# A full-semester course on "Version control of code and data using Git and DataLad"

Symposium on "Reproducible Research: Education and Teaching Formats Reports From the Reproducibility Networks" at QUEST Center Berlin, Germany

Talk by Dr. Lennart Wittkuhn

License: CC BY 4.0

2023-05-11

# About

## Me

👷 **Position:** PostDoc & Lab Manager at University of Hamburg & MPI for Human Development Berlin

🎓 **Education:** BSc Psychology & MSc Cognitive Neuroscience (TU Dresden), PhD Psychology (FU Berlin)

🔬 **Research:** I study the role of fast neural memory reactivation ("replay") in the human brain using fMRI

🕸️ **GRN:** Member of the MPIB's working group on research data management & open science (GRN member)

🔗 **Contact:** You can connect with me via email, Twitter, Mastodon, GitHub or LinkedIn

ℹ️ **Info:** Find out more about my work on my website, Google Scholar and ORCiD

## This presentation

🏗️ **WIP:** The presented teaching project is work in progress!

💻 **Slides:** Slides are publicly available at lennartwittkuhn.com/ddlitlab-presentation

📦 **Software:** Reproducible slides built with Quarto and deployed to GitHub Pages using GitHub Actions

⭕ **Source:** Source code is publicly available on GitHub at github.com/lnnrtwttkhn/ddlitlab-presentation

🙏 **Contact:** I am happy for any feedback or suggestions via email or GitHub issues. Thank you!

# Why we need version control ...

... for **code** (text files)



... for **data** (binary files)



© Jorge Cham (phdcomics.com)

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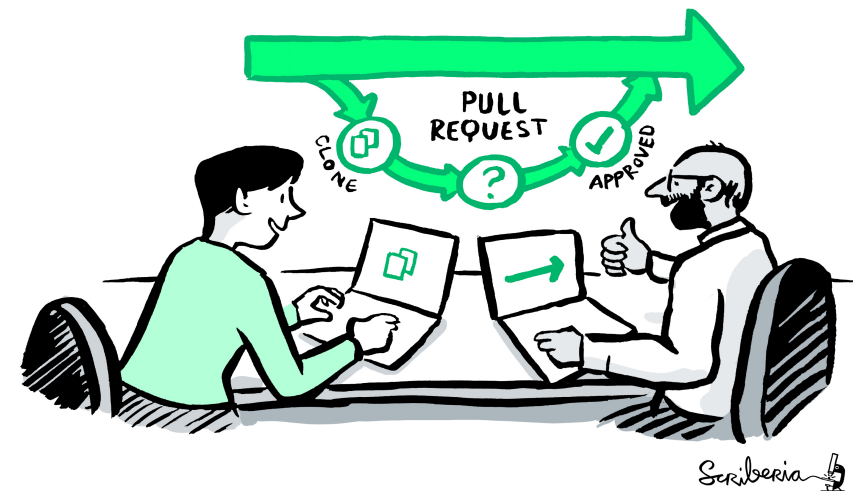
*"Version control is a systematic approach to record changes made in a [...] set of files, over time. This allows you and your collaborators to track the history, see what changed, and recall specific versions later [...]"* (Turing Way)

- 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

Version control of code and data using Git and DataLad

# What are git and DataLad?

git-scm.com

datalad.org

- most popular version control system

- free, open-source command-line tool

- graphical user interfaces exist, e.g., GitKraken

- standard tool for most (all?) software developers

- 100 million GitHub users [1]

- "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]

- graphical user interface exists: DataLad Gooey

# Course details

## Overview

🛈 Full-semester seminar (~ 12 sessions of 90 mins)

📅 Winter semester 2023/24 (October to January)

🏛 University of Hamburg (virtual option TBD)

👥 MSc and PhD students (research focus)

