

Reproducible Research Data Management with DataLad

6th RDM-Workshop 2024 on Research Data Management in the Max Planck Society

 Slides |  Source

License [CC BY SA 4.0](#) DOI [10.5281/zenodo.10844705](https://doi.org/10.5281/zenodo.10844705)

Dr. Lennart Wittkuhn 

wittkuhn@mpib-berlin.mpg.de

Max Planck Research Group NeuroCode, Max Planck Institute for Human Development, Berlin, Germany

Max Planck UCL Centre for Computational Psychiatry and Ageing Research, Berlin, Germany

Institute of Psychology, University of Hamburg, Germany

March 20, 2024

About



Dr. Lennart Wittkuhn

✉ wittkuhn@mpib-berlin.mpg.de

🏠 <https://lennartwittkuhn.com/>

🐙 Mastodon 🗨️ GitHub 🌐 LinkedIn

About me

👤 I am a **Postdoctoral Research Data Scientist** in Cognitive Neuroscience at the [Institute of Psychology](#) at the [University of Hamburg](#) (PI: Nicolas Schuck)

🎓 **BSc Psychology & MSc Cognitive Neuroscience** (TU Dresden), **PhD Cognitive Neuroscience** (Max Planck Institute for Human Development)

🧠 I study **the role of fast neural memory reactivation** in the human brain, applying **machine learning** and **computational modeling** to **fMRI** data

</> I am passionate about **computational reproducibility**, **research data management**, **open science** and tools that improve the scientific workflow

📄 Find out more about my work on [my website](#), [Google Scholar](#) and [ORCID](#)

About this presentation

🖥️ Slides: <https://lennartwittkuhn.com/talk-mps-fdm-2024>

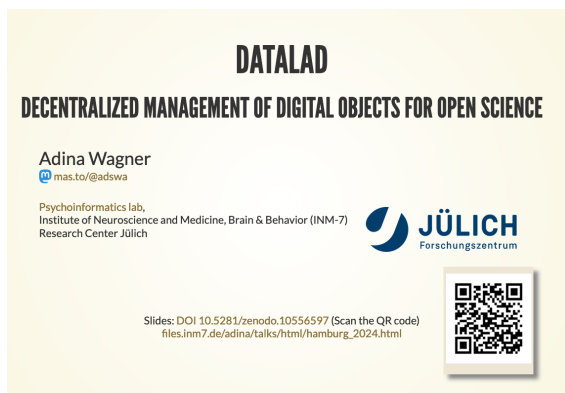
🗄️ Source: <https://github.com/lennartwittkuhn/talk-mps-fdm-2024>

📄 Software: Reproducible slides built with [Quarto](#) and deployed to [GitHub Pages](#) using [GitHub Actions](#) for continuous integration & deployment

📄 License: Creative Commons Attribution-ShareAlike 4.0 ([CC BY-SA 4.0](#))

🗨️ Contact: Feedback or suggestions via [email](#) or [GitHub issues](#). Thank you!

Acknowledgements and further reading



Slides and presentations by Dr. Adina Wagner and the DataLad team, e.g., “DataLad - Decentralized Management of Digital Objects for Open Science” ([Wagner 2024](#))



The **DataLad Handbook** by Wagner et al. ([2022](#)) is a comprehensive educational resource for data management with DataLad.

Slides by Wagner ([2024](#)) (CC BY-SA 4.0)

DataLad Handbook
(CC BY-SA 4.0)

Papers

- Wilson et al. (2014). [Best practices for scientific computing](#). *PLOS Biology*.
- Wilson et al. (2017). [Good enough practices in scientific computing](#). *PLOS Computational Biology*.
- Lowndes et al. (2017). [Our path to better science in less time using open data science tools](#). *Nature Ecology Evolution*.

Talks

- Richard McElreath (2020). [Science as amateur software development](#). YouTube
- Russ Poldrack (2020). [Toward a Culture of Computational Reproducibility](#). YouTube

... and many more!

Agenda

1. Scientific workflows with DataLad

1.1 Version Control

1.2 Modularity & Linking

1.3 Provenance

1.4 Collaboration & Interoperability

2. Integrating DataLad with MPS infrastructure (and beyond)

2.1 Keeper

2.2 Edmond

2.3 ownCloud / nextCloud

2.4 GIN

2.5 GitLab

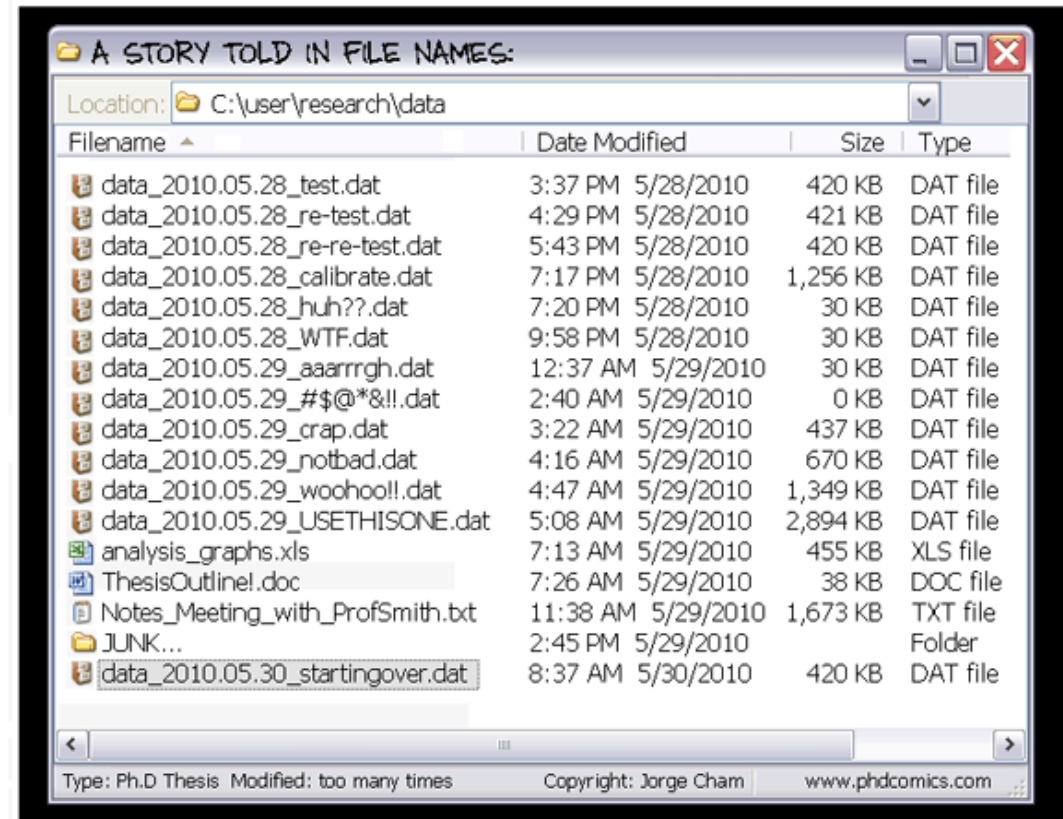
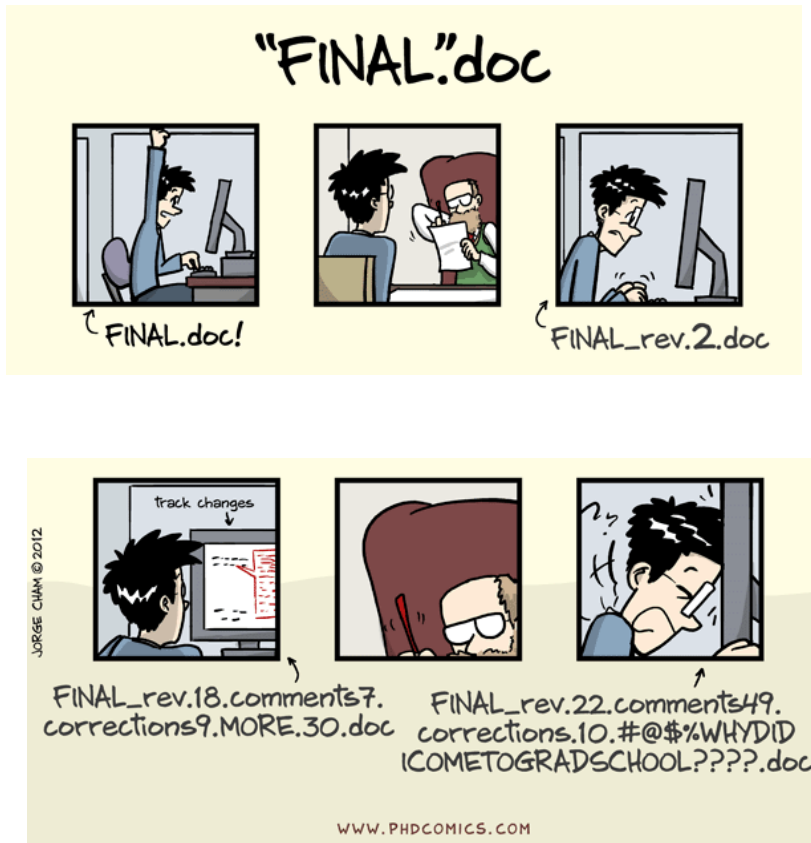
3. Summary & Discussion

Scientific building blocks are not static
We need version control

Why we need version control

... for code (text files)

... for data (binary files)



© Jorge Cham (phdcomics.com)

If everything is relevant, track everything.

