

> git annex

GIT/GIT-ANNEX/DATALAD/FORGEJO-ANEKSAJO

A pragmatic data collaboration ecosystem

2025-03-21 | Matthias Riße | ICE-4

Who am I?

- Matthias RiBe
- MSc in Applied Mathematics and Computer Science and MaTSE
- Part of the IT team at ICE-4
- Software developer, Sysadmin, Data Management guy, Technical support person, ...
- Maintainer of Forgejo-aneksajo
- Generally an open source enthusiast

Overview

Vision

Problem and Solution

The ecosystem and how it works

Forgejo-aneksajo

Vision

- Every publication should have an accompanying publishable artifact that fully describes the data, code and computational environment used to produce it
- Anyone¹ should be able to go from e.g. a plot in a paper back through the processes used to produce it, be able to follow any intermediate steps, and arrive at the data, or even the data acquisition process
- Questions like this need to be answerable:
 - Where did the raw data come from?
 - What processing steps were done?
 - Which software was used, and in what version?
- This "log" should be a natural by-product of the creation process

¹with sufficient determination

The issue

```
code/
├── code_final/
│   ├── final_2/
│   │   ├── main_script_fixed.py
│   │   └── takethisscriptformostthingsnow.py
│   ├── utils_new.py
│   ├── main_script.py
│   ├── utils_new.py
│   ├── utils_2.py
│   └── main_analysis_newparameters.py
└── main_script_DONTUSE.py
data/
├── data_updated/
│   └── dataset1/
│       └── datafile_a
├── dataset1/
│   └── datafile_a
├── outputs/
│   ├── figures/
│   │   ├── figures_new.py
│   │   └── figures_final_forreal.py
│   ├── important_results/
│   ├── random_results_file.tsv
│   ├── results_for_paper/
│   ├── results_for_paper_revised/
│   └── results_new_data/
├── random_results_file.tsv
└── random_results_file_v2.tsv
[...]
```

The solution

Use a version control system for everything:

- VCS for code: git
- VCS for data: git-annex/DataLad
- VCS for computational environments: still git, but use a package manager that can uniquely define an environment for each commit: pixi, nix, guix, ...

→ every aspect is part of a self-contained repository, i.e. an artifact that fully describes the research results

Git is a decent version control system:

- tracks the evolution of content (a hierarchical directory structure) through annotated snapshots in time (commits)
- allows concurrent diverging developments (branches) ...
- ... as well as reconciliation of those developments (merge, rebase, fast-forward, cherry-pick, ...)
- distributed and local-first:
 - works offline, e.g. on an airplane or a field campaign
 - good sync support to keep different "clones" up-to-date
- essentially a standard tool already

In summary:

- git is a good decentralized, distributed, versioned database ...
- ... but it is bad at handling large files

Enter git-annex:

- a layer of indirection: git repository is an index of "annexed files"
- git stores filenames and some metadata, the actual content is in a distributed key-value store (in `.git/annex/objects`)
- each clone of a repository can contain none, some, or all of the annexed files
- git-annex provides commands and mechanism to modify what is stored where (get, drop, copy, move, sync preferred content ...)
- special (non-git) remotes incorporate other storage platforms (web, torrent, S3, webdav, rsync, rclone, ..., or [build your own](#))

And DataLad:

- terminology: datasets are git-annex-enabled git repositories with a dataset ID
- streamlines more common workflows
 - updating the index (i.e. push/pull of the git repository, downloading files, uploading new or changed files)
- provenance tracking of dataset mutations
 - What command was run to create some output(s) from some input(s)?
- Fantastic documentation: <https://handbook.data-lad.org/>
- Better support for datasets as building blocks
 - nested, reusable datasets