🧠 Psychology and Cognitive Neuroscience

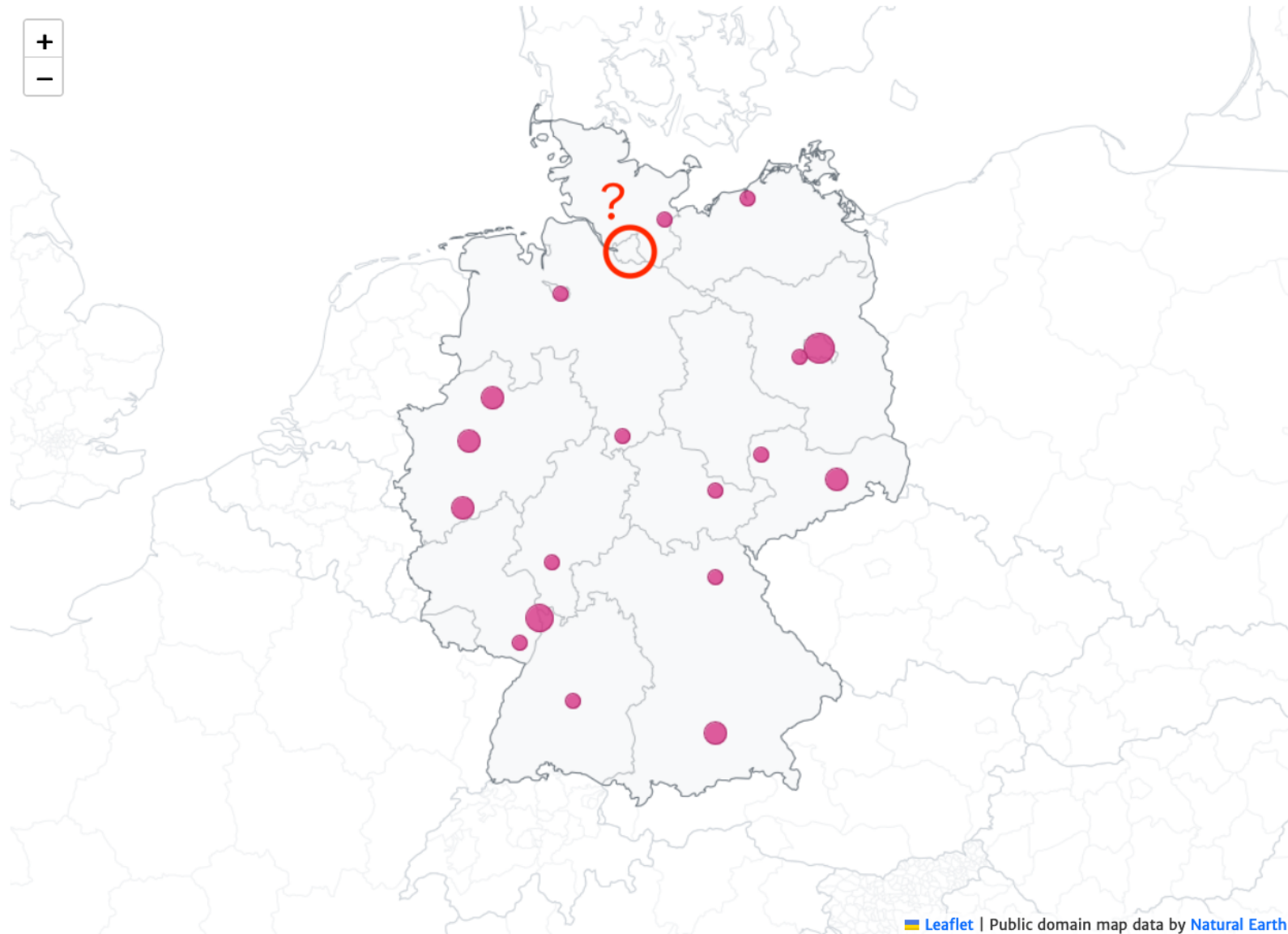## Milestones

⭐ Project received funding!

⭐ Project website online!

⭐ Course website online!

→ Next up: Prepare course!

## Implementation and Tools

🗔 **Impulse lectures** & **live demonstrations**

</> **Code-along** & **exercises** (individual and group)

💬 **Discussions** on reproducibility, open code & data

☁ Fixed computational environments on JupyterHub

>_ Focus on **command-line** interaction

🖱 Alternative use via **Graphical User Interfaces**

📋 **Quizzes** & continuous **evaluation** (in R Shiny)

♻ Reuse quiz & evaluation data as **example datasets**

▶▶ **Follow-up** research projects in summer semester

👥 Support by **research** and **teaching assistants**

🎁 Materials shared as **Open Educational Resources**

🔨 Integration with GRN, Carpentries Incubator, etc.?

# Local GRN node in Hamburg?



from reproducibilitynetwork.de/members
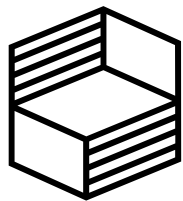
Version control of code and data using Git and DataLad

# Thank you!

## Funding & Support



Digital and Data Literacy in Teaching Lab (DDLitLab), an initiative by the ISA-Zentrum at University of Hamburg



Stiftung Innovation in der Hochschullehre

## People



Prof. Dr. Nicolas Schuck
(UHH & MPIB)

Carolin Scharfenberg
(UHH DDLitLab)

## Contact

✉ lennart.wittkuhn@uni-hamburg.de

⌂ lennartwittkuhn.com

🐦 Twitter  Ⓜ Mastodon  GitHub  🔗 LinkedIn

# Footnotes

1. (Source: Wikipedia)

2. see DataLad dataset of 80TB / 15 million files from the Human Connectome Project (see details)