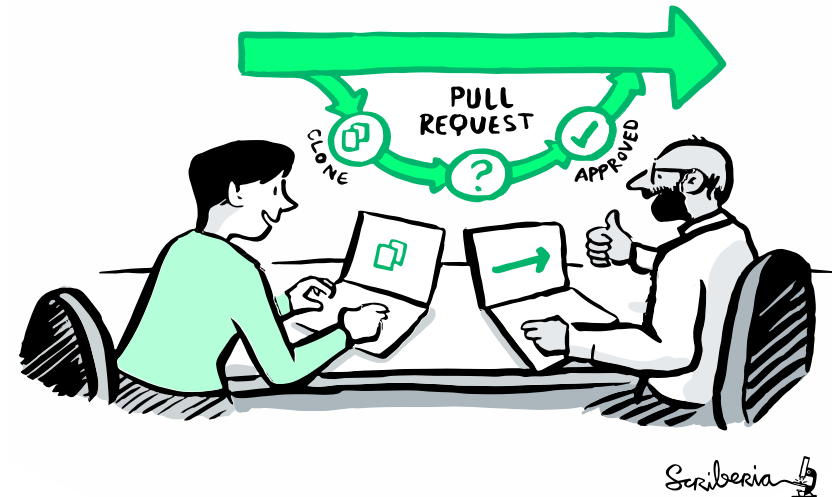
Reproducible Research Data Management with DataLad

What is version control?

“Version control is a systematic approach to record changes made in a [...] set of files, over time. This allows you and your collaborators to track the history, see what changed, and recall specific versions later [...]” (Turing Way)

- 📁 keep track of changes in a directory (a “repository”)
- 🔑 take snapshots (“commits”) of your repo at any time
- 👥 know the history: what was changed when by whom
- 🔄 compare commits and go back to any previous state
- 🔗 work on parallel “branches” & flexibly “merge” them

- 📤 “push” your repo to a “remote” location & share it
- 🔄 share repos on platforms like GitHub or GitLab
- 👥 work together on the same files at the same time
- ✍️ others can read, copy, edit and suggest changes
- 📄 make your repo public and openly share your work



by The Turing Way Community and Scriberia (2024) (CC BY 4.0)

by The Turing Way Community and Scriberia (2024) (CC BY 4.0)

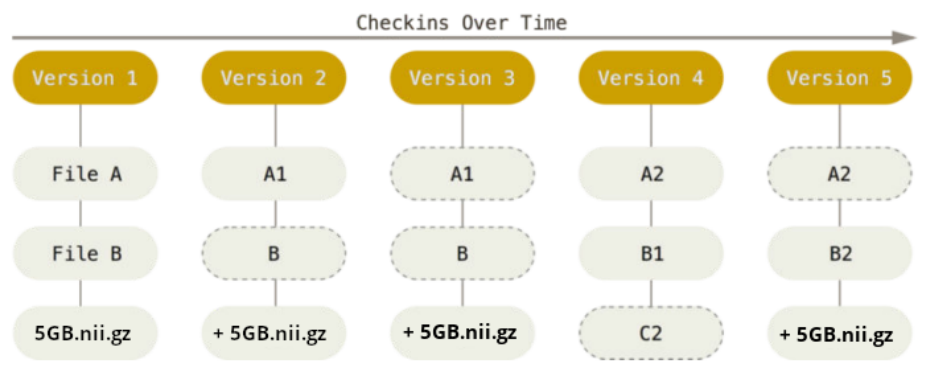
What are Git and DataLad?



git-scm.com (by Jason Long; [CC BY 3.0 Unported](https://creativecommons.org/licenses/by/3.0/))

- most popular version control system
- free, [open-source](#) command-line tool
- graphical user interfaces exist, e.g., [GitKraken](#)
- standard tool in the software industry
- 100 million [GitHub](#) users ¹

Sadly, Git does not handle large files well.

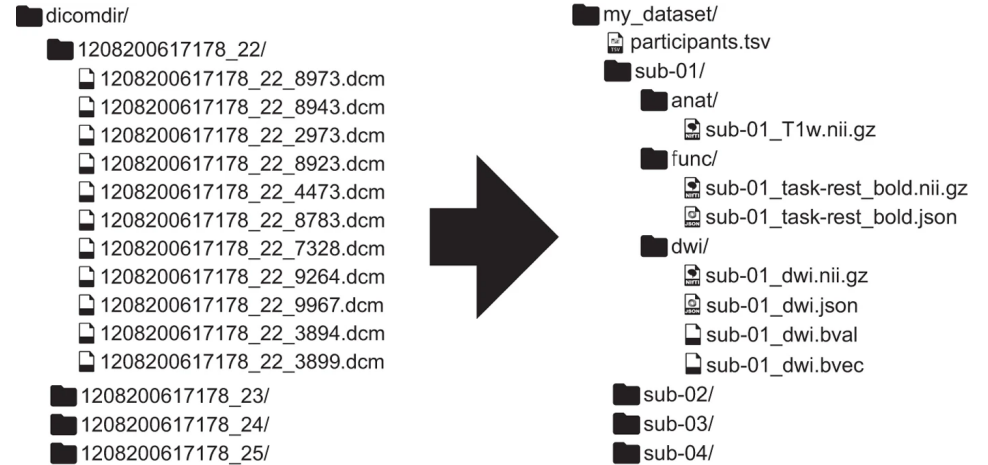
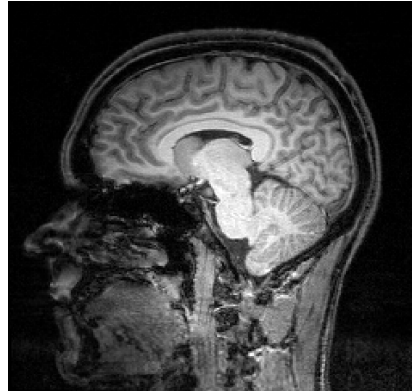


datalad.org (from the [DataLad Handbook](#) by Wagner et al. (2022); [CC BY-SA 4.0](https://creativecommons.org/licenses/by-sa/4.0/))

- “Git for (large) data”
- free, [open-source](#) command-line tool
- builds on top of [Git](#) and [git-annex](#)
- **allows to version control arbitrarily large datasets** ²
- [DataLad Python API](#): Use DataLad in your Python code
- graphical user interface exists: [DataLad Gooney](#)

Example Dataset: Brain Imaging Data

Single subject epoch (block) auditory fMRI activation data



```
1 mkdir neuro-data
2 wget https://www.fil.ion.ucl.ac.uk/spm/download/data/MoAEpilot/MoAEpilot
3 -O neuro-data.zip
4 unzip neuro-data.zip -d neuro-data
5 rm neuro-data.zip
6 cd neuro-data
7 mv MoAEpilot/* .
8 rm -R MoAEpilot
```

Brain Imaging Data Structure (BIDS) Gorgolewski et al. (2016) (CC BY 4.0)

```
1 tree
2 .
3 |--- CHANGES
4 |--- README
5 |--- dataset_description.json
6 |--- sub-01
7 |   |--- anat
8 |   |   |--- sub-01_T1w.nii
9 |   |   |--- func
10 |   |       |--- sub-01_task-auditory_bold.nii
11 |   |       |--- sub-01_task-auditory_events.tsv
12 |--- task-auditory_bold.json
13
14 4 directories, 7 files
```

Dataset from Functional Imaging Laboratory, UCL Queen Square Institute of Neurology, London, UK (Source)

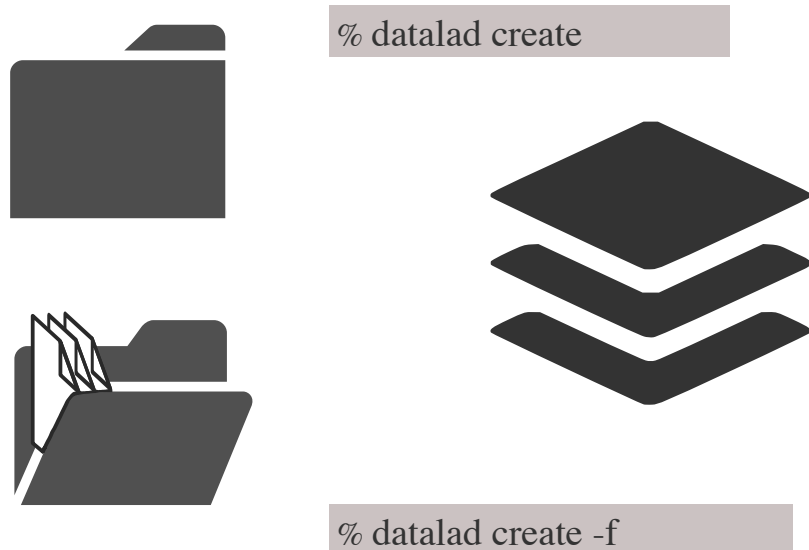


Brain Imaging Data Structure (BIDS) Gorgolewski et al. (2016) (CC BY 4.0)

Reproducible Research Data Management with DataLad

Version Control with DataLad

create new, empty datasets to populate...