Datasets are building blocks

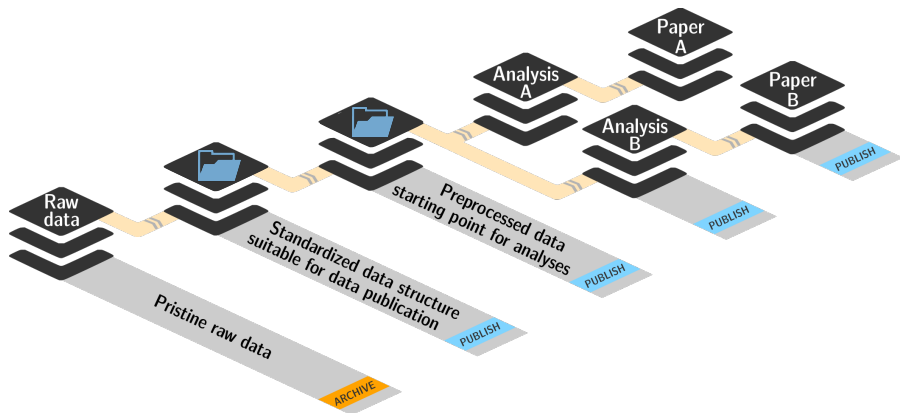


Figure:

<https://handbook.data-lad.org/en/latest/basics/101-127-yoda.html>

Live demo

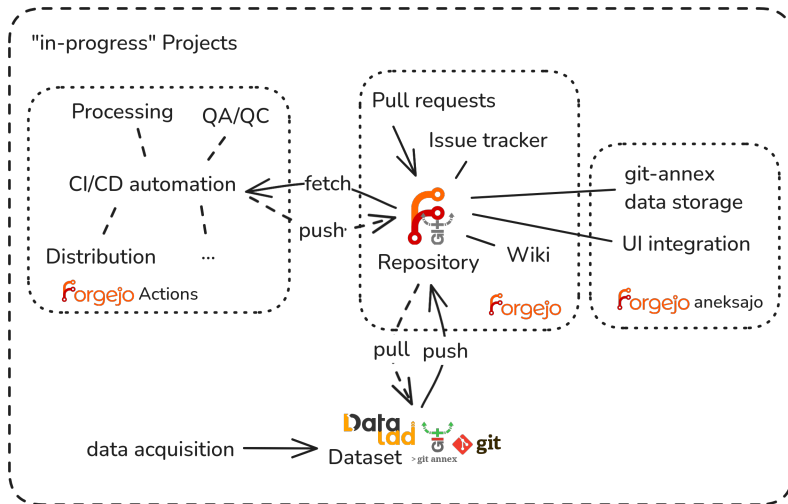
```
# To start off, we will clone the datasets repository ...
$ datalad clone git@atriis.fz-juelich.de:m.risse/era5-sat-ts-analysis-paper.git
install(ok): /tmp/tmp.wor9L9LSFn/era5-sat-ts-analysis-paper (dataset)

# ... and cd into it.
$ cd era5-sat-ts-analysis-paper

# Since we will be digging a bit deeper into it, we can prefetch all contained sub-datasets.
$ datalad get -n -r .
[INFO ] Ensuring presence of Dataset(/tmp/tmp.wor9L9LSFn/era5-sat-ts-analysis-paper) to get /tmp/tmp.wor9L9LSFn/era5-sat-ts-analysis-paper
install(ok): /tmp/tmp.wor9L9LSFn/era5-sat-ts-analysis-paper/inputs/era5-sat-ts-analysis (dataset)
install(ok): /tmp/tmp.wor9L9LSFn/era5-sat-ts-analysis-paper/inputs/era5-sat-ts-analysis/inputs/era5-sat-timeseries (dataset)
install(ok): /tmp/tmp.wor9L9LSFn/era5-sat-ts-analysis-paper/inputs/era5-sat-ts-analysis/inputs/era5-sat-timeseries/inputs/era5-2t (dataset)
action summary:
  install (ok: 3)

# This is what we are working with:
$ ls --color=always -l
total 420
-rw-rw-r-- 1 icg149 icg149   29 Mär 20 13:54 all.tikzdefs
-rw-rw-r-- 1 icg149 icg149  102 Mär 20 13:54 all.tikzdefs.license
-rw-rw-r-- 1 icg149 icg149  693 Mär 20 13:54 all.tikzstyles
-rw-rw-r-- 1 icg149 icg149  102 Mär 20 13:54 all.tikzstyles.license
-rw-rw-r-- 1 icg149 icg149 83509 Mär 20 13:54 climate-ml.bib
-rw-rw-r-- 1 icg149 icg149   102 Mär 20 13:54 climate-ml.bib.license
drwxrwxr-x 2 icg149 icg149  4096 Mär 20 13:54 code
-rw-rw-r-- 1 icg149 icg149  4489 Mär 20 13:54 fhacThesis.cls
-rw-rw-r-- 1 icg149 icg149   567 Mär 20 13:54 Flake.lock
-rw-rw-r-- 1 icg149 icg149  2326 Mär 20 13:54 Flake.nix
lrwxrwxrwx 1 icg149 icg149   128 Mär 20 13:54 Follen_kolloquium.pdf -> .git/annex/objects/Kg/Pk/MD5E-s4052039--1a437138c79caf63904c774f57dcd13.pdf/MD5E-s4052039--1a437138c79caf63904c774f57dcd13.pdf
```

