# The DataLad Handbook

A flexible, extendable & reusable open source framework for user-focused and workflow-driven software documentation

Wagner, A. S.<sup>1</sup>, Waite, L. K.<sup>1</sup>, Waite, A.Q.<sup>1</sup>, Reuter, N.<sup>1</sup> Poldrack, B.<sup>1</sup>, Poline, J-B.<sup>3</sup>, Kadelka, T.<sup>1</sup>, Markiewicz, C. J.<sup>4</sup>, Vavra, P.<sup>5</sup>, Paas, L. K.<sup>1</sup>, Herholz, P.<sup>3</sup>, Mochalski, L.<sup>1</sup>, Wiersch, L.<sup>1</sup>, Kraljevic, N.<sup>1</sup>, Heckner, M.<sup>1</sup>, Chormai, P.<sup>6</sup>, Halchenko, Y.O.<sup>7</sup> & Hanke, M.<sup>1,2</sup>

<sup>1</sup>Institute of Neuroscience and Medicine, Juelich Research Centre; <sup>2</sup>Institute of Systems Neuroscience, Heinrich Heine University Düsseldorf; <sup>3</sup>McGill University, Montreal; <sup>4</sup>Stanford University, Stanford; <sup>5</sup>Department of Biological Psychology, OvGU Magdeburg, <sup>6</sup>Max Planck School of Cognition, Leipzig; <sup>7</sup>Dartmouth College, Hanover (NH);

# Background

DataLad (Halchenko, Hanke et al.) is a comprehensive data management tool and can help to solve various data management problems such as file size independent version control, data sharing, data storage and backup, computationally reproducible data analysis, or metadata management.

But: The functionality of any tool remains inaccessible or unknown if it is not sufficiently documented.

#### Complex software needs accessible user documentation!

The DataLad Handbook is a comprehensive documentational resource that fulfills the needs of different software user types independent of background: —