.... or transform existing directories into datasets

from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

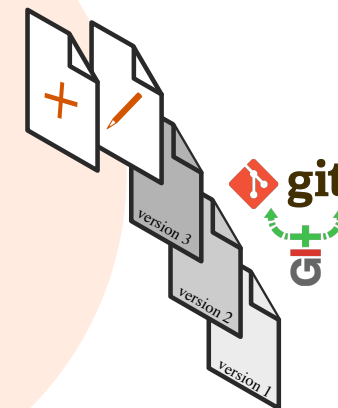
```
1 datalad create neuro-data
```

► View output

modify the dataset

save changes

% datalad save

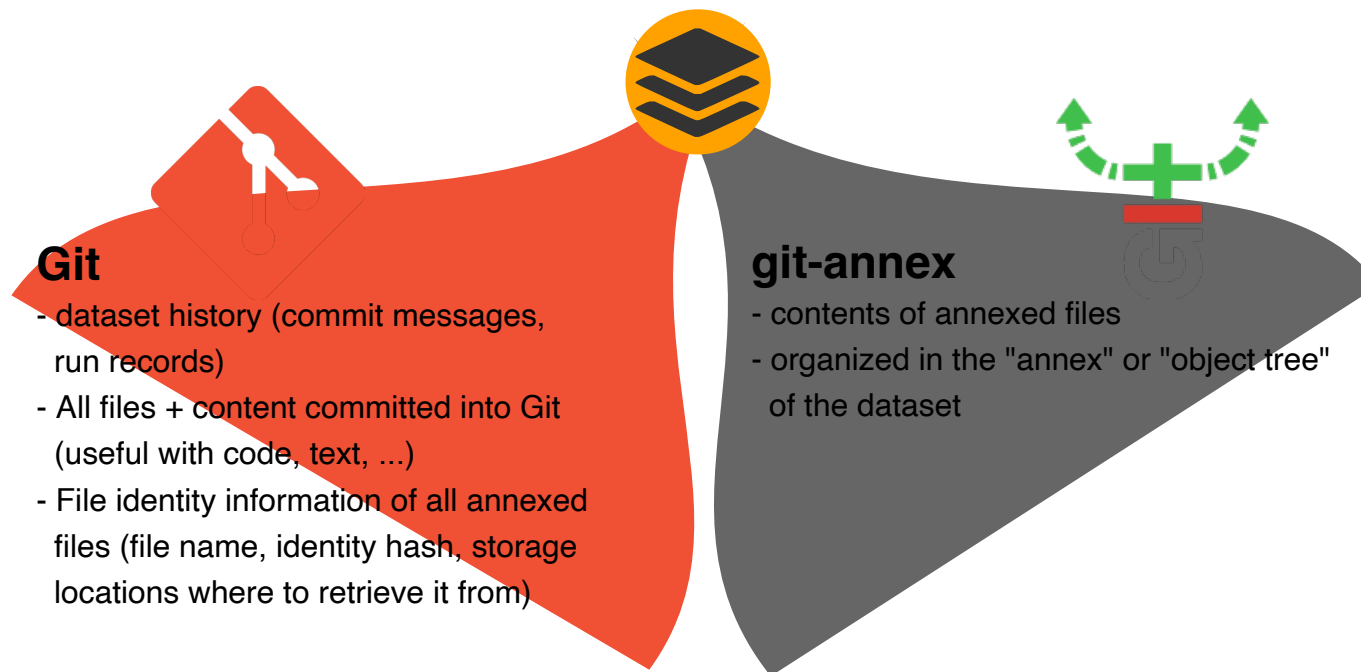


from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

```
1 datalad save -m "save neuro data"
```

► View output

Data in DataLad datasets are either stored in Git or git-annex



from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

Git

- handles **small files** well (text, code)
- file contents are in Git history and will be **shared**
- Shared with every dataset clone
- Useful: small, non-binary, frequently modified files

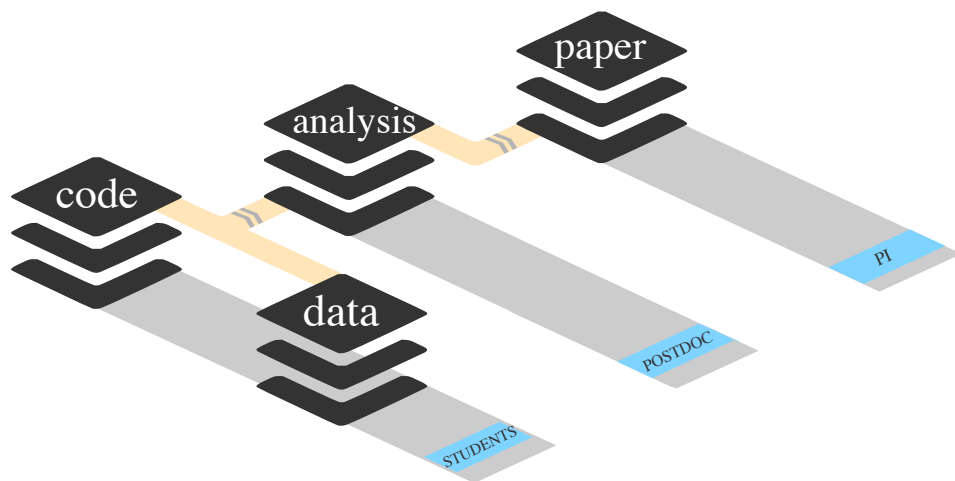
git-annex

- handles **all** types and sizes of files well
- file contents are in the annex, not necessarily shared
- **Can be kept private** on a per-file level
- Useful: Large files, private files

Science is build from modular units
We need modularity and linking

Version control beyond single repositories

Research as a sequence

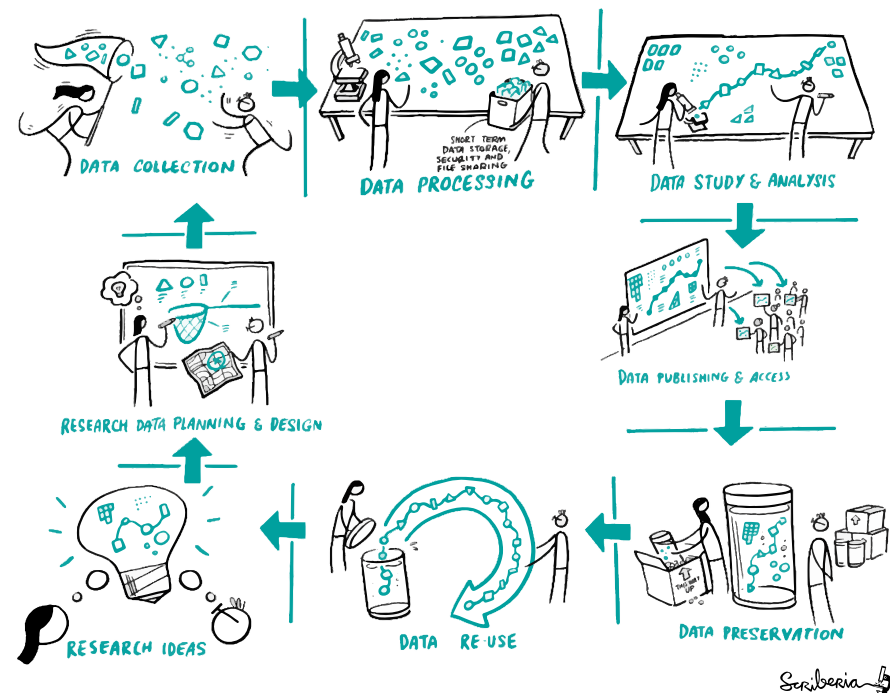


from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

- Prior works (code development, empirical data, etc.) are combined to produce results with goal of a publication
- Aggregation across time and contributors
- Aiming for (but often failing) to be reproducible
- Often, there is one big project folder

A single repository is not enough!