Forgejo-aneksajo



Forgejo-aneksajo



- a data collaboration platform based on Forgejo
 - Forgejo is a git forge like GitHub or GitLab
- one possible storage location for git-annex repositories
- datasets are repositories and have some common features:
 - issue trackers for discussion and coordination
 - pull requests for change proposals and review
 - CI support for quality assurance, common processing steps, or other automation
 - tags and releases
 - a wiki
 - ...
- organisations provide groups and users can have different access levels in orgs and repositories

Live showcase



The screenshot shows the homepage of the ATRIS (ATmospheric Research Data Information System) website. The page has a dark blue background. At the top left, there is a navigation bar with the text 'ATRIS Explore Help' and a small logo. At the top right, there is a 'Sign in' link. The main content area features a large, stylized logo consisting of a white 'A' and a blue 'R' with circular nodes. Below the logo, the text 'ATRIS' is displayed in a large, white, sans-serif font. Underneath, the full name 'ATmospheric Research Data Information System' is written in a smaller, white, sans-serif font. A paragraph of text describes the system as a data management system designed for comprehensive and reproducible management of atmospheric research data, based on **Forgejo**. It mentions that it keeps track of changes to contents and organization of files and provides secure remote access to hosted data by utilizing **git**, **git-annex**, and **DataLad**. Another paragraph states that external users need to apply for data access authorization. A final paragraph mentions that the service is hosted by the Institute of Climate and Energy Systems - Stratosphere (ICE-4), Forschungszentrum Jülich, and provides an email contact: iek7-rdm@fz-juelich.de. At the bottom of the page, there is a footer with the text 'Powered by Forgejo Page: Tms Template: Tms' on the left and a list of links: 'English Licenses API About Contact Legal Notice Privacy Policy' on the right.

Figure: <https://atris.fz-juelich.de>

Summary

- git/git-annex/DataLad/Forgejo-aneksajo extend established tools and practices from software development to data projects
- You can version control code, data, and software environments within one unified system
- Fits into "legacy" projects as long as you are dealing with files and directories
- Repositories become self-contained research artifacts
- Forgejo-aneksajo provides a mature collaboration platform backed by a large developer community

Some links

- DataLad: <https://www.datalad.org/>
- DataLad Handbook: <https://handbook.datalad.org/>
- Distribits, a community conference on distributed data management (23.10. - 25.10.2025): <https://distribits.live/>
- git-annex: <https://git-annex.branchable.com/>
- git: <https://git-scm.com/>
- Forgejo-aneksajo:
<https://codeberg.org/forgejo-aneksajo/forgejo-aneksajo>

Thank you for your attention!