



<http://git-annex.branchable.com>

git-annex

is a distributed version control system for large files developed by Joey Hess. `git-annex` enriches regular `git` repositories with meta-information and facilities for retrieving and/or sending data from/to remote locations. Actual data itself is not directly committed to `git`, making `git-annex` repositories lightweight and easily manageable, while still providing access to all data upon request.

git-annex features

Data versioning

Every file is identified by its checksum (*e.g.*, SHA256), thus guaranteeing unambiguous versioning

Special remotes

Besides accessing data from other `git` remotes (*e.g.*, via `ssh`), data can be offloaded to the cloud (*e.g.*, Amazon S3 and Glacier), or just some remote computer (via `ssh` or `rsync`)

Association with web origins

Every file can have a list of URLs from which it can be downloaded

Encryption

Data can be uploaded to remotes encrypted, guaranteeing privacy even when offloading sensitive data

Synchronization

`git-annex assistant` provides Dropbox-like operation: keep your data files in sync, or automatically offloaded, backed-up, *etc.* across multiple hosts

Semi-compatible with GitHub

As a `git-annex` repository is just a `git` repository, it can be pushed to GitHub. The content of data files, however, need to be stored in some *Special remote* or have associated download URLs

git-annex 101

Just learn few new commands in addition to stock `git`. <https://git-annex.branchable.com/walkthrough> provides more details and examples

```
git init; git annex init "repo name"
```

Initialize both `git` and `git-annex`

```
git annex add my_cool_big_file
```

Add file(s) under control of `git-annex`

```
git commit -a -m "message"
```

Commit as you usually do with `git`

```
git annex move my_cool_big_file --to usbdrive
```

Offload the file to a repository you named “usbdrive”

```
git annex whereis my_cool_big_file
```

Check where file is available from

```
git annex get my_cool_big_file
```

Fetch file back to the local repository

Neuroimaging datasets to play with

We have made some datasets already available as `git-annex` repositories, so `git clone ...`

```
http://psydata.ovgu.de/forrest\_gump/.git/
```

Unique 7T MRI, fMRI, DTI dataset acquired during rich auditory stimulation. See

<http://studyforrest.org> for more information

```
git://github.com/datalad/nih--videocast
```

`git-annex` “mirror” of <http://videocast.nih.gov>, created and updated by `datalad crawl`

```
http://data.pympva.org/datasets/haxby2001/.git/
```

Seminal work by Haxby et al. (2001) for testing and demonstrating MVPA techniques

How to install git-annex

```
apt-get install git-annex
```

See <https://git-annex.branchable.com/install> for more (OS X, Windows, Linux, Android)

How to get support

On Debian systems

```
reportbug git-annex
```

Community support

<http://git-annex.branchable.com/bugs>

IRC

#git-annex at OFTC network



<http://datalad.org>

DataLad

aims to simplify and thus facilitate delivery and sharing of scientific data by establishing a federated data distribution. While initially aiming to deliver public neuroimaging datasets, DataLad will be easy to adopt more neuroscience data or other fields of endeavor.

DataLad FAQ

Federated?

It is impractical to distribute data through *classical* distribution mechanisms, where content is contained within packages available from the central location (or its mirrors). DataLad will only collect, unify, monitor, and expose through convenient interfaces data available across a wide range of data providers

Distributed?

DataLad uses distributed version control `Git` and built on top of it `Git-annex` for data logistics. `Git-annex` enables *distributed* operation where clones of the datasets could be made available across multiple sites and media without losing track of data and meta-information (such as versioning)

Planned dataset coverage

OpenfMRI.org

curated fMRI (and EEG) datasets

HumanConnectome.org

anatomical, functional, diffusion MRI data from 1,200 subjects

CRCNS.org

curated electrophysiological and neuroimaging datasets

INDI

a collation of various datasets and initiatives (functional connectome, *etc.*)

Planned integration

We will expose and interface to the datasets available from

XNAT

widely used imaging informatics platform used by [HumanConnectome.org](#), [NITRC-IR](#), [OpenfMRI](#) and others

COINS

web-based neuroimaging and neuropsychology software suite hosting many neuroimaging datasets

Academic Torrents

collection of academic datasets delivered as [Torrents](#)

NeuroDebian

Through the joint venture with the [NeuroDebian](#) project, DataLad will also expose itself as a [Debian](#) APT repository, making it possible to *install* and *upgrade* datasets using conventional tools such as `apt` and `aptitude`. Data *installation* would become as easy as software installation.

How could you help *both of us*?

Sharing scientific data is not yet as easy as it could and should be. Adhering to the following guidelines could help you to avoid unnecessary burden, thus making sharing easier and thus more rewarding. Overall motto is *Be Ready* and ...

Stay legit:

clear up and state ahead ownership (copyright) and (public domain dedication) license for your dataset. Provision public sharing in your consent forms **BEFORE** the data collection begins:

<https://open-brain-consent.readthedocs.org>

Keep detail:

keep original detail – copies of acquisition protocols, exam cards, and the DICOMs (not only NIFTIs)

Be comprehensible:

adhere to a homogeneous files structure, adopt and extend if necessary some standard (e.g., [openfMRI](#)). Consider providing dataset descriptor (e.g., [W3C HCLS Dataset descriptors](#) or <http://dataprotocols.org/data-packages>)

Prepare to be reproduced:

analyze already pre-processed anonymized data

Version your data:

even data close to bare origin might be *screwed* and require versioning. You could

- use [Git-annex](#) for your data
- consistently version older versions with the date in the suffix, e.g., how [1000genomes](#) project does
- turn on versioning for your [AWS S3](#) bucket(s)

Think about longevity:

Deposit datasets to some public community repository (e.g., [OpenfMRI](#), [figshare](#))

How could you help *DataLad*?

We have started the development (within both `git-annex` and `DataLad`) of necessary features and hope to deliver an initial functional prototype later this year. Meanwhile we would appreciate if you

Contribute

<http://github.com/datalad/datalad/pulls>

Follow&Share

Twitter: <http://twitter.com/datalad>

Google+: <http://plus.google.com/+DataladOrg>

Blog: <http://datalad.org>

Discuss

<https://groups.google.com/forum/#!forum/datalad>

Complain&Suggest

<http://github.com/datalad/datalad/issues>

We are interested in use-cases, interesting datasets, feedback on design decisions, alpha-users, etc.

Acknowledgements

This project is co-funded by the US National Science Foundation ([1429999](#)) and the German Federal Ministry of Education and Research (BMBF 01GQ1411)



SPONSORED BY THE