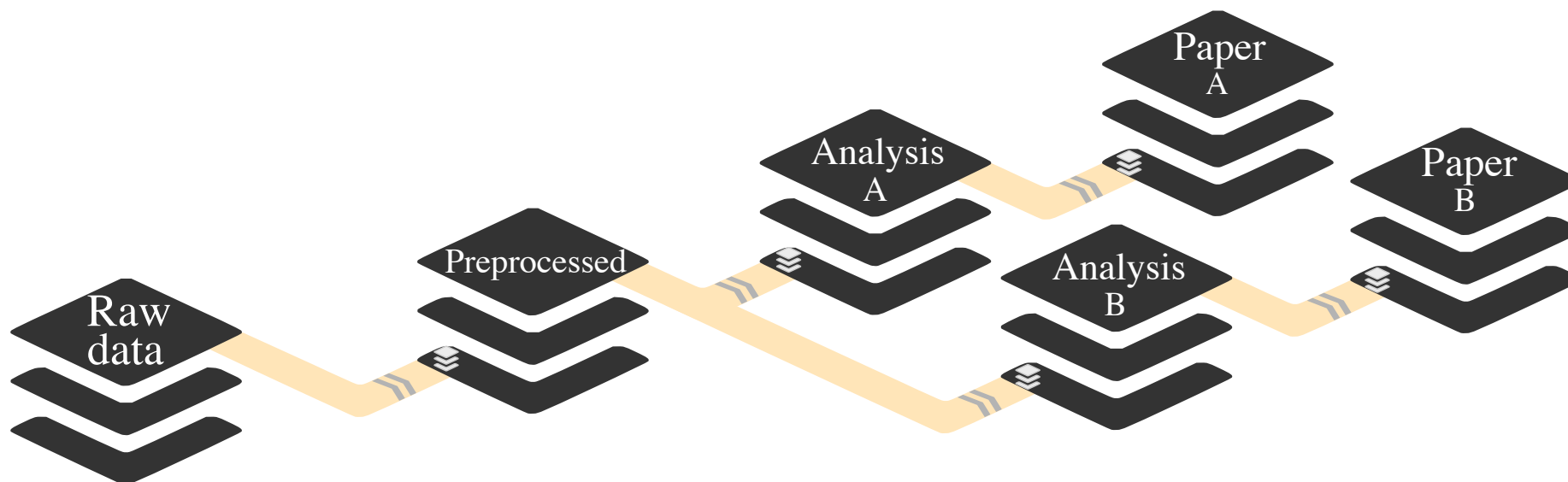
Research as a cycle



by The Turing Way Community and Scriberia (2024) (CC BY 4.0)

- Develop scientific outputs as modular but linked units
- Independently update and develop data sources
- Manage access to public / private datasets

Nesting of modular DataLad datasets

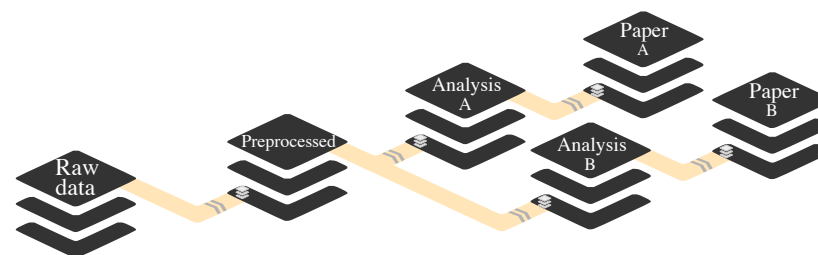


Nest modular datasets to create a linked hierarchy of datasets, and enable recursive operations throughout the hierarchy

from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

- seamless nesting of modular datasets in hierarchical super-/sub-dataset relationships
- based in Git submodules, but mono-repo feel thanks to recursive operations
- overcomes scaling issues with large amounts of files (Example: [Human Connectome Project](#))
- modularizes research components for transparency, reuse and access management

Example: Intuitive data analysis structure



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

from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

First, let's create a new data analysis dataset:

```
1 datalad create -c yoda myanalysis
2 [INFO ] Creating a new annex repo at /tmp/myanalysis
3 [INFO ] Scanning for unlocked files (this may take some time)
4 [INFO ] Running procedure cfg_yoda
5 [INFO ] == Command start (output follows) =====
6 [INFO ] == Command exit (modification check follows) =====
7 create(ok): /tmp/myanalysis (dataset)
```

-c **yoda** initializes useful structure (details [here](#)):

```
1 tree
2 .
3 |— CHANGELOG.md
4 |— README.md
5 |— code
6 |   |— README.md
7 2 directories, 3 files
```

We install analysis input data as a subdataset to the dataset:

```
1 datalad clone -d . https://github.com/datalad-handbook/iris_data.git input/
2 [INFO ] Remote origin not usable by git-annex; setting annex-ignore
3 install(ok): input (dataset)
4 add(ok): input (dataset)
5 add(ok): .gitmodules (file)
6 save(ok): . (dataset)
7 action summary:
8   add (ok: 2)
9   install (ok: 1)
10  save (ok: 1)
```

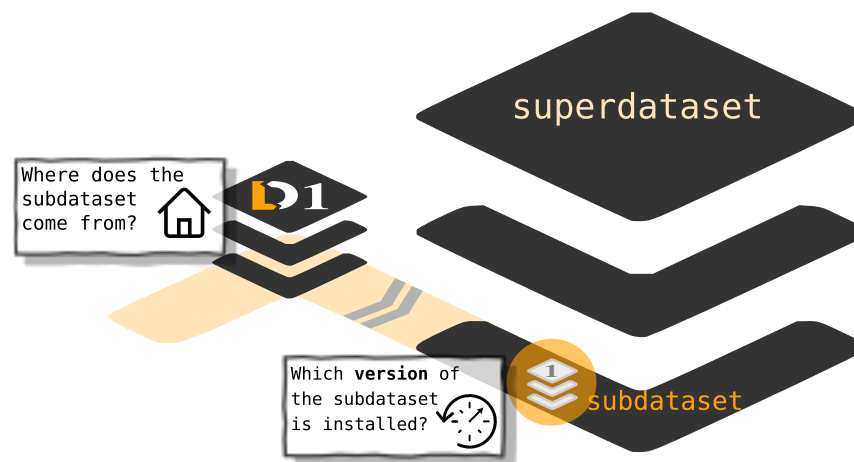
input is a regular folder inside **myanalysis**

```
1 tree
2 .
3 |— CHANGELOG.md
4 |— README.md
5 |— code
6 |   |— README.md
7 |— input
8 |   |— iris.csv
9 3 directories, 4 files
```

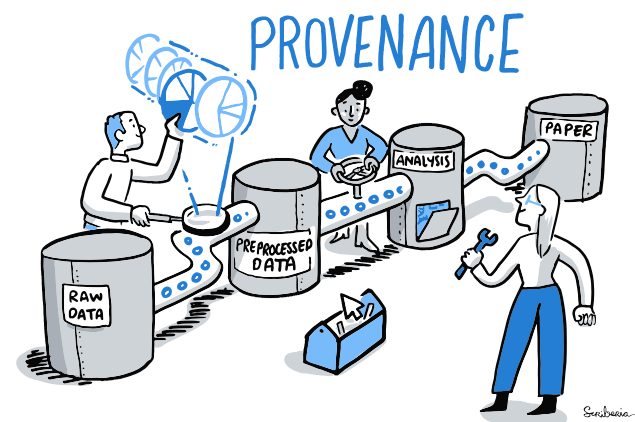
Modular units with clear provenance

```
1 git diff HEAD~1
2 diff --git a/.gitmodules b/.gitmodules
3 new file mode 100644
4 index 0000000..fc69c84
5 --- /dev/null
6 +++ b/.gitmodules
7 @@ -0,0 +1,5 @@
8 +[submodule "input"]
9 +   path = input
10 +   url = https://github.com/datalad-handbook/iris_data.git
11 +   datalad-id = 5800e71c-09f9-11ea-98f1-e86a64c8054c
12 +   datalad-url = https://github.com/datalad-handbook/iris_data.git
13 diff --git a/input b/input
14 new file mode 160000
15 index 0000000..b9eb768
16 --- /dev/null
17 +++ b/input
18 @@ -0,0 +1 @@
19 +Subproject commit b9eb768c145e4a253d619d2c8285e540869d2021
```

- We know *exactly* where the subdataset comes from
- We know *exactly* which version of the subdataset is installed
- We can develop and update each subdataset independently



from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)



by The Turing Way Community and Scriberia (2024) (CC BY 4.0)

Science is exploratory and iterative
We need provenance

Establishing provenance with DataLad

`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.