**! Trainees /** learn the tool

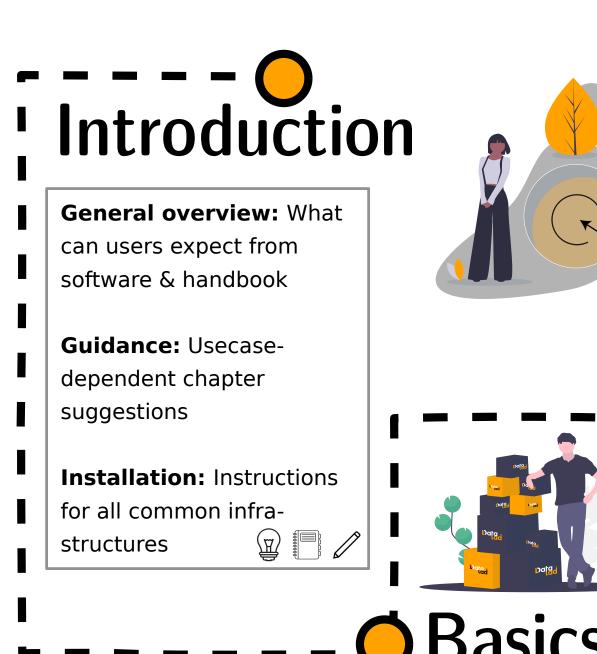
**→Planners** 

assess applicability of the tool

**Teachers** 

teach how the tool is used

#### Book structure



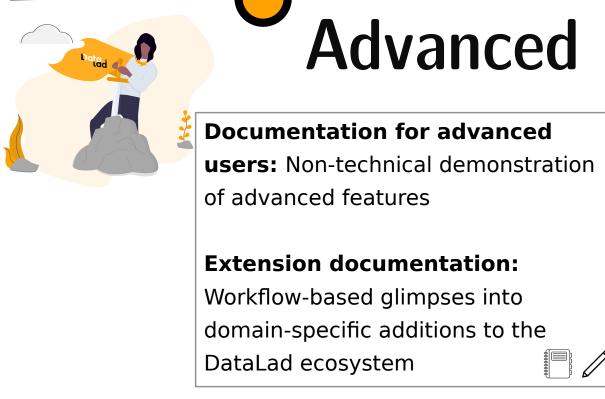


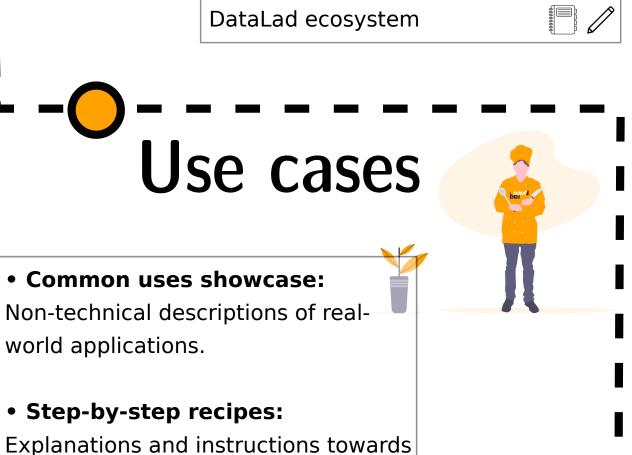
Domain-agnostic tutorial: Narrative-based codealong course with tested code snippets **Basic software skills:** Provides a broad

exploration of the software in a continuous, projectbased workflow **Trial and error:** Common errors are explicitly

demonstrated in the safe-space of a tutorial

Optional advanced information: Toggleable or custom sections contain extra information. This keeps the visible information consise, but allows for exploration of advanced contents







needed.

the described real-world application.

Links to required chapters where



### Technical infrastructure

• Flexible, extendable & reusable open source infrastructure: Python-based, written in RST markup using Sphinx, hosted on readthedocs.org, illustrations by undraw.io, source code on GitHub, continuous integration using Travis & readthedocs.org







- Multiple formats: HTML, PDF, EPUB
- Custom Sphinx extensions (github.com/mih/autorunrecord) allow code execution and record code output. The handbook in itself is a framework for workflow testing.
- Simple audience tailoring: Different branches or tags can be rendered simultaneously, allowing dedicated URLs for different content. Example: Stand-alone branch for instituteinternal workflows
- CC-BY-SA: Feel free to use the handbook infrastructure for your documentation project, e.g. Princeton Handbook for Reproducible Neuroimaging (Brooks et. al)

## Conjunct software & userdocumentation development

#### **Advantages**

- Higher rates of bug detection
- User-based documentation efforts uncover deficiencies of technical docs and user experience
- Workflow-based demos highlight API inconsistencies
- Documentation challenges facilitate software development

#### **Caveats**

- User-documentation does not replace technical docs
- Premature feature documentation: helpful for feedback & software dev facilitation, but increases documentation workload
- Separate software and user-docs rely on synchronized release management. Otherwise, unreleased functionality is documented publicly

### Contributing



☼ Create pull request

**Handbook** 

- GitHub-based development allows different contribution types
- Low-barrier contributions: General improvements, feature requests, feedback. High-barrier contributions (for advanced users): Content contributions, technical infrastructure
- "Basics/Advanced": Discussions on order/emphasis, feature requests
- "Use cases": Users contribute their DataLad workflows
- Technical infrastructure and visuals: Contributions to artwork dataset or handbook support software

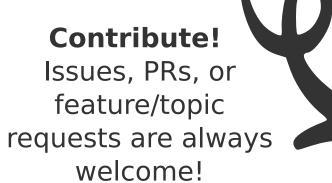
All contributions are reviewed by the DataLad core developer team

### Community and acknowledgement

- Credit is given for commit- and not commit-based contributions
- Co-authorship (PDF/EPUB + each Zenodo release), recognition with allcontributors-bot (allcontributors. org, following The Turing Way project; the-turing-way.netlify.org)
- Future directions: Presence in Hackathons/Hacktoberfest/ etc.
- Goal: Users share their individual workflows as use cases
- Current contributor count: 26

**Want to learn** more? Find the source code

on GitHub









SPONSORED BY THE and Research 01GQ1112 01GQ1411