

813 Accesses | 23 Altmetric | Metrics



PREREQUISITES: TERMINAL

DataLad can be used from the command line

datalad create mydataset

... or with its Python API

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

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

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

PREREQUISITES: TERMINAL

What is a commit hash?

http://etc.ch/7YEk

A trending Twitter topic	0
40-character string identifying changes	0
Drugs consumed at a wedding ceremony	0

Direct Poll

0 votes - 0 participants

datalad-hub.inm7.de

Unix terminal cheatsheet (incl. Windows equivalents)

PREREQUISITES: INSTALLATION AND CONFIGURATION

Your installed version of DataLad should be at least 0.16

```
datalad --version 0.17.2
```

 DataLad relies on Git to create a revision history with detailed information on what was changes, when, and how. Therefore, you should tell Git who you are and configure a Git identity (name and email). Find out if an identity is set by running either of:

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

\$ datalad configuration get user.name user.email
Adina Wagner
adina.wagner@t-online.de

Set a Git identity using either of

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

\$ datalad configuration --scope global \
set user.name="Adina Wagner"
\$ datalad configuration --scope global \
set user.email="adina.wagner@t-online.de"

Allow brand-new DataLad functionality:

```
datalad configuration --scope global set datalad.extensions.load=next
```

Find installation and configuration instructions at handbook.datalad.org

PREREQUISITES: USING DATALAD

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

```
Each datalad
datalad [--GLOBAL-OPTION <opt. flag spec.>] COMMAND [ARGUMENTS] [--OPTION <opt. flag spec>]
                                                                                                                                       invocation can
                                                                                                                                       have two sets of
                                                          OPTIONS
                                                              -d/--dataset
                                                                                Dataset location: path to root, or ^ for superdataset
-c KEY=VALUE
                                                                                                                                       options: general
  Set config variables (overrides configurations in files)
                                                              -D/--description A location description (e.g., "my backup server")
                                                                                                                                       options are given
-f/--output-format default|json|json_pp|tailored
                                                              -f/--force
                                                                                Force execution of a command (Dangerzone!)
                                                                                                                                       first, command-
  Specify the format for command result renderung
                                                                                A description about a change made to the dataset
                                                                                                                                       specific ones go
-l/--log-level critical|error|warning|info|debug
                                                              -r/--recursive
                                                                                Perform an operation recursively across subdatasets
                                                                                                                                       after the
  Set logging verbosity level
                                                              -R/--recursion-limit <n> Limit recursion to n subdataset levels
                                                                                                                                       subcommand.
```

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

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

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

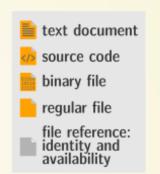
If everything is important...

...track everything!

http://datalad.org

EXHAUSTIVE TRACKING OF RESEARCH COMPONENTS



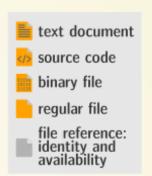


Well-structured datasets (using community standards), and portable computational environments — and their evolution — are the precondition for reproducibility

```
# turn any directory into a dataset
# with version control
# file content of any size
% datalad create <directory>
% datalad save
```

CAPTURE COMPUTATIONAL PROVENANCE



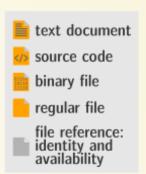


Which data was needed at which version, as input into which code, running with what parameterization in which computional environment, to generate an outcome?

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

EXHAUSTIVE CAPTURE ENABLES PORTABILITY



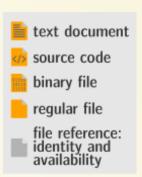


Precise identification of data and computational environments, combined for provenance records form a comprehensive and portable data structure, capturing all aspects of an investigation.

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

REPRODUCIBILITY STRENGTHENS TRUST





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>
```

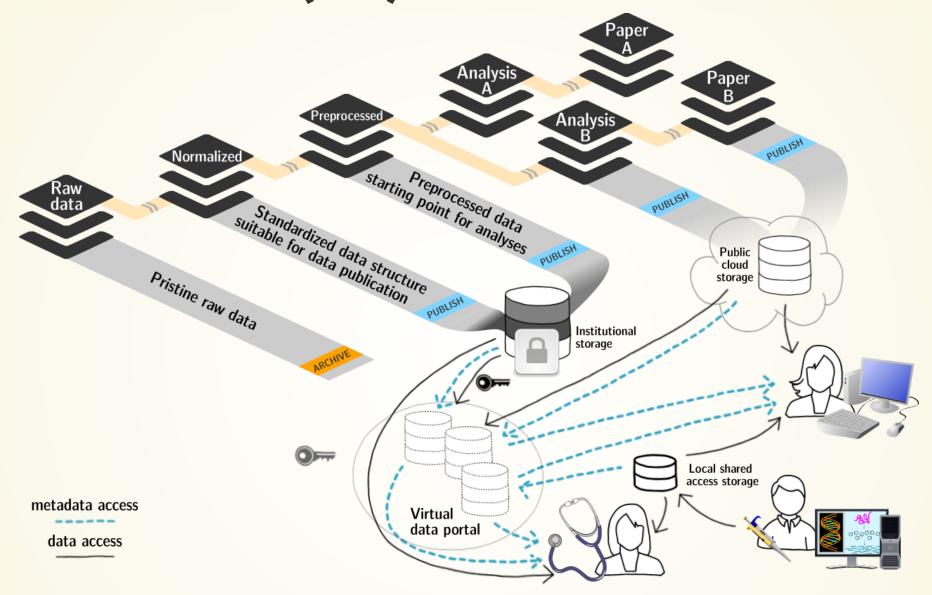
ULTIMATE GOAL: (RE-)USABILITY



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>
```

DATALAD: MANAGE (CO-)EVOLUTION OF DIGITAL OBJECTS



Consume, create, curate, analyze, publish, and query data with full provenance capture and "universal" metadata support.

DataLad is free and open source (MIT-licensed). http://datalad.org

Halchenko, Meyer, Poldrack, ... & Hanke, M. (2021). DataLad: distributed system for joint management of code, data, and their relationship. Journal of Open Source Software, 6(63), 3262.

LET'S TRY...