from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

- Enshrine the analysis in a script and record code execution together with input data, output files and software environment in the execution command
- Result: Machine readable record about which data, code and software produced a result how, when and why
- Use the unique identifier (hash) of the execution record to have a machine recompute and verify past work

```
1 datalad containers-run \  
2 --message "Time series extraction from Locus Coeruleus"  
3 --container-name nilearn \  
4 --input 'mri/*_bold.nii' \  
5 --output 'sub-*/LC_timeseries_run-*.csv' \  
6 "python3 code/extract_lc_timeseries.py"  
7  
8 -- Git commit --  
9 commit 5a7565a640ff6de67e07292a26bf272f1ee4b00e  
10 Author: Adina Wagner adina.wagner@t-online.de  
11 AuthorDate: Mon Nov 11 16:15:08 2019 +0100  
12  
13 [DATALAD RUNCMD] Time series extraction from Locus Coeruleus  
14 === Do not change lines below ===  
15 {  
16 "cmd": "singularity exec --bind {pwd} .datalad/environments/nilearn  
17 "dsid": "92ealfaa-632a-11e8-af29-a0369f7c647e",  
18 "inputs": [  
19 "mri/*.bold.nii.gz",  
20 ".datalad/environments/nilearn.simg"  
21 ],  
22 "outputs": ["sub-*/LC_timeseries_run-*.csv"],  
23 ...  
24 }  
25 ^^^ Do not change lines above ^^^
```

```
1 datalad rerun 5a7565a640ff6de67  
2 [INFO ] run commit 5a7565a640ff6de67; (Time series extraction from  
3 [INFO ] Making sure inputs are available (this may take some time)  
4 get(ok): mri/sub-01_bold.nii (file)  
5 [...]  
6 [INFO ] == Command start (output follows) =====  
7 [INFO ] == Command exit (modification check follows) =====  
8 add(ok): sub-01/LC_timeseries_run-*.csv(file)
```

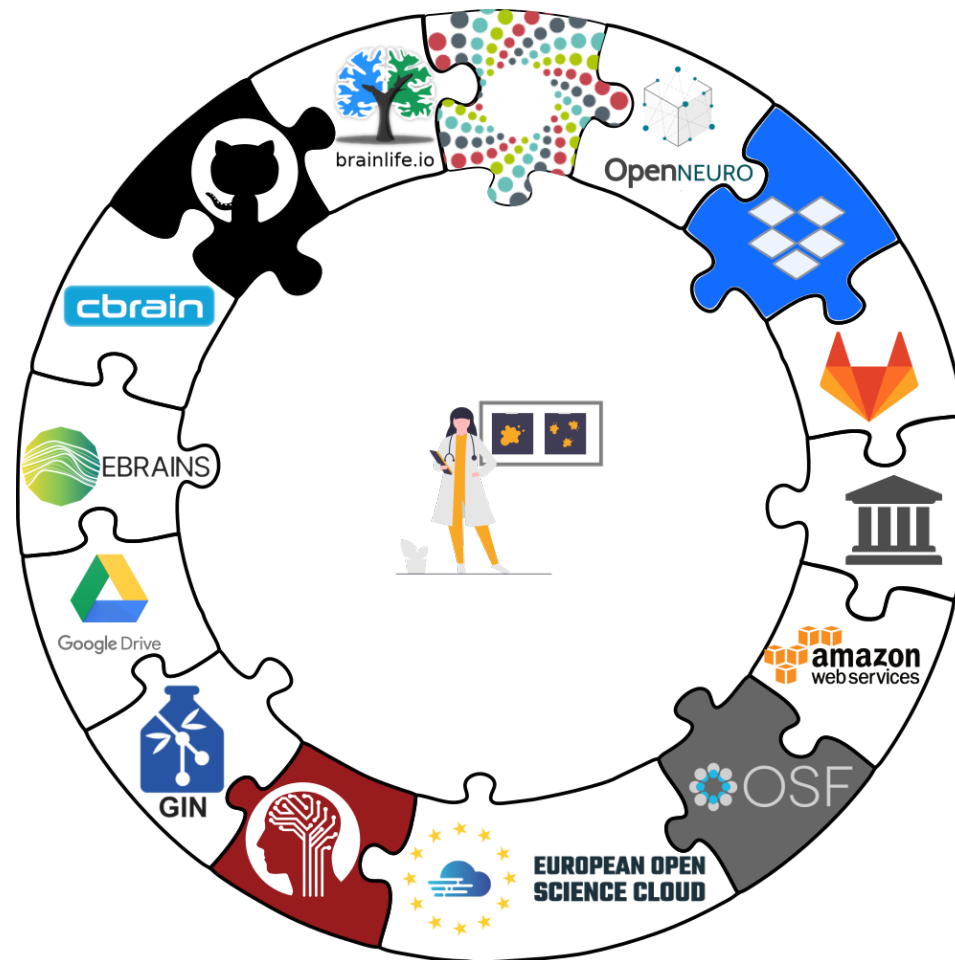
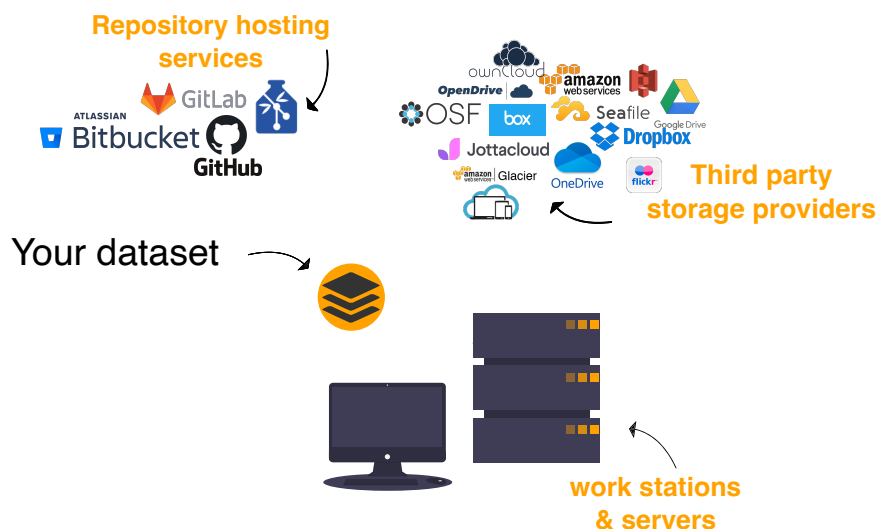
Science is collaborative & distributed

We need interoperability & transport logistics

Data sharing and collaboration with DataLad

“I have a dataset on my computer.
How can I share it or collaborate on it?”

Challenge: Scientific workflows are idiosyncratic across institutions / departments / labs / any two scientists



```

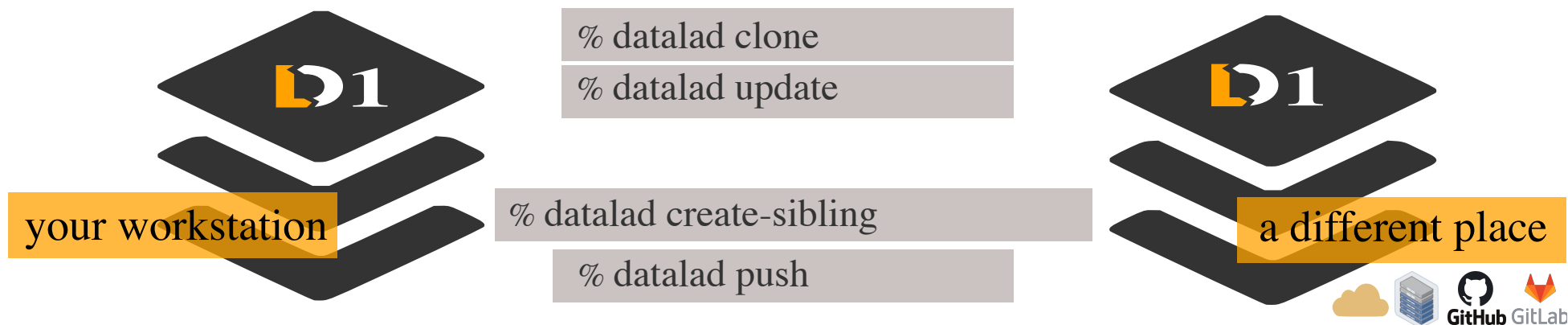
1 tree
2 .
3 |--- CHANGES
4 |--- README
5 |--- dataset_description.json
6 |--- sub-01
7 |   |--- anat
8 |   |   |--- sub-01_T1w.nii
9 |   |   |--- func
10 |   |       |--- sub-01_task-auditory_bold.nii
11 |   |       |--- sub-01_task-auditory_events.tsv
12 |--- task-auditory_bold.json
13
14 4 directories, 7 files
  
```

from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

Share data like code

- With DataLad, you can **share data like you share code**: As version-controlled datasets via repository hosting services
- DataLad datasets can be cloned, pushed and updated from and to a wide range of remote hosting services

Consume existing datasets and stay up-to-date

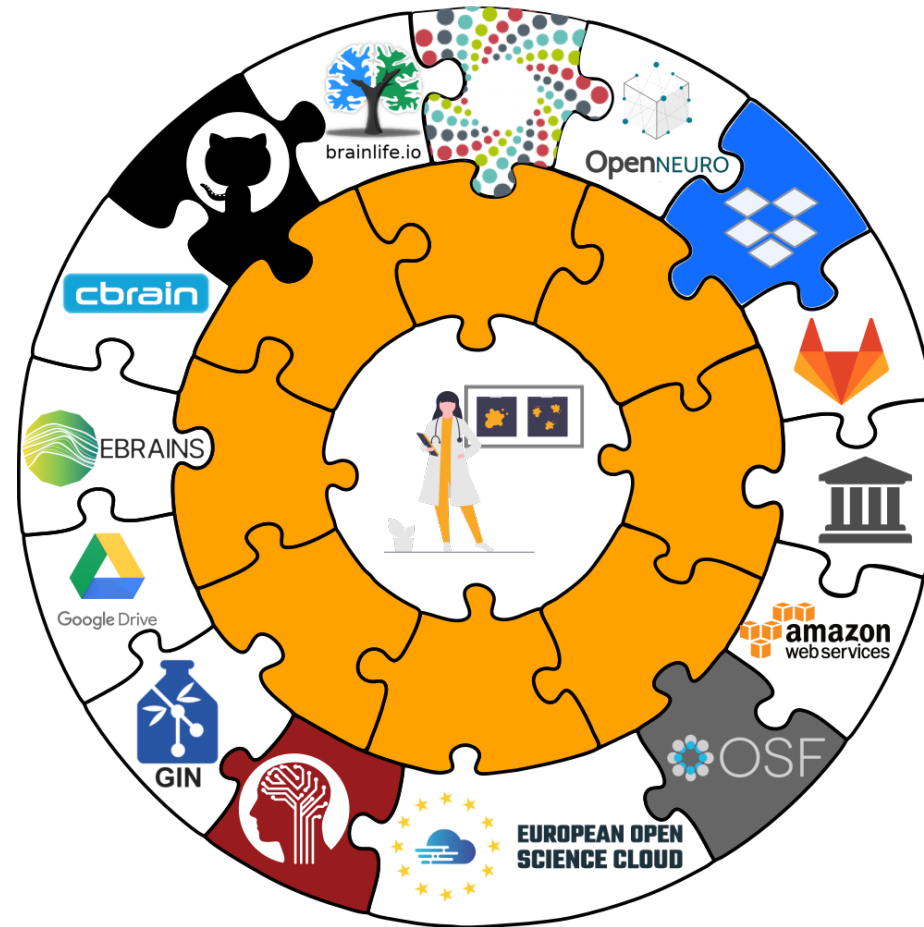


Create sibling datasets to publish to or update from

from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

Interoperability with a range of hosting services

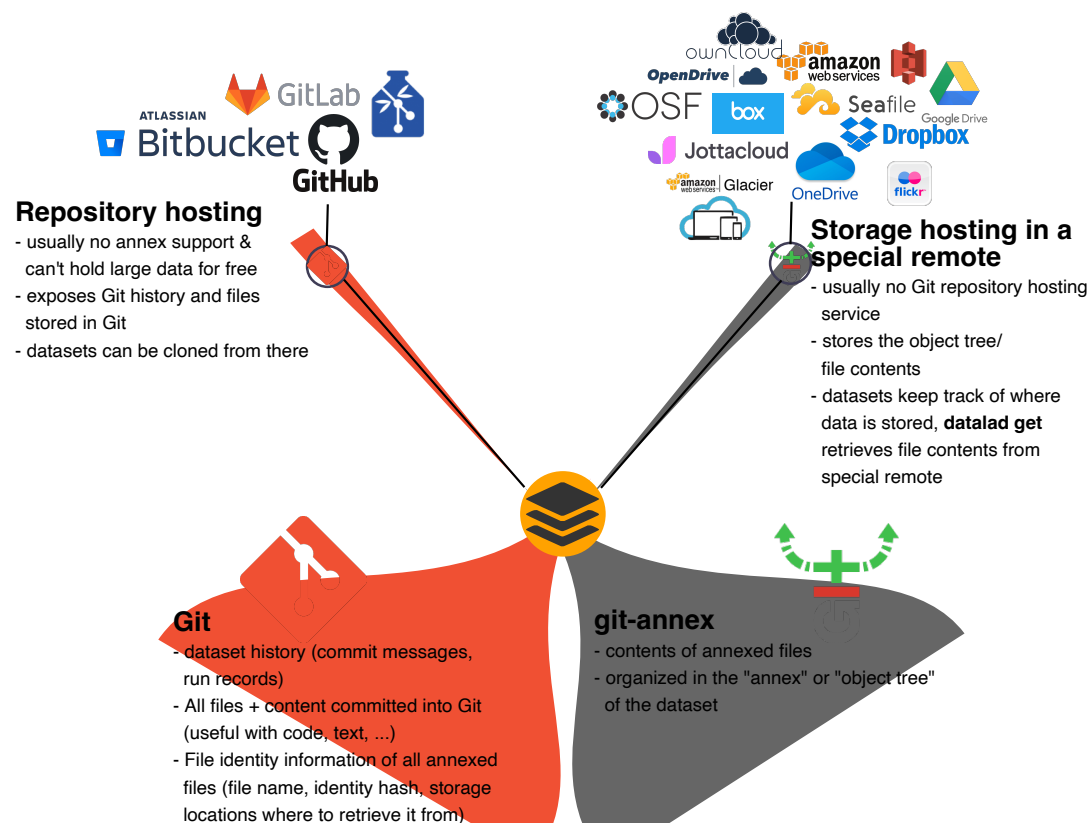
DataLad is built to maximize interoperability and streamline routines across hosting services and storage technology



see DataLad Handbook: [“Beyond shared infrastructure”](#)

Separate content in Git vs. git-annex behind the scenes

- DataLad datasets are exposed via private or public repositories on a repository hosting service (e.g., GitLab or GitHub)
- Data can't be stored in the repository hosting service but can be kept in almost any third party storage
- Publication dependencies automate interactions between both paces



from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

Special cases

Repositories with annex support

Repositories with annex support

- examples: GIN (gin.g-node.org), GitLab instances with enabled annex support
- can hold large data for free
- exposes Git history and all files + content
- datasets can be cloned from there



from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)



- Easy: Only one remote repository
- Examples: [GIN](#), GitLab with annex support

gin.g-node.org

Special remotes with repositories

Publishing both components to a special remote

- performed with DataLad extensions (e.g., `datalad-osf`), or helpers (e.g., `git-remote-rclone`)
- publishes the Git repository (in two files) in conjunction with annexed data
- allows cloning from a special remote



from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)



- Flexible: Full history or single snapshot
- Examples: [DataLad-OSF](#)

[DataLad-OSF](#)

Have access to more data than you have disk-space

Cloned datasets are lean.

```
1 datalad clone git@gin.g-node.org:/lnnrtwttkhn/neuro-data.git
2 install(ok): /tmp/neuro-data (dataset)
3 cd neuro-data && du -sh
4 212K
```

“Metadata” (file names, availability) are present ...

```
1 tree
2 .
3 |— CHANGES
4 |— README
5 |— dataset_description.json
6 |— sub-01
7 |   |— anat
8 |   |   |— sub-01_T1w.nii
9 |   |   |— func
10 |   |       |— sub-01_task-auditory_bold.nii
11 |   |       |— sub-01_task-auditory_events.tsv
12 |— task-auditory_bold.json
13
14 4 directories, 7 files
```

... but no file content:

```
1 open README
2 The file /tmp/README does not exist.
```

File contents can be retrieved on demand:

```
1 datalad get .
2 get(ok): CHANGES (file) [from origin...]
3 get(ok): README (file) [from origin...]
4 get(ok): dataset_description.json (file) [from origin...]
5 get(ok): sub-01/anat/sub-01_T1w.nii (file) [from origin...]
6 get(ok): sub-01/func/sub-01_task-auditory_bold.nii (file) [from origin...]
7 get(ok): sub-01/func/sub-01_task-auditory_events.tsv (file) [from origin...]
8 action summary:
9 get (ok: 6)
```

Let’s check the dataset size again:

```
1 du -sh
2 49M
```

Drop file content that is not needed:

```
1 datalad drop .
2 drop(ok): CHANGES (file) [locking origin...]
3 drop(ok): README (file) [locking origin...]
4 drop(ok): dataset_description.json (file) [locking origin...]
5 drop(ok): sub-01/anat/sub-01_T1w.nii (file) [locking origin...]
6 drop(ok): sub-01/func/sub-01_task-auditory_bold.nii (file) [locking origin...]
7 drop(ok): sub-01/func/sub-01_task-auditory_events.tsv (file) [locking origin...]
8 drop(ok): . (directory)
9 action summary:
10 drop (ok: 7)
```

When files are dropped, only “metadata” stays behind, and files can be re-obtained on demand.

Data sharing using DataLad and data infrastructure of the Max Planck Society

Sharing DataLad datasets via Keeper



“A free service for all Max Planck employees and project partners with **more than 1TB of storage per user** for your research data.”

keeper.mpg.de

✨ Features of Keeper ✨

- > 1 TB per Max Planck employee (and expandable):
- based on cloud-sharing service [Seafile](#)
- data hosted on MPS servers
- configurable as a [DataLad special remote](#)

1. Configure [rclone](#)

```
1 rclone config create neuro-data seafile \  
2 url https://keeper.mpg.de/ user wittkuhn@mpib-berlin.mpg.de \  
3 library neuro-data pass supersafepassword
```

2. Create a library on Keeper and a Keeper sibling

```
1 git annex initremote keeper type=external externaltype=rclone \  
2 chunk=50MiB encryption=none target=neuro-data
```

3. Push dataset to Keeper

```
1 datalad push --to keeper
```

A screenshot of the Keeper web interface. At the top, there are navigation buttons: Upload, New, Share, Archive, Certify, and Metadata. A search bar is on the right with the text "Search files in this library". Below the navigation is a breadcrumb trail: "Libraries / neuro-data / git-annex". The main content is a table with columns for "Name", "Size", and "Last Update". Each row represents a library and includes a checkbox, a star icon, and a folder icon. The libraries listed are: 2de, 179, 972, c60, deb, f2f, and fba. All "Last Update" entries are "a few seconds ago".

<input type="checkbox"/>	Name ^	Size	Last Update
<input type="checkbox"/>	☆ 📁 2de		a few seconds ago
<input type="checkbox"/>	☆ 📁 179		a few seconds ago
<input type="checkbox"/>	☆ 📁 972		a few seconds ago
<input type="checkbox"/>	☆ 📁 c60		a few seconds ago
<input type="checkbox"/>	☆ 📁 deb		a few seconds ago
<input type="checkbox"/>	☆ 📁 f2f		a few seconds ago
<input type="checkbox"/>	☆ 📁 fba		a few seconds ago

Sharing DataLad datasets via Edmond



edmond.mpg.de

“Edmond is a research data repository for Max Planck researchers. It is the place to store completed datasets of research data with open access.”

✨ Features of Edmond ✨

- based on [Dataverse](#), hosted on MPS servers
- use is free of charge
- no storage limitation (on datasets or individual files)
- flexible licensing

Two modes:

1. **annex mode** (default): non-human readable representation of the dataset that includes Git history and annexed data
2. **filetree mode**: human readable single snapshot of your dataset “as it currently is” that does not include history of annexed files (but Git history)

1. Create a Dataverse sibling for Edmond:

```
1 data-lad add-sibling-dataverse https://edmond.mpg.de/ \  
2 doi:10.17617/3.8LDVXX --mode filetree
```

2. Push dataset to Edmond / Dataverse

```
1 data-lad push --to dataverse
```

The screenshot shows the Edmond Dataverse interface for a dataset named "neuro-data". The dataset is in "Draft" and "Unpublished" status. The description is "Wittkuhn, Lennart, 2024, 'neuro-data', <https://doi.org/10.17617/3.8LDVXX>, Edmond, DRAFT VERSION". The license is "Public Domain" (CC0 1.0). The dataset is described as "Single subject epoch (block) auditory fMRI activation data". The interface includes tabs for "Files", "Metadata", "Terms", and "Versions". The "Files" tab is active, showing a list of files with columns for file name, size, and date deposited. The files listed are: "CHANGES" (Plain Text, 71 B, deposited Mar 20, 2024), "dataset_description.json" (JSON, 204 B, deposited Mar 20, 2024), "README" (Plain Text, 1.3 KB, deposited Mar 20, 2024), and "task-2D-auditory_bold.json" (JSON, 466 B, deposited Mar 20, 2024). The interface also includes a search bar, a filter by "File Type" and "Access", and buttons for "Upload Files", "Sort", "Edit Files", and "Download".

Sharing DataLad datasets via ownCloud / Nextcloud



owncloud.gwdg.de



nextcloud.com

DataLad NEXT extension allows to push / clone DataLad datasets to / from ownCloud & Nextcloud (via WebDAV)

ownCloud GWDG: “50 GByte default storage space per user; flexible increase possible upon request”

✨ Features of ownCloud and Nextcloud ✨

- data privacy compliant alternative to Google Drive, Dropbox, etc. (usually hosted on-site)
- provided by your institution, so free to use
- supports private and public repositories
- can be used together with external collaborators
- expose datasets for regular download without DataLad

1. Create a WebDAV sibling:

```
1 datalad create-sibling-webdav --dataset . \  
2   --name owncloud-gwdg --mode filetree \  
3   'https://owncloud.gwdg.de/remote.php/nonshib-webdav/neuro-data'
```

2. Push dataset to ownCloud

```
1 datalad push --to owncloud-gwdg
```

The screenshot shows the ownCloud file manager interface. The top bar displays 'Files' and 'GWDG ownCloud' with a search icon and the user 'wittkuhn@mpib-berlin.mpg.de'. The main area shows a directory view for 'neuro-data' containing the following files and folders:

Name	Size	Modified
.datalad	< 1 KB	seconds ago
sub-01	48.8 MB	2 minutes ago
.gitattributes	< 1 KB	3 minutes ago
CHANGES	< 1 KB	3 minutes ago
dataset_description.json	< 1 KB	3 minutes ago
README	1 KB	3 minutes ago
task-auditory_bold.json	< 1 KB	3 minutes ago

At the bottom of the file list, it shows '2 folders and 5 files' and a total size of '48.8 MB'.

Access the dataset [on owncloud.gwdg.de](https://owncloud.gwdg.de)

Sharing DataLad datasets via GIN



“GIN is [...] a web-accessible repository store of your data **based on git and git-annex** that you can access securely anywhere you desire while keeping your data in sync, backed up and easily accessible [...]”

gin.g-node.org

✨ Features of GIN ✨

- free to use and open-source (could be hosted within your institution; for more details, see [here](#))
- currently unlimited storage capacity and no restrictions on individual file size
- supports private and public repositories
- publicly funded by the Federal Ministry of Education and Research (BMBF; details [here](#))
- servers on German land (Munich, Germany; cf. GDPR)
- provides Digital Object Identifiers (DOIs) (details [here](#)) and allows free licensing (details [here](#))

1. Create a GIN sibling

```
1 datalad siblings add --dataset . \  
2 --name gin --url git@gin.g-node.org:/lnnrtwttkhn/neuro-data.git
```

2. Push dataset to GIN

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

File	Commit Hash	Action	Time
.datalad	e9ea13e2bf	[DATALAD] new dataset	22 minutes ago
sub-01	9916e58287	save neuro data	21 minutes ago
.gitattributes	e9ea13e2bf	[DATALAD] new dataset	22 minutes ago
CHANGES	9916e58287	save neuro data	21 minutes ago
README	9916e58287	save neuro data	21 minutes ago
dataset_description.json	9916e58287	save neuro data	21 minutes ago
task-auditory_bold.json	9916e58287	save neuro data	21 minutes ago

Access the dataset [on GIN](#)

Sharing DataLad datasets via GitLab



“GitLab is open source software to collaborate on code. Manage git repositories with fine-grained access controls that keep your code secure.”

gitlab.com

GitLab for Max Planck employees

- hosted by GWDG: gitlab.gwdg.de
- hosted by your institute, e.g., git.mpib-berlin.mpg.de

✨ Features of GitLab ✨

- free to use and open-source
- several MPS instances available (see above)
- supports private and public repositories
- use project management infrastructure (merge requests, issue boards, etc.) for your dataset projects

1. Create a GitLab sibling

```
1 datalad siblings add --dataset . --name gitlab \  
2 --url git@git.mpib-berlin.mpg.de:wittkuhn/neuro-data.git
```

2. Push dataset *metadata* to GitLab

```
1 datalad push --to gitlab
```

The screenshot shows the GitLab interface for a repository named 'neuro-data'. At the top, it indicates 'Lennart Wittkuhn / neuro-data' and shows 2 commits, 2 branches, 0 tags, and 19 KIB project storage. A recent commit titled 'save neuro data' by Lennart Wittkuhn is shown, authored 35 minutes ago with commit ID 9916e582. Below the commit, there are navigation options for 'main' branch and 'neuro-data' directory. A list of files and folders is displayed with their last commit and update times:

Name	Last commit	Last update
└ .datalad	[DATALAD] new dataset	36 minutes ago
└ sub-01	save neuro data	35 minutes ago
└ .gitattributes	[DATALAD] new dataset	36 minutes ago
└ CHANGES	save neuro data	35 minutes ago
└ README	save neuro data	35 minutes ago
└ dataset_description.json	save neuro data	35 minutes ago
└ task-auditory_bold.json	save neuro data	35 minutes ago

Access the dataset [on GitLab](#)

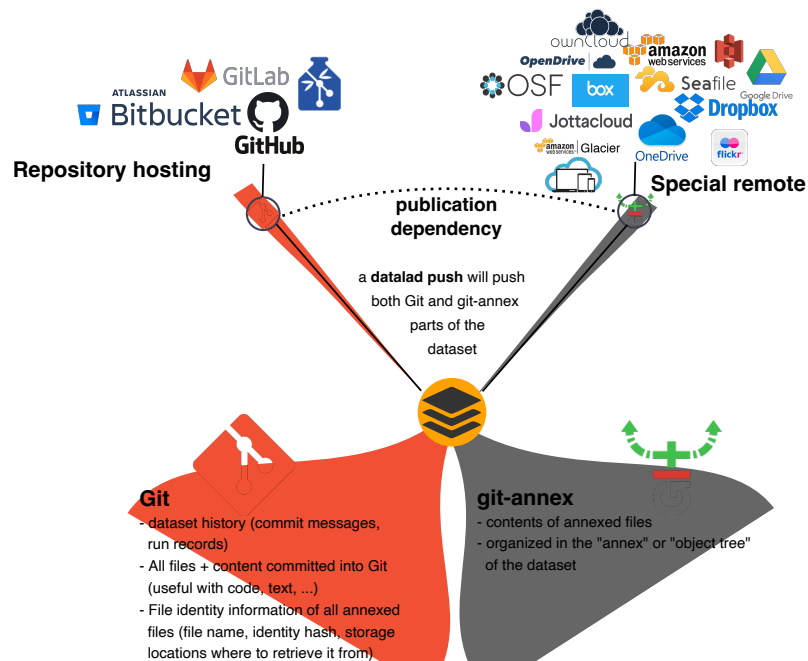
Publish and consume datasets like source code

Datasets can *comfortably* live in multiple locations:

```

1 datalad siblings
2  :: here(+) [git]
3  :: owncloud-gwdg(+) [git]
4  :: dataverse(+) [dataverse]
5  :: gin(+) [https://gin.g-node.org/lnnrtwttkhn/neuro-data (git)]
6  :: keeper(+) [rclone]
7  :: gitlab(-) [git@git.mpib-berlin.mpg.de:wittkuhn/neuro-data.git (git)]
    
```

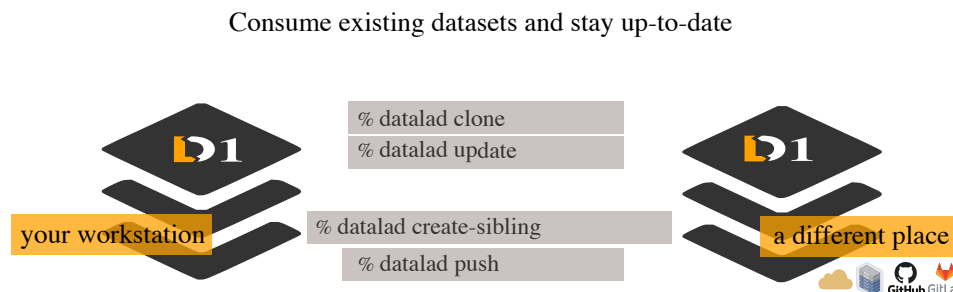
Publication dependencies automate update in all places:



```

1 datalad siblings configure --name gitlab --publish-depends SIBLING
    
```

Redundancy: DataLad gets data from available sources



Create sibling datasets to publish to or update from

from the [DataLad Handbook](#) by Wagner et al. (2022) (CC BY-SA 4.0)

Clone the dataset from GitLab:

```

1 datalad clone https://git.mpib-berlin.mpg.de/wittkuhn/neuro-data
    
```

Access to special remotes needs to be configured:

```

1 [INFO ] access to 3 dataset siblings keeper, dataverse-storage,
2 owncloud-gwdg-storage not auto-enabled, enable with:
3 | datalad siblings -d "/tmp/neuro-data" enable -s SIBLING
    
```

DataLad retrieves data from available sources (here, GIN):

```

1 datalad get .
2 get(ok): CHANGES (file) [from gin-src...]
3 get(ok): README (file) [from gin-src...]
4 [...]
    
```

Summary

Summary and discussion

Science is complex

- Scientific units are not static: We need version control
- Science is modular: We need to link modular datasets
- Science is iterative: We need to establish provenance
- Science is collaborative and distributed: We want to share our work and integrate with diverse infrastructure

DataLad: Decentralized management of digital objects for open science

- DataLad can version control arbitrary datasets
- DataLad links modular version-controlled datasets
- DataLad establishes provenance and transparency
- DataLad integrates with diverse infrastructure

🌟 Towards science as distributed, open-source **software knowledge development** 🌟 (cf. McElreath, [2020](#), [2023](#))

Develop *everything* like source code

- Code and data *management* using **Git** and **DataLad** (free, open-source command-line tools)
- Code and data *sharing* via flexible repository hosting services (**GitLab**, **GitHub**, **GIN**, etc.)
- Code and data *storage* on various infrastructure (**GIN**, **OSF**, **S3**, **Keeper**, **Dataverse**, and many more!)
- Project-related communication (ideas, problems, discussions) via **issue boards** on GitLab / GitHub etc.
- Transparent contributions to code and data via **merge requests** on GitLab (i.e., pull requests on GitHub)
- **Reproducible procedures** using `datalad run`, `rerun`, and `containers-run` commands (also [Make](#) etc.)
- *Reproducible computational environments using software containers (e.g., Docker, Apptainer, etc.)*

Overview of learning resources

Learn Git

- [“Pro Git”](#) by Scott Chacon & Ben Straub
- [“Happy Git and GitHub for the useR”](#) by Jenny Bryan, the STAT 545 TAs & Jim Hester
- [“Version Control”](#) by The Turing Way
- [“Version Control with Git”](#) by The Software Carpentries
- [“Version control”](#) (chapter 3 of “Neuroimaging and Data Science”) by Ariel Rokem & Tal Yarkoni

Learn DataLad

- [“Datalad Handbook”](#) by the DataLad team / Wagner et al., 2022, *Zenodo*
- [“Research Data Management with DataLad”](#) | Recording of a full-day workshop on YouTube
- [Datalad on YouTube](#) | Recorded workshops, tutorials and talks on DataLad

Learn both (disclaimer: shameless plug 🙌)

Full-semester course on [“Version control of code and data using Git and DataLad”](#) (v2.0!) in summer semester 2024 at University of Hamburg (generously funded by the [Digital and Data Literacy in Teaching Lab](#) program) with many open educational resources ([online guide](#), quizzes and exercises)

References

- Gorgolewski, Krzysztof J., Tibor Auer, Vince D. Calhoun, R. Cameron Craddock, Samir Das, Eugene P. Duff, Guillaume Flandin, et al. 2016. “The Brain Imaging Data Structure, a Format for Organizing and Describing Outputs of Neuroimaging Experiments.” *Scientific Data* 3 (1). <https://doi.org/10.1038/sdata.2016.44>.
- The Turing Way Community, and Scriberia. 2024. “Illustrations from the Turing Way: Shared Under CC-BY 4.0 for Reuse,” January. <https://doi.org/10.5281/ZENODO.3332807>.
- Wagner, Adina S. 2024. “DataLad: Decentralized Management of Digital Objects for Open Science.” *Zenodo*, January. <https://doi.org/10.5281/ZENODO.10556597>.
- Wagner, Adina S., Laura K. Waite, Kyle Meyer, Marisa K. Heckner, Tobias Kadelka, Niels Reuter, Alexander Q. Waite, et al. 2022. “The DataLad Handbook,” April. <https://doi.org/10.5281/ZENODO.6463273>.

Thank you!



Dr. Lennart Wittkuhn

✉ wittkuhn@mpib-berlin.mpg.de

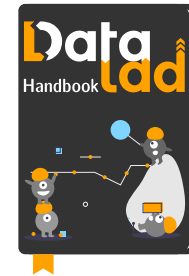
🏠 <https://lennartwittkuhn.com/>

🐙 Mastodon 🗨️ GitHub 🌐 LinkedIn



Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

MAX PLANCK INSTITUTE
FOR HUMAN DEVELOPMENT



🖥️ Slides: <https://lennartwittkuhn.com/talk-mps-fdm-2024>

🔗 Source: <https://github.com/lennartwittkuhn/talk-mps-fdm-2024>

📄 Software: Reproducible slides built with [Quarto](#) and deployed to [GitHub Pages](#) using [GitHub Actions](#) for continuous integration & deployment

📄 License: Creative Commons Attribution-ShareAlike 4.0 ([CC BY-SA 4.0](#))

🗨️ Contact: Feedback or suggestions via [email](#) or [GitHub issues](#). Thank